### 9.5.3   Starting Hive Shell

To start Hive, go to the installation path of Hive and type as below:

```
[root@volgalnx005 ~]# hive

Logging initialized using configuration in jar:file:/root/Desktop/VMDATA/Hive/hi
ve/lib/hive-common-0.14.0.jar!/hive-log4j.properties
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/root/Desktop/VMDATA/Hadoop/hadoop/share/hadoo
p/common/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/root/Desktop/VMDATA/Hive/hive/lib/hive-jdbc-0
.14.0-standalone.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
hive>
```

The sections have been designed as follows:

*Objective:*          What is it that we are trying to achieve here?
*Input (optional):*   What is the input that has been given to us to act upon?
*Act:*                The actual statement/command to accomplish the task at hand.
*Outcome:*            The result/output as a consequence of executing the statement.

### 9.5.4   Database

A database is like a container for data. It has a collection of tables which houses the data.

**Objective:** To create a database named "STUDENTS" with comments and database properties.

**Act:**

> **CREATE DATABASE IF NOT EXISTS STUDENTS COMMENT 'STUDENT Details'**
> **WITH DBPROPERTIES ('creator' = 'JOHN');**

**Outcome:**

```
hive> CREATE DATABASE IF NOT EXISTS STUDENTS COMMENT 'STUDENT Details' WITH DBPROPERTIES ('creato
r' = 'JOHN');
OK
Time taken: 0.536 seconds
hive>
```

**Explanation of the syntax:**

**IF NOT EXIST:** It is an optional clause. The create database statement with "IF Not EXISTS" clause creates a database if it does not exist. However, if the database already exists then it will notify the user that a database with the same name already exists and will not show any error message.

**COMMENT:** This is to provide short description about the database.

**WITH DBPROPERTIES:** It is an optional clause. It is used to specify any properties of database in the form of (key, value) separated pairs. In the above example, "Creator" is the "Key" and "JOHN" is the value.

We can use "**SCHEMA**" in place of "**DATABASE**" in this command.

**Note:** We have not specified the location where the Hive database will be created. By default all the Hive databases will be created under default warehouse directory (set by the property hive.metastore.warehouse. dir) as /user/hive/warehouse/database_name.db. But if we want to specify our own location, then the LOCATION clause can be specified. This clause is optional.

**Objective:** To display a list of all databases.

**Act:**

```
SHOW DATABASES;
```

**Outcome:**

```
hive> SHOW DATABASES;
OK
students
Time taken: 0.082 seconds, Fetched: 22 row(s)
hive>
```

By default, **SHOW DATABASES** lists all the databases available in the metastore. We can use "**SCHEMAS**" in place of "**DATABASES**" in this command. The command has an optional "Like" clause. It can be used to filter the database names using regular expressions such as "*" , "?", etc.

SHOW DATABASES LIKE "Stu*"

SHOW DATABASES like "Stud???s"

**Objective:** To describe a database.

**Act:**

```
DESCRIBE DATABASE STUDENTS;
```

**Note:** Shows only DB name, comment, and DB directory.

**Outcome:**

```
hive> DESCRIBE DATABASE STUDENTS;
OK
students        STUDENT Details hdfs://volgalnx010.ad.infosys.com:9000/user/hive/warehouse/studen
ts.db   root    USER
Time taken: 0.03 seconds, Fetched: 1 row(s)
hive>
```

**Objective:** To describe the extended database.

**Act:**

```
DESCRIBE DATABASE EXTENDED STUDENTS;
```

**Note:** Shows DB properties also.

**Outcome:**

```
hive> DESCRIBE DATABASE EXTENDED STUDENTS;
OK
students        STUDENT Details hdfs://volgalnx010.ad.infosys.com:9000/user/hive/warehouse/studen
ts.db   root    USER    {creator=JOHN}
Time taken: 0.027 seconds, Fetched: 1 row(s)
hive> █
```

**DESCRIBE DATABASE EXTENDED** shows database's properties given under DBPROPERTIES argument at the time of creation.

We can use "**SCHEMA**" in place of "**DATABASE**", "**DESC**" in place of "**DESCRIBE**" in this command.

**Objective:** To alter the database properties.

**Act:**

**ALTER DATABASE STUDENTS SET DBPROPERTIES ('edited-by' = 'JAMES');**

**Note:** We can use the "ALTER DATABASE" command to
- Assign any new (key, value) pairs into DBPROPERTIES.
- Set owner user or role to the Database.

In Hive, it is not possible to unset the DB properties.

