Q:Will the reducer work or not if you use "Limit 1" in any HiveQL query?

If we use a simple select * query then the reducer will not work as it will create a fetch task only. But if we use other queries such as group by or other aggregations then the reducer will work.

Q:Suppose I have installed Apache Hive on top of my Hadoop cluster using default metastore configuration. Then, what will happen if we have multiple clients trying to access Hive at the same time?

As the default metastore configuration in hive uses derby db only one connection is supported at a time. If we try to access using multiple clients it will give error. If we want multiple concurrent clients we have to use a standalone metastore, i.e. Local or remote metastore configuration in Apache Hive.

Q:Suppose, I create a table that contains details of all the transactions done by the customers: CREATE TABLE transaction_details (cust_id INT, amount FLOAT, month STRING, country STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY ','; Now, after inserting 50,000 records in this table, I want to know the total revenue generated for each month. But, Hive is taking too much time in processing this query. How will you solve this problem and list the steps that I will be taking in order to do so?

As hive is taking too much time in processing this query we can partition the table transaction_details on month column so now when query runs it only loads the data from the partitions and does not scan the full data.

Steps to Partition the table:

1.CREATE TABLE transaction_details_partitioned (cust_id INT, amount FLOAT, country STRING)

PARTITIONED BY (month STRING)

ROW FORMAT DELIMITED

FIELDS TERMINATED BY ',';

2.SET hive.exec.dynamic.partition = true; SET hive.exec.dynamic.partition.mode = nonstrict;

3.Load the data into partitioned table INSERT OVERWRITE TABLE transaction_details_partitioned PARTITION (month) SELECT cust id, amount, country, month FROM transaction_details;

Q:How can you add a new partition for the month December in the above partitioned table?

ALTER TABLE transaction_details_partitioned ADD PARTITION (month='Dec') LOCATION '/transaction_details_partitioned;

Q:I am inserting data into a table based on partitions dynamically. But, I received an error – FAILED ERROR IN SEMANTIC ANALYSIS: Dynamic partition strict mode requires at least one static partition column. How will you remove this error?

SET hive.exec.dynamic.partition = true;

SET hive.exec.dynamic.partition.mode = nonstrict;

Q:Suppose, I have a CSV file – 'sample.csv' present in '/temp' directory with the following entries:

id first_name last_name email gender ip_address
How will you consume this CSV file into the Hive warehouse using built-in SerDe?

CREATE EXTERNAL TABLE sample_csv
(id int,
first_name string,
last_name string,
email string,
gender string,
ip_address string)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
STORED AS TEXTFILE LOCATION '/temp';

Now we can execute queries on this table.

Q:Suppose, I have a lot of small CSV files present in the input directory in HDFS and I want to create a single Hive table corresponding to these files. The data in these files are in the format: {id, name, e-mail, country}. Now, as we know, Hadoop performance degrades when we use lots of small files.

So, how will you solve this problem where we want to create a single Hive table for lots of small files without degrading the performance of the system?

One can use the SequenceFile format which will group these small files together to form a single sequence file. The steps that will be followed in doing so are as follows:

We need to create a sequence file to resolve this issue.

Steps to create a sequence file

• Create a temporary table:

```
CREATE TABLE temp_table (
id INT,
name STRING,
email STRING,
country STRING)
DELIMITED TERMINATED BY ','
STORED AS TEXTFILE;
```

Load the data into temp_table:

LOAD DATA INPATH '/input' INTO TABLE temp_table;

Create a table that will store data in SequenceFile format:

```
CREATE TABLE seq_file_table (id INT, name STRING, email STRING, country STRING)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS SEQUENCEFILE;
```

• Load the data from the temporary table into the sample_seqfile table:

INSERT OVERWRITE TABLE seg file table SELECT * FROM temp table;

```
LOAD DATA LOCAL INPATH 'Home/country/state/'
OVERWRITE INTO TABLE address;
The following statement failed to execute. What can be the cause?
```

As we are loading data from local, the path should contain a file path not a directory path.

