PRACTICAL 7C

Partitioning and Bucketing

Partitioning the Table

Apache Hive is an open source data warehouse system used for querying and analyzing large datasets. Data in Apache Hive can be categorized into Table, Partition, and Bucket. The table in Hive is logically made up of the data being stored.

Hive provides way to categories data into smaller directories and files using partitioning or/and bucketing/clustering in order to improve performance of data retrieval queries and make them faster.

Main difference between Partitioning and Bucketing is that partitioning is applied directly on the column value and data is stored within directory named with column value whereas bucketing is applied using hash function on the column value MOD function with the number of buckets to store data in specific bucket file.

Hive table partition is a way to split a large table into smaller logical tables based on one or more partition keys. These smaller logical tables are not visible to users and users still access the data from just one table.

Partition eliminates creating smaller tables, accessing, and managing them separately.

To create a Hive table with partitions, you need to use PARTITIONED BY clause along with the column you wanted to partition and its type. Let's create a table and Load the CSV file.

The data file that I am using to explain partitions can be downloaded from GitHub, It's a simplified zipcodes codes where I have RecordNumber, Country, City, Zipcode, and State columns. I will be using State as a partition column.

Load Data into Partition Table

Download the <u>zipcodes.CSV from GitHub</u>, upload it to HDFS, and finally load the CSV file into a partition table.

Load Data into Partition Table

Download the <u>zipcodes.CSV from GitHub</u>, upload it to HDFS, and finally load the CSV file into a partition table.

Show All Partitions on Hive Table

After loading the data into the Hive partition table, you can use SHOW PARTITIONS command to see all partitions that are present.

```
hive> load data local inpath '/home/cloudera/Documents/zipcode.csv' into table zipcodes;
Loading data to table default.zipcodes
Table default.zipcodes stats: [numFiles=1, totalSize=591]
0K
Time taken: 0.538 seconds
hive> select * from zipcodes;
0K
NULL
        Country City
                         NULL
                                 State
                PARC PARQUE
        US
                                 704
                                          PR
                                          704
        US
                PASEO COSTA DEL SUR
10
        US
                BDA SAN LUIS
                                 709
                                          PR
        US
                                          76166
                                                  ΤX
61391
                CINGULAR WIRELESS
        US
61392
                FORT WORTH
                                 76177
                                          TX
61393
        US
                FT WORTH
                                  76177
                                          ΤX
                URB EUGENE RICE 704
                                          PR
        US
39827
        US
                MESA
                         85209
                                 ΑZ
39828
        US
                MESA
                         85210
                                 ΑZ
49345
        US
                HILLIARD
                                 32046
                                          FL
49346
                HOLDER 34445
        US
                                 FL
49347
        US
                HOLT
                         32564
                                 FL
                                 34487
49348
        US
                HOMOSASSA
                                          FL
        US
                SECT LANAUSSE
                                 704
                                          PR
54354
        US
                SPRING GARDEN
                                 36275
                                          AL
54355
        US
                SPRINGVILLE
                                 35146
                                          AL
54356
        US
                SPRUCE PINE
                                 35585
                                          AL
76511
        US
                ASH HILL
                                 27007
                                          NC
76512
        US
                ASHEB0R0
                                 27203
                                          NC
76513
        US
                ASHEB0R0
                                 27204
                                          NC
NULL
        NULL
                NULL
                        NULL
                                 NULL
Time taken: 0.345 seconds, Fetched: 22 row(s)
```

```
hive> create table zipcode(RecordNumber int,Country string,City string,Zipcode int) PARTITIONED BY(State string);
OK
Time taken: 0.053 seconds
hive> set hive.exec.dynamic.partition.mode=nonstrict;
```

Add New Partition to the Hive Table

A new partition can be added to the table using the ALERT TABLE statement, you can also specify the location where you wanted to store partition data on HDFS.