**Outcome:**

```
hive> ALTER DATABASE STUDENTS SET DBPROPERTIES ('edited-by' = 'JAMES');
OK
Time taken: 0.086 seconds
hive> █
```

In the above example, the ALTER DATABASE command is used to assign new ('edited-by' = 'JAMES') pair into DBPROPERTIES. This can be verified by using the 'describe *extended*'.

Hive>  DESCRIBE DATABASE Student EXTENDED

**Objective:** To make the database as current working database.

**Act:**

**USE STUDENTS;**

**Outcome:**

```
hive> USE STUDENTS;
OK
Time taken: 0.02 seconds
hive> █
```

There is no command to show the current database, but use the below command statement to keep printing the current database name as suffix in the command line prompt.

set hive.cli.print.current.db=true;

**Objective:** To drop database.

**Act:**

> **DROP DATABASE STUDENTS;**

**Note:** Hive creates database in the warehouse directory of Hive as shown below:

Contents of directory /user/hive/warehouse

Goto : /user/hive/warehouse    go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| students.db | dir | | | | 2015-02-24 21:50 | rwxr-xr-x | root | supergroup |

Now assume that the database "STUDENTS" has 10 tables within it. How do we delete the complete database along with the tables contained therein?

Use the command:

**DROP DATABASE STUDENTS CASCADE;**

By default the mode is RESTRICT which implies that the database will NOT be dropped if it contains tables.

**Note:** The complete syntax is as follows:

**DROP DATABASE [IF EXISTS] database_name [RESTRICT | CASCADE]**

## 9.5.5   Tables

Hive provides two kinds of table:

1. Internal or Managed Table
2. External Table

### 9.5.5.1   Managed Table

1. Hive stores the Managed tables under the warehouse folder under Hive.
2. The complete life cycle of table and data is managed by Hive.
3. When the internal table is dropped, it drops the data as well as the metadata.

When you create a table in Hive, by default it is internal or managed table. If one needs to create an external table, one will have to use the keyword "EXTERNAL".

**Objective:** To create managed table named 'STUDENT'.

**Act:**

> **CREATE TABLE IF NOT EXISTS STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';**

**Outcome:**

```
hive> CREATE TABLE IF NOT EXISTS STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMAT DELIMITED F
IELDS TERMINATED BY '\t';
OK
Time taken: 0.355 seconds
hive>
```

**Objective:** To describe the "STUDENT" table.

**Act:**

**DESCRIBE STUDENT;**

**Outcome:**

```
hive> DESCRIBE STUDENT;
OK
rollno                  int
name                    string
gpa                     float
Time taken: 0.163 seconds, Fetched: 3 row(s)
hive>
```

**Note:** Hive creates managed table in the warehouse directory of Hive as shown below:

Contents of directory /user/hive/warehouse/students.db

Goto : /user/hive/warehouse/student  go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| student | dir | | | | 2015-02-24 22:03 | rwxr-xr-x | root | supergroup |

Go back to DFS home

**Local logs**

Log directory

Hadoop, 2015.

**To check whether an existing table is managed or external, use the below syntax:**

DESCRIBE FORMATTED tablename;

It displays complete metadata of a table. You will see one row called table type which will display either MANAGED_TABLE OR EXTERNAL_TABLE.

DESCRIBE FORMATTED STUDENT;

### 9.5.5.2  External or Self-Managed Table

1. When the table is dropped, it retains the data in the underlying location.
2. **External** keyword is used to create an external table.
3. **Location** needs to be specified to store the dataset in that particular location.

**Objective:** To create external table named 'EXT_STUDENT'.

**Act:**

**CREATE EXTERNAL TABLE IF NOT EXISTS EXT_STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' LOCATION '/STUDENT_INFO;**

**Outcome:**

```
hive> CREATE EXTERNAL TABLE IF NOT EXISTS EXT_STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMA
T DELIMITED FIELDS TERMINATED BY '\t'  LOCATION '/STUDENT_INFO';
OK
Time taken: 0.123 seconds
hive>
```

**Note:** Hive creates the external table in the specified location.

### 9.5.5.3 Loading Data into Table from File

**Objective:** To load data into the table from file named student.tsv.

**Act:**

**LOAD DATA LOCAL INPATH '/root/hivedemos/student.tsv' OVERWRITE INTO TABLE EXT_STUDENT;**

**Note:** Local keyword is used to load the data from the local file system. In this case, the file is copied. To load the data from HDFS, remove local key word from the statement. In this case, the file is moved from the original location.