Eg: LOAD DATA LOCAL INPATH 'Home/country/state/file_name.txt' OVERWRITE INTO TABLE address;

Q:ls it possible to add 100 nodes when we already have 100 nodes in Hive? If yes, how?

Yes, we can add the nodes by following the below steps:

- Step 1: Take a new system; create a new username and password
- **Step 2**: Install SSH and with the master node setup SSH connections
- **Step 3**: Add ssh public_rsa id key to the authorized_keys file

Step 4: Add the new DataNode hostname, IP address, and other details in /etc/hosts slaves file:

192.168.1.102 slave3.in slave3

Step 5: Start the DataNode on a new node

Step 6: Login to the new node like suhadoop or:

ssh -X hadoop@192.168.1.103

Step 7: Start HDFS of the newly added slave node by using the following command:

./bin/hadoop-daemon.sh start data node

Hive Practical Questions:

Hive Join operations

Create a table named CUSTOMERS(ID | NAME | AGE | ADDRESS | SALARY)

Create a Second table ORDER(OID | DATE | CUSTOMER_ID | AMOUNT)

Now perform different joins operations on top of these tables (Inner JOIN, LEFT OUTER JOIN ,RIGHT OUTER JOIN ,FULL OUTER JOIN)

cloudera@quickstart:~

```
hive> select * from customers;
OK
customers.id
                                              customers.address
               customers.name customers.age
                                                                     customers.salary
                              Gurgaon 100000
       Shashank
                                              786665
       Sudhanshu
                              Bangalore
                    Nagpur 12231
       Prajwal 23
       Vansh 15
                      Nagpur
       Prakash 35
                      Nagpur 121212
                       Delhi
                              23322
       Ram
Time taken: 0.166 seconds, Fetched: 6 row(s)
hive>
```

cloudera@quickstart:~

INNER JOIN

hive> select c.id,c.name,o.oid,o.order_date,o.amount from

- > orders o
- > join customers c
- > on o.customer id = c.id;

cloudera@quickstart:~

```
c.id
       c.name o.oid
                      o.order date o.amount
1
       Shashank
                      1
                              2022-08-10
                                             10000.0
       Sudhanshu
                      2
                              2022-09-01
                                             890.0
3
       Prajwal 3
                                     90.0
                      2022-01-01
4
       Vansh 4
                      2022-08-01
                                      8700.0
5
2
                                      1000.0
       Prakash 5
                      2022-08-10
       Sudhanshu
                      6
                              2022-09-01
                                             890.0
3
       Prajwal 7
                      2022-01-01 90.0
       Vansh 8
                      2022-08-01
                                     8700.0
Time taken: 86.685 seconds, Fetched: 8 row(s)
hive>
```

Left Outer Join

hive> Select c.id,c.name,o.oid,o.order_date,o.amount

- > from customers c
- > left join orders o
- > on c.id = o.customer_id;

cloudera@quickstart:~

```
c.id
       c.name o.oid o.order date o.amount
       Shashank
                               2022-08-10
                                               10000.0
2
       Sudhanshu
                       2
                               2022-09-01
                                               890.0
2
       Sudhanshu
                       6
                               2022-09-01
                                               890.0
3
       Prajwal 3
                       2022-01-01
                                       90.0
3
       Prajwal 7
                       2022-01-01
                                       90.0
4
       Vansh 4
                       2022-08-01
                                       8700.0
4
       Vansh
                       2022-08-01
                                       8700.0
5
       Prakash 5
                       2022-08-10
                                       1000.0
6
       Ram NULL
                       NULL
                              \mathtt{NULL}
Time taken: 43.445 seconds, Fetched: 9 row(s)
hive>
```