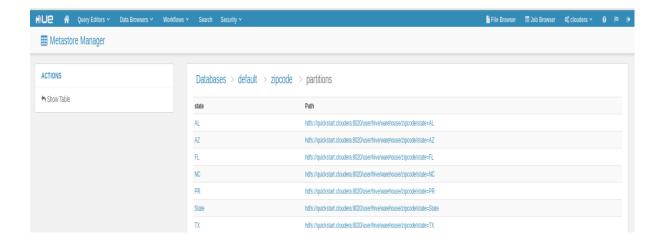
```
hive> insert overwrite table zipcode PARTITION(State) SELECT RecordNumber,Country,City,Zipcode,State from zipcodes;

Query ID = cloudera_0220322184444_4c8a90la-bbde-4aal-8c04-26e6bc3e38aa
Total_jobs = 3
Launching_Job 1 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
Starting_Job = job_1647952873179_0801, Tracking_URL = http://quickstart.cloudera:8088/proxy/application_1647952873179_0801/
KRIL Command = Jusr/Lib/hadoopybin/hadoopy job - kill job_1647952873179_0801
Hadoop_Job information for Stage=1 number of mappers: 1; number of reducers: 0
222:03-22 la3-44:34,826 Stage=1 num p= 100%, reduce = 0%, Cumulative CPU 1.0 sec
MaphReduce Total cumulative CPU Lime: 1 seconds 0 mec
Ended Dob = job_1647952873179_0801
Stage=3 is filtered on by condition resolver.
Stage=4 is selected by condition resolver.
Stage=3 is filtered on by condition resolver.
Stage=3 is filtered on by condition resolver.
Stage=4 is not able default.zipcode partition (state=mll)

Tate taken for load dynamic partitions (state=mll)

Loading partition (state=FR)
Loading partition (state=FR)
Loading partition (state=FR)
Loading partition (state=FR)
Loading partition (state=RR)
Loading partition (state=R
```

From the below image we can see that 6 partition have been created based on the name of the States.



Bucketing the Table

Hive Bucketing is a way to split the table into a managed number of clusters with or without partitions. With partitions, Hive divides(creates a directory) the table into smaller parts for every distinct value of a column whereas with bucketing you can specify the number of buckets to create at the time of creating a Hive table.

Load Data into Bucket

Loading/inserting data into the Bucketing table would be the same as inserting data into the table.

```
MapReduce Total cumulative CPU time: 35 seconds 950 msec
Ended Job = job_1646966376578_0003

Loading data to table default.zipcodes_bucket partition (state=null)

Time taken for load dynamic partitions: 3203

Loading partition {state= HIVE_DEFAULT_PARTITION_}

Loading partition {state=FR}

Loading partition {state=FR}

Loading partition {state=AZ}

Loading partition {state=AZ}

Loading partition {state=AZ}

Loading partition {state=AZ}

Loading partition {state=NEV}

Loading partition {state=NEV}

Loading partition {state=NEV}

Loading partition {state=AL}

Time taken for adding to write entity: 1

Partition default.zipcodes bucket{state=AL} stats: [numFiles=32, numRows=3, totalSize=83, rawDataSize=80]

Partition default.zipcodes bucket{state=AL} stats: [numFiles=32, numRows=2, totalSize=40, rawDataSize=88]

Partition default.zipcodes_bucket{state=FL} stats: [numFiles=32, numRows=3, totalSize=91, rawDataSize=87]

Partition default.zipcodes_bucket{state=FL} stats: [numFiles=32, numRows=3, totalSize=72, rawDataSize=69]

Partition default.zipcodes_bucket{state=PL} stats: [numFiles=32, numRows=3, totalSize=121, rawDataSize=69]

Partition default.zipcodes_bucket{state=State} stats: [numFiles=32, numRows=3, totalSize=19, rawDataSize=16]

Partition default.zipcodes_bucket{state=State} stats: [numFiles=32, numRows=3, totalSize=19, rawDataSize=18]

Partition default.zipcodes_bucket{state=MIVE_DEFAULT_PARTITION__) stats: [numFiles=32, numRows=3, rawDataSize=24, rawDataSize=28]

Partition default.zipcodes_bucket{state=_HIVE_DEFAULT_PARTITION__) stats: [numFiles=32, numRows=3, rawDataSize=24, rawDataSize=22]

MapReduce Jobs Launched:

Stage-Stage-1: Map: 1 Reduce: 32 Cumulative CPU: 35.95 sec HDFS Read: 119079 HDFS Write: 2102 SUCCESS

Total MapReduce CPU Time Spent: 35 seconds 950 msec

OK

Time taken: 204.824 seconds

hive>
```

Altering the table: Renaming the State name AL to 'NY'

```
hive> alter table zipcode partition(State='AL') rename to partition(State='NY');
DK
Time taken: 0.325 seconds
hive>
```

Now we can see from the below image ,the state name 'AL' is renamed to 'NY'.