**Outcome:**

```
hive> LOAD DATA LOCAL INPATH  '/root/hivedemos/student.tsv' OVERWRITE INTO TABLE EXT_STUDENT;
Loading data to table students.ext_student
Table students.ext_student stats: [numFiles=0, numRows=0, totalSize=0, rawDataSize=0]
OK
Time taken: 5.034 seconds
hive>
```

Hive loads the file in the specified location as shown below:

**Contents of directory /STUDENT_INFO**

Goto :/STUDENT_INFO      go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|---|---|---|---|---|---|---|---|---|
| student.tsv | file | 121 B | 3 | 128 MB | 2015-02-24 22:19 | rw-r--r-- | root | supergroup |

Go back to DFS home

**Local logs**

Log directory

Hadoop, 2015.

File: /STUDENT_INFO/student.tsv

Goto: /STUDENT_INFO          go

*Go back to dir listing*
*Advanced view/download options*

```
1001   John   3.0
1002   Jack   4.0
1003   Smith  4.5
1004   Jamesh 4.2
1005   Joshi  3.5
1006   Alex   4.0
1007   David  4.2
1008   Scott  3.9
```

Let us understand the difference between INTO TABLE and OVERWRITE TABLE with an example:

Assume the "EXT_STUDENT" table already had 100 records and the "student.tsv" file has 10 records. After issuing the LOAD DATA statement with the INTO TABLE clause, the table "EXT_STUDENT" will contain 110 records; however, the same LOAD DATA statement with the OVERWRITE clause will wipe out all the former content from the table and then load the 10 records from the data file.

### 9.5.5.4 *Collection Data Types*

**Objective:** To work with collection data types.

**Input:**
1001,John,Smith:Jones,Mark1!45:Mark2!46:Mark3!43
1002,Jack,Smith:Jones,Mark1!46:Mark2!47:Mark3!42

**Act:**

> **CREATE TABLE STUDENT_INFO (rollno INT,name String, sub ARRAY<STRING>,marks MAP<STRING,INT>)**
>
> **ROW FORMAT DELIMITED FIELDS TERMINATED BY ',''**
>
> **COLLECTION ITEMS TERMINATED BY ':'**
>
> **MAP KEYS TERMINATED BY '!';**
>
> **LOAD DATA LOCAL INPATH '/root/hivedemos/studentinfo.csv' INTO TABLE STUDENT_INFO;**

**Outcome:**

```
hive> CREATE TABLE STUDENT_INFO (rollno INT,name String, sub ARRAY<STRING>,marks MAP<STRING,FLOAT>)
    > ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
    > COLLECTION ITEMS TERMINATED BY ':'
    > MAP KEYS TERMINATED BY '!';
OK
Time taken: 0.112 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/studentinfo.csv' INTO TABLE STUDENT_INFO;
Loading data to table students.student_info
Table students.student_info stats: [numFiles=1, totalSize=109]
OK
Time taken: 0.397 seconds
hive>
```

### 9.5.5.5   Querying Table

**Objective:** To retrieve the student details from "EXT_STUDENT" table.

**Act:**

```
SELECT * from EXT_STUDENT;
```

**Outcome:**

```
hive> select * from EXT_STUDENT;
OK
1001    John    3.0
1002    Jack    4.0
1003    Smith   4.5
1004    Scott   4.2
1005    Joshi   3.5
1006    Alex    4.5
1007    David   4.2
1008    James   4.0
1009    John    3.0
1010    Joshi   3.5
Time taken: 0.054 seconds, Fetched: 10 row(s)
hive>
```

**Objective:** Querying Collection Data Types.

**Act:**

```
SELECT * from STUDENT_INFO;
SELECT NAME,SUB FROM STUDENT_INFO;
// To retrieve value of Mark1
SELECT NAME, MARKS['Mark1'] from STUDENT_INFO;
// To retrieve subordinate (array) value
SELECT NAME,SUB[0] FROM STUDENT_INFO;
```

**Outcome:**

```
hive> SELECT * from STUDENT_INFO;
OK
1001    John    ["Smith","Jones"]       {"Mark1":45,"Mark2":46,"Mark3":43}
1002    Jack    ["Smith","Jones"]       {"Mark1":46,"Mark2":47,"Mark3":42}
Time taken: 0.044 seconds, Fetched: 2 row(s)

hive> SELECT NAME,SUB FROM STUDENT_INFO;
OK
John    ["Smith","Jones"]
Jack    ["Smith","Jones"]
Time taken: 0.061 seconds, Fetched: 2 row(s)
hive>

hive> SELECT NAME, MARKS['Mark1']  from STUDENT_INFO;
OK
John    45
Jack    46
Time taken: 0.06 seconds, Fetched: 2 row(s)
hive>

hive> SELECT NAME,SUB[0] FROM STUDENT_INFO;
OK
John    Smith
Jack    Smith
Time taken: 0.071 seconds, Fetched: 2 row(s)
hive>
```

## 9.5.6 Partitions

In Hive, the query reads the entire dataset even though a where clause filter is specified on a particular column. This becomes a bottleneck in most of the MapReduce jobs as it involves huge degree of I/O. So it is necessary to reduce I/O required by the MapReduce job to improve the performance of the query. A very common method to reduce I/O is data partitioning.

Partitions split the larger dataset into more meaningful chunks.

We will try to understand partitioning with the help of a simple example.

**PICTURE THIS...**

"XYZ enterprise" has a wide customer base spread across several states in the US. The data has been fed to the central system. The senior leadership team would like to get a report providing the statewise Sales %.

The IT team has proposed two options to help service this request:

1. Run the query with the where clause of each state name (such as where StateName = "A"). This will mean that the entire dataset is scanned/read through for each and every state. If we have data for 25 states, the dataset is read 25 times (once for each state). Assuming we have 5 million records, it would mean that 5 million records are read 25 times. This can be a major performance bottle neck and will worsen as the dataset grows.

2. The second option stated here can help alleviate this problem. The IT team can start off with creating 25 folders (one for each state) and instruct to have the data pertaining to states place into the folder of the respective states. This arrangement will greatly help at the time of querying the data. To get the Sales % of a particular state, the folder belonging to that state only has to be scanned/read through.

This intelligent way of grouping data during data load is termed as PARTITIONING in Hive.

This brings us to the next question.

If you are looking through the folder for state = "A", is there any possibility of you coming across data for state = "B" or state = "C"? Think through! I am sure your answer will be No! And we know the reason.

A point to note here is that as we create the partitions, there is no need to include the partitioned column along with the other columns of the dataset as this is something that is automatically taken care of "BY PARTITIONED" clause.

In our example above, we will refrain from adding the partitioned column along with other columns of the dataset and trust Hive to automatically manage this.

Partition is of two types:

1. **STATIC PARTITION:** It is upon the user to mention the partition (the segregation unit) where the data from the file is to be loaded.

2. **DYNAMIC PARTITION:** The user is required to simply state the column, basis which the partitioning will take place. Hive will then create partitions basis the unique values in the column on which partition is to be carried out.

Points to consider as you create partitions:

1. STATIC PARTITIONING implies that the user controls everything from defining the PARTITION column to loading data into the various partitioned folders.

2. As in our example above, if STATIC partition is done over the STATE column and assume by mistake the data for state "B" is placed inside the partition for state "A", our query for data for state "B" is

bound to return zilch records. The reason is obvious. A Select fired on STATIC partition just takes into consideration the partition name, and does not consider the data held inside the partition.

3. DYNAMIC PARTITIONING means Hive will intelligently get the distinct values for partitioned column and segregate data into respective partitions. There is no manual intervention.

By default, dynamic partitioning is enabled in Hive. Also by default it is strictly implying that one is required to do one level of STATIC partitioning before Hive can perform DYNAMIC partitioning inside this STATIC segregation unit.

In order to go with full dynamic partitioning, we have to set below property to non-strict in Hive.

hive> set hive.exec.dynamic.partition.mode=nonstrict

### 9.5.6.1 Static Partition

Static partitions comprise columns whose values are known at compile time.

**Objective:** To create static partition based on "gpa" column.

**Act:**

**CREATE TABLE IF NOT EXISTS STATIC_PART_STUDENT (rollno INT, name STRING) PARTITIONED BY (gpa FLOAT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';**

**Outcome:**

```
hive> CREATE TABLE IF NOT EXISTS STATIC_PART_STUDENT(rollno INT,name STRING) PARTITIONED BY (gpa FLOAT) RO
W FORMAT DELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.105 seconds
hive>
```

**Objective:** Load data into partition table from table.

**Act:**

**INSERT OVERWRITE TABLE STATIC_PART_STUDENT PARTITION (gpa =4.0) SELECT rollno, name from EXT_STUDENT where gpa=4.0;**

**Outcome:**

```
hive> INSERT OVERWRITE TABLE STATIC_PART_STUDENT PARTITION (gpa =4.0) SELECT rollno,name from EXT_STUDENT
where gpa=4.0;
Query ID = root_20150224230404_4500d58a-cb21-4912-ba40-788e5cf8f9da
Total jobs = 3
```

Hive creates the folder for the value specified in the partition.

Contents of directory /user/hive/warehouse/students.db

Goto [/user/hive/warehouse/student] go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| static_part_student | dir | | | | 2015-02-24 23:04 | rwxr-xr-x | root | supergroup |
| student | dir | | | | 2015-02-24 22:03 | rwxr-xr-x | root | supergroup |
| student_info | dir | | | | 2015-02-24 22:54 | rwxr-xr-x | root | supergroup |

Go back to DFS home

**Local logs**

Log directory

Contents of directory /user/hive/warehouse/students.db/static_part_student

Goto : /user/hive/warehouse/student  go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| gpa=4.0 | dir | | | | 2015-02-24 23:04 | rwxr-xr-x | root | supergroup |

Go back to DFS home

**Local logs**

Log directory

File: /user/hive/warehouse/students.db/static_part_student/gpa=4.0/000000_0

Goto : /user/hive/warehouse/student  go

Go back to dir listing
Advanced view/download options

```
1002    Jack
1008    James
```

**Objective:** To add one more static partition based on "gpa" column using the "alter" statement.

**Act:**

> ALTER TABLE STATIC_PART_STUDENT ADD PARTITION (gpa=3.5);
>
> INSERT OVERWRITE TABLE STATIC_PART_STUDENT PARTITION (gpa =4.0) SELECT rollno,name from EXT_STUDENT where gpa=4.0;

**Outcome:**

```
hive> ALTER TABLE STATIC_PART_STUDENT ADD PARTITION (gpa=3.5);
OK
Time taken: 0.166 seconds
hive>
```

Contents of directory /user/hive/warehouse/students.db/static_part_student

Goto : /user/hive/warehouse/student  go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| gpa=3.5 | dir | | | | 2015-02-24 23:09 | rwxr-xr-x | root | supergroup |
| gpa=4.0 | dir | | | | 2015-02-24 23:11 | rwxr-xr-x | root | supergroup |

Go back to DFS home

### 9.5.6.2 Dynamic Partition

Dynamic partition have columns whose values are known only at Execution Time.

**Objective:** To create dynamic partition on column date.

**Act:**

> CREATE TABLE IF NOT EXISTS DYNAMIC_PART_STUDENT(rollno INT,name STRING) PARTITIONED BY (gpa FLOAT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';

**Outcome:**

```
hive> CREATE TABLE IF NOT EXISTS DYNAMIC_PART_STUDENT(rollno INT,name STRING) PARTITIONED BY (gpa FLOAT) R
OW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.166 seconds
hive>
```

**Objective:** To load data into a dynamic partition table from table.

**Act:**

> **SET hive.exec.dynamic.partition = true;**
>
> **SET hive.exec.dynamic.partition.mode = nonstrict;**
>
> **Note:** The dynamic partition strict mode requires at least one static partition column. To turn this off, set hive.exec.dynamic.partition.mode=nonstrict
>
> **INSERT OVERWRITE TABLE DYNAMIC_PART_STUDENT PARTITION (gpa) SELECT rollno,name,gpa from EXT_STUDENT;**

**Outcome:**

Contents of directory /user/hive/warehouse/students.db/dynamic_part_student

Goto : /user/hive/warehouse/student  go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|------|------|------|-------------|------------|-------------------|------------|-------|-------|
| gpa=3.0 | dir | | | | 2015-02-24 23:16 | rwxr-xr-x | root | supergroup |
| gpa=3.5 | dir | | | | 2015-02-24 23:16 | rwxr-xr-x | root | supergroup |
| gpa=4.0 | dir | | | | 2015-02-24 23:16 | rwxr-xr-x | root | supergroup |
| gpa=4.2 | dir | | | | 2015-02-24 23:16 | rwxr-xr-x | root | supergroup |
| gpa=4.5 | dir | | | | 2015-02-24 23:16 | rwxr-xr-x | root | supergroup |

Go back to DFS home

**Note:** Create partition for all values.

## 9.5.7 Bucketing

Bucketing is similar to partition. However, there is a subtle difference between partition and bucketing. In a partition, you need to create partition for each unique value of the column. This may lead to situations where you may end up with thousands of partitions. This can be avoided by using Bucketing in which you can limit the number of buckets that will be created. A bucket is a file whereas a partition is a directory.

**Objective:** To learn the concept of bucket in hive.

**Act:**

> **CREATE TABLE IF NOT EXISTS STUDENT (rollno INT,name STRING,grade FLOAT)**
> **ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';**
> **LOAD DATA LOCAL INPATH '/root/hivedemos/student.tsv' INTO TABLE STUDENT;**
>
> Set below property to enable bucketing.
>
> **set hive.enforce.bucketing=true;**

// To create a bucketed table having 3 buckets

**CREATE TABLE IF NOT EXISTS STUDENT_BUCKET (rollno INT,name STRING,grade FLOAT)**

**CLUSTERED BY (grade) into 3 buckets;**

// Load data to bucketed table

**FROM STUDENT**

**INSERT OVERWRITE TABLE STUDENT_BUCKET**

**SELECT rollno,name,grade;**

// To display content of first bucket

**SELECT DISTINCT GRADE FROM STUDENT_BUCKET**

**TABLESAMPLE(BUCKET 1 OUT OF 3 ON GRADE);**

## Outcome:

```
hive> CREATE TABLE IF NOT EXISTS STUDENT (rollno INT,name STRING,grade FLOAT)
    > ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.101 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/student.tsv' INTO TABLE STUDENT;
Loading data to table book.student
Table book.student stats: [numFiles=1, totalSize=145]
OK
Time taken: 0.536 seconds
hive>
```

```
hive> set hive.enforce.bucketing=true;
hive>
```

```
hive> CREATE TABLE IF NOT EXISTS STUDENT_BUCKET (rollno INT,name STRING,grade FLOAT)
    > CLUSTERED BY (grade) into 3 buckets;
OK
Time taken: 0.101 seconds
hive>
```

```
hive> FROM STUDENT
    > INSERT OVERWRITE TABLE STUDENT_BUCKET
    > SELECT rollno,name,grade;
```

## 3 buckets have been created as shown below:

Contents of directory /user/hive/warehouse/book.db/student_bucket

Goto : /user/hive/warehouse/book.d| go

Go to parent directory

| Name | Type | Size | Replication | Block Size | Modification Time | Permission | Owner | Group |
|---|---|---|---|---|---|---|---|---|
| 000000_0 | file | 59 B | 3 | 128 MB | 2015-03-10 22:29 | rw-r--r-- | root | supergroup |
| 000001_0 | file | 59 B | 3 | 128 MB | 2015-03-10 22:29 | rw-r--r-- | root | supergroup |
| 000002_0 | file | 28 B | 3 | 128 MB | 2015-03-10 22:29 | rw-r--r-- | root | supergroup |

Go back to DFS home

## Local logs

Log directory

Hadoop, 2015.

```
hive>
    > SELECT DISTINCT GRADE FROM STUDENT_BUCKET
    > TABLESAMPLE(BUCKET 1 OUT OF 3 ON GRADE);
OK
4.0
4.2
Time taken: 21.117 seconds, Fetched: 2 row(s)
hive>
```

## 9.5.8 Views

In Hive, view support is available only in version starting from 0.6. Views are purely logical object.

**Objective:** To create a view table named "STUDENT_VIEW".

**Act:**

**CREATE VIEW STUDENT_VIEW AS SELECT rollno, name FROM EXT_STUDENT;**

**Outcome:**

```
hive> CREATE VIEW STUDENT_VIEW AS SELECT rollno,name FROM EXT_STUDENT;
OK
Time taken: 0.606 seconds
hive>
```

**Objective:** Querying the view "STUDENT_VIEW".

**Act:**

**SELECT * FROM STUDENT_VIEW LIMIT 4;**

**Outcome:**

```
hive> SELECT * FROM STUDENT_VIEW LIMIT 4;
OK
1001    John
1002    Jack
1003    Smith
1004    Scott
Time taken: 0.279 seconds, Fetched: 4 row(s)
hive>
```

**Objective:** To drop the view "STUDENT_VIEW".

**Act:**

**DROP VIEW STUDENT_VIEW;**

**Outcome:**

```
hive> DROP VIEW STUDENT_VIEW;
OK
Time taken: 0.452 seconds
hive>
```

### 9.5.9   Sub-Query

In Hive, sub-queries are supported only in the FROM clause (Hive 0.12). You need to specify name for sub-query because every table in a FROM clause has a name. The columns in the sub-query select list should have unique names. The columns in the subquery select list are available to the outer query just like columns of a table.

**Objective:** Write a sub-query to count occurrence of similar words in the file.

**Act:**

**CREATE TABLE docs (line STRING);**

**LOAD DATA LOCAL INPATII '/root/hivedemos/lines.txt' OVERWRITE INTO TABLE docs;**

**CREATE TABLE word_count AS**

**SELECT word, count(1) AS count FROM**

**(SELECT explode (split (line, ' ')) AS word FROM docs) w**

**GROUP BY word**

**ORDER BY word;**

**SELECT * FROM word_count;**

**Outcome:**

```
hive> CREATE TABLE docs (line STRING);
OK
Time taken: 0.118 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/lines.txt' OVERWRITE INTO TABLE docs;
Loading data to table students.docs
Table students.docs stats: [numFiles=1, numRows=0, totalSize=91, rawDataSize=0]
OK
Time taken: 2.697 seconds
hive>
```

```
hive> CREATE TABLE word_count AS
    > SELECT word, count(1) AS count FROM
    > (SELECT explode (split (line, ' ')) AS word FROM docs) w
    > GROUP BY word
    > ORDER BY word:

hive> SELECT * FROM word_count;
OK
Hadoop  2
Hive    2
Introducing     1
Introduction    1
Pig     1
Session 3
Welcome 1
to      2
Time taken: 0.062 seconds, Fetched: 8 row(s)
hive>
```

**Note:** The explode() function takes an array as input and outputs the elements of the array as separate rows.

**In Hive 0.13, sub-queries are supported in the where clause as well.**

## 9.5.10 Joins

Joins in Hive is similar to the SQL Join.