RIGHT OUTER JOIN

hive> Select c.id,c.name,o.oid,o.order_date,o.amount

- > from orders o
- > right join customers c
- > on o.id = c.customer_id;

```
OK
c.id
       c.name o.oid
                       o.order date o.amount
1
                       1
       Shashank
                               2022-08-10
                                             10000.0
2
                       2
       Sudhanshu
                              2022-09-01
                                              890.0
2
       Sudhanshu
                      6
                              2022-09-01
                                              890.0
3
       Prajwal 3
                                      90.0
                      2022-01-01
3
       Prajwal 7
                      2022-01-01
                                      90.0
4
       Vansh 4
                       2022-08-01
                                      8700.0
4
       Vansh
               8
                       2022-08-01
                                      8700.0
5
       Prakash 5
                       2022-08-10
                                      1000.0
               NULL
                       NULL
                             NULL
       Ram
Time taken: 37.616 seconds, Fetched: 9 row(s)
hive>
```

FULL OUTER JOIN

hive> Select c.id,c.name,o.oid,o.order_date,o.amount

- > from orders o
- > full join customers c
- > on o.customer_id = c.id;

cloudera@quickstart:~

```
OK
c.id
        c.name o.oid
                        o.order date
                                        o.amount
1
                        1
        Shashank
                                2022-08-10
                                                10000.0
2
        Sudhanshu
                        6
                                2022-09-01
                                                890.0
                                2022-09-01
        Sudhanshu
                        2
                                                890.0
3
        Prajwal 7
                        2022-01-01
                                        90.0
3
        Prajwal 3
                        2022-01-01
                                        90.0
4
                        2022-08-01
        Vansh
                                        8700.0
4
        Vansh
              4
                        2022-08-01
                                        8700.0
5
        Prakash 5
                        2022-08-10
                                        1000.0
        Ram NULL
                        NULL
                                NULL
Time taken: 93.762 seconds, Fetched: 9 row(s)
hive>
```

Download a data from the given location -

https://archive.ics.uci.edu/ml/machine-learning-databases/00360/

1. Create a hive table as per given schema in your dataset

```
hive>
hive> create table bank(
           age int,
           job string,
           marital string,
           education string,
           default string,
           balance int,
           housing string,
           loan string,
           contact string,
           day int,
           month string,
           duration int,
           campaign int,
           padays int,
         previous int,
           poutcome string,
           y string)
           row format serde 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
           with SERDEPROPERTIES (
           "separatorChar" = "\;",
           "quoteChar" = "\""
           stored as textfile;
OK
Time taken: 0.268 seconds
```

2. try to place a data into table location

```
cloudera@quickstart:~
```

```
hive> load data local inpath 'file:///home/cloudera/bank.csv' into table bank;
Loading data to table hive_class_b1.bank
Table hive_class_b1.bank stats: [numFiles=1, totalSize=461474]
OK
Time taken: 0.618 seconds
hive>
```

3. Perform a select operation.

```
Courter@Quickstart-
hives Select * from bank limit 10;

OK

Dank.age bank.job bankmarital bank.education bank.poutcome age job marital education default balance housing loan contact day month on cellular 19 oct 79 1 -1 0 unknown no

30 unemployed married perindry no 1787 no no cellular 11 may 220 1 339 4 failure no
33 services married secondary no 4789 yes yes cellular 11 may 220 1 339 4 failure no
35 management single tertiary no 1350 yes no cellular 16 apr 185 1 330 1 failure no
30 management married tertiary no 1476 yes yes unknown 3 jun 199 4 -1 0 unknown no
35 management single tertiary no 1476 yes yes unknown 5 may 226 1 -1 0 unknown no
36 self-employed married tertiary no 747 no no cellular 23 feb 141 2 176 3 failure no
36 self-employed married secondary no 147 yes no cellular 18 may 341 1 330 2 other no
39 management single tertiary no 307 yes no cellular 23 feb 141 2 176 3 failure no
47 no no cellular 18 may 341 1 330 2 other no
48 no cellular 19 yes no yes no cellular 19 yes no yes no yes no yes no yes no yes no y
```

