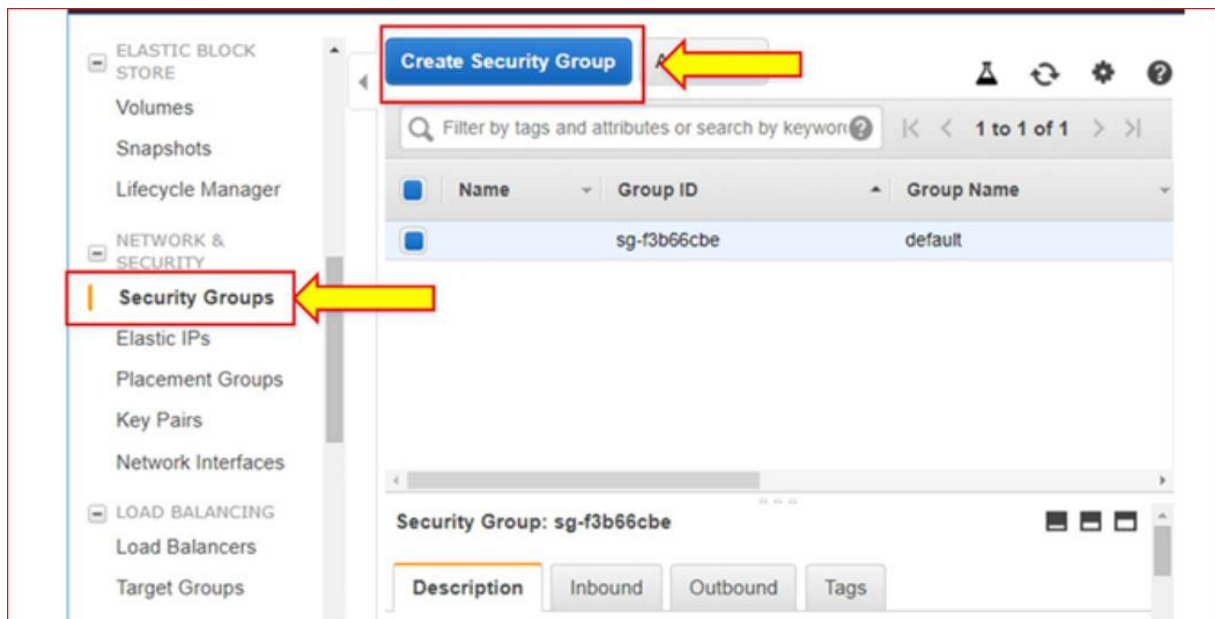