---

**Objective:** To create JOIN between Student and Department tables where we use RollNo from both the tables as the join key.

**Act:**

**CREATE TABLE IF NOT EXISTS STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';**

**LOAD DATA LOCAL INPATH '/root/hivedemos/student.tsv' OVERWRITE INTO TABLE STUDENT;**

**CREATE TABLE IF NOT EXISTS DEPARTMENT(rollno INT,deptno int,name STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';**

**LOAD DATA LOCAL INPATH '/root/hivedemos/department.tsv' OVERWRITE INTO TABLE DEPARTMENT;**

**SELECT a.rollno, a.name, a.gpa, b.deptno FROM STUDENT a JOIN DEPARTMENT b ON a.rollno = b.rollno;**

**Outcome:**

```
hive> CREATE TABLE IF NOT EXISTS STUDENT(rollno INT,name STRING,gpa FLOAT) ROW FORMAT D
ELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.115 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/student.tsv' OVERWRITE INTO TABLE STUDENT
;
Loading data to table students.student
Table students.student stats: [numFiles=1, numRows=0, totalSize=145, rawDataSize=0]
OK
Time taken: 0.723 seconds
hive>
```

```
hive> CREATE TABLE IF NOT EXISTS DEPARTMENT(rollno INT,deptno int,name STRING) ROW FORM
AT DELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.099 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/department.tsv' OVERWRITE INTO TABLE DEPA
RTMENT;
Loading data to table students.department
Table students.department stats: [numFiles=1, numRows=0, totalSize=120, rawDataSize=0]
OK
Time taken: 0.442 seconds
hive>
```

```
hive> SELECT a.rollno, a.name, a.gpa, b.deptno  FROM STUDENT  a JOIN DEPARTMENT b ON a.
rollno = b.rollno;
OK
1001    John    3.0     101
1002    Jack    4.0     102
1003    Smith   4.5     103
1004    Scott   4.2     104
1005    Joshi   3.5     105
1006    Alex    4.5     101
1007    David   4.2     104
1008    James   4.0     102
Time taken: 115.282 seconds, Fetched: 8 row(s)
hive>
```

### 9.5.11 Aggregation

Hive supports aggregation functions like avg, count, etc.

**Objective:** To write the average and count aggregation functions.

**Act:**

```
SELECT avg(gpa) FROM STUDENT;
SELECT count(*) FROM STUDENT;
```

**Outcome:**

```
hive> SELECT avg(gpa) FROM STUDENT;

OK
3.839999961853027
Time taken: 28.639 seconds, Fetched: 1 row(s)
hive>
```

```
hive> SELECT count(*) FROM STUDENT;

OK
10
Time taken: 26.218 seconds, Fetched: 1 row(s)
hive>
```

### 9.5.12 Group By and Having

Data in a column or columns can be grouped on the basis of values contained therein by using "Group By". "Having" clause is used to filter out groups NOT meeting the specified condition.

**Objective:** To write group by and having function.

**Act:**

```
SELECT rollno, name,gpa FROM STUDENT GROUP BY rollno,name,gpa HAVING gpa >
4.0;
```

**Outcome:**

```
1003    Smith    4.5
1004    Scott    4.2
1006    Alex     4.5
1007    David    4.2
Time taken: 78.972 seconds, Fetched: 4 row(s)
hive>
```

## 9.6 RCFILE IMPLEMENTATION

**RCFile** (Record Columnar File) is a data placement structure that determines how to store relational tables on computer clusters.

**Objective:** To work with RCFILE Format.

**Act:**

> CREATE TABLE STUDENT_RC( rollno int, name string,gpa float ) STORED AS RCFILE;
> INSERT OVERWRITE table STUDENT_RC SELECT * FROM STUDENT;
> SELECT SUM(gpa) FROM STUDENT_RC;

**Outcome:**

```
hive> CREATE TABLE STUDENT_RC( rollno int, name string,gpa float ) STORED AS RCFILE;
OK
Time taken: 0.093 seconds
hive>
```

```
hive> INSERT OVERWRITE table STUDENT_RC SELECT * from STUDENT;
```

```
hive> SELECT SUM(gpa) from STUDENT_RC;

OK
38.39999961853027
Time taken: 25.41 seconds, Fetched: 1 row(s)
hive>
```

**Note:** Stores the data in column oriented manner.

File: /user/hive/warehouse/students.db/student_rc/000000_0

Goto : /user/hive/warehouse/student  go

Go back to dir listing
Advanced view download options

RCFHive.is.rcfile.column.number3V◆◆◆◆◆◆7◆◆7◆◆◆
◆..    ◆◆◆((◆◆◆◆◆◆◆◆◆◆◆◆◆JohnJackSmithScottJoshiAlexDavidJamesJohnJoshi◆◆◆◆◆◆◆ff◆˙◆◆◆ff◆◆◆◆˙

## 9.7  SERDE

SerDe stands for Serializer/Deserializer.

1. Contains the logic to convert unstructured data into records.
2. Implemented using Java.
3. Serializers are used at the time of writing.
4. Deserializers are used at query time (SELECT Statement).

Deserializer interface takes a binary representation or string of a record, converts it into a java object that Hive can then manipulate. Serializer takes a java object that Hive has been working with and translates it into something that Hive can write to HDFS.

**Objective:** To manipulate the XML data.

**Input:**

<employee>  <empid>1001</empid>  <name>John</name>  <designation>Team  Lead</designation>
</employee>
<employee>  <empid>1002</empid>  <name>Smith</name>  <designation>Analyst</designation>
</employee>

**Act:**

CREATE TABLE XMLSAMPLE(xmldata string);

LOAD DATA LOCAL INPATH '/root/hivedemos/input.xml' INTO TABLE XMLSAMPLE;

CREATE TABLE xpath_table AS

SELECT xpath_int(xmldata,'employee/empid'),

xpath_string(xmldata,'employee/name'),

xpath_string(xmldata,'employee/designation')

FROM xmlsample;

SELECT * FROM xpath_table;

**Outcome:**

```
hive> CREATE TABLE XMLSAMPLE(xmldata string);
OK
Time taken: 0.244 seconds
hive>
```

```
hive> LOAD DATA LOCAL INPATH '/root/hivedemos/input.xml' INTO TABLE XMLSAMPLE;
Loading data to table students.xmlsample
Table students.xmlsample stats: [numFiles=1, totalSize=194]
OK
Time taken: 0.889 seconds
hive>
```

```
hive> CREATE TABLE xpath_table AS
    > SELECT xpath_int(xmldata,'employee/empid'),
    > xpath_string(xmldata,'employee/name'),
    > xpath_string(xmldata,'employee/designation')
    > FROM xmlsample;

hive> SELECT * FROM xpath_table;
OK
1001    John    Team Lead
1002    Smith   Analyst
Time taken: 0.064 seconds, Fetched: 2 row(s)
hive>
```

## 9.8   USER-DEFINED FUNCTION (UDF)

In Hive, you can use custom functions by defining the User-Defined Function (UDF).

**Objective:** Write a Hive function to convert the values of a field to uppercase.

**Act:**

```
package com.example.hive.udf;
import org.apache.hadoop.hive.ql.exec.Description;
import org.apache.hadoop.hive.ql.exec.UDF;
@Description(
  name="SimpleUDFExample")
```

```java
public final class MyLowerCase extends UDF {
  public String evaluate(final String word) {
    return word.toLowerCase();
  }
}
```

**Note:** Convert this Java Program into Jar.

**ADD JAR /root/hivedemos/UpperCase.jar;**

**CREATE TEMPORARY FUNCTION touppercase AS 'com.example.hive.udf.MyUpperCase';**

**SELECT TOUPPERCASE(name) FROM STUDENT;**

**Outcome:**

```
hive> ADD JAR /root/hivedemos/UpperCase.jar;
Added [/root/hivedemos/UpperCase.jar] to class path
Added resources: [/root/hivedemos/UpperCase.jar]
hive> CREATE TEMPORARY FUNCTION touppercase AS 'com.example.hive.udf.MyUpperCase';
OK
Time taken: 0.014 seconds
hive>
```

```
hive> Select touppercase (name) from STUDENT;
OK
JOHN
JACK
SMITH
SCOTT
JOSHI
ALEX
DAVID
JAMES
JOHN
JOSHI
Time taken: 0.061 seconds, Fetched: 10 row(s)
hive>
```

## REMIND ME

- Hive is a Data Warehousing tool.
- Hive is used to query structured data built on top of Hadoop.
- Hive provides HQL (Hive Query Language) which is similar to SQL.
- A Hive database contains several tables. Each table is constituted of rows and columns. In Hive, tables are stored as a folder and partition tables are stored as a sub-directory.
- Bucketed tables are stored as a file.

## POINT ME (BOOKS)

- Programming Hive, Jason Rutherglen, O'Reilly Publication.