4. Fetch the result of the select operation in your local as a csv file .

```
Cloudera@quickstart ~]$ hive -e "select * from hive_class_bl.bank" > ~/output.csv

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties

OK

Time taken: 0.632 seconds, Fetched: 4522 row(s)

[cloudera@quickstart ~]$ a

AirQualityUcl.csv country_wise_latest.csv

Desktop enterprise-deployment.json kerberos output.csv sales_order_data.csv test2.txt

array_data.csv covid_19_clean_complete.csv

Downloads full grouped.csv map_data.csv Pictures sample_data.csv usa_county_wise.csv

Dank.csv customers.csv

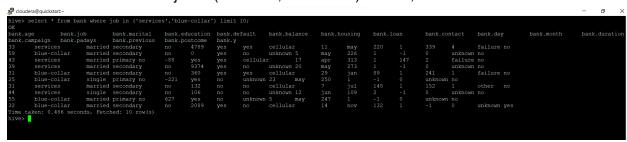
cloudera_manager day wise.csv celipse hive-heatalog-core-0.14.0.jar Music

mapi.py dept_data.csv worldometer_data.csv

[cloudera@quickstart ~]$ employee.csv json_data.json orders.csv sales_data.csv Templates worldometer_data.csv
```

- 5. Perform group by operation . hive> Select marital,count(1) cnt
 - > from bank
 - > group by marital;

7. Perform filter operation at least 5 kinds of filter examples . select * from bank where job in ('services', 'blue-collar') limit 10;



hive> select * from bank where education like 'primary' limit 10;

```
## Courage Cou
```

8. show and example of regex operation

cloudera@quickstart:~

```
hive> alter table bank rename to bank_data;
OK
Time taken: 0.295 seconds
hive> alter table bank_data change column loan loan string after balance;
OK
Time taken: 0.387 seconds
hive>
```

10.Drop table bank_data;

12 . order by operation .

```
oK
secondary 2306
tertiary 1350
primary 678
unknown 187
education 1
Time taken: 72.597 seconds, Fetched: 5 row(s)
hive>
```

14 . sorting operation you have to perform .

```
cloudera@quickstart:~
```

```
hive> Select marital, count(1) cnt
    > from bank data
    > group by marital
    > sort by cnt desc;
Query ID = cloudera 20220916034646 34497e96-b40a-4895-8d2a-4c4754d4a1b5
Total iobs = 2
```

```
Total MapReduce CPU Time Spent: 16 seconds 870 msec
OK
secondary
                2306
tertiary
                1350
primary 678
unknown 187
education
Time taken: 58.047 seconds, Fetched: 5 row(s)
```

15. distinct operation you have to perform. Select distinct job from bank data;

```
dickstart:~
hive> select distinct job from bank_data;
Query ID = cloudera_20220916034848_93f03fa3-ea4d-442b-8e37-5e0562a51cb0
Total jobs = 1
 Launching Job 1 out of 1
 Number of reduce tasks not specified. Estimated from input data size: 1
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
set mapreduce.job.reduces=<number>
Starting Job = job_1663312363320_0013, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1663312363320_0013/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1663312363320_0013
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
Hadoop Job Information for Stage-1: number of mappers: 1; number of reducers: 1
2022-09-16 03:48:29,469 Stage-1 map = 0%, reduce = 0%
2022-09-16 03:48:38,390 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 7.93 sec
2022-09-16 03:48:47,942 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 10.06 sec
MapReduce Total cumulative CPU time: 10 seconds 60 msec
Ended Job = job_1663312363320_0013
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 10.06 sec HDFS Read: 470444 HDFS Write: 126 SUCCESS Total MapReduce CPU Time Spent: 10 seconds 60 msec
admin.
 blue-collar
 nousemaid
 iob
management
 self-employed
 services
 student
 Time taken: 31.514 seconds, Fetched: 13 row(s)
 hive>
```

- 16 . like an operation you have to perform . Select * from bank_data where job like 'student';
- 17 . union operation you have to perform .

18 . table view operation you have to perform .



Q:Create a python application that connects to the Hive database for extracting data, creating sub tables for data processing, drops temporary tables.fetch rows to python itself into a list of tuples and mimic the join or filter operations

https://github.com/PrajwalMahale/hive/blob/main/main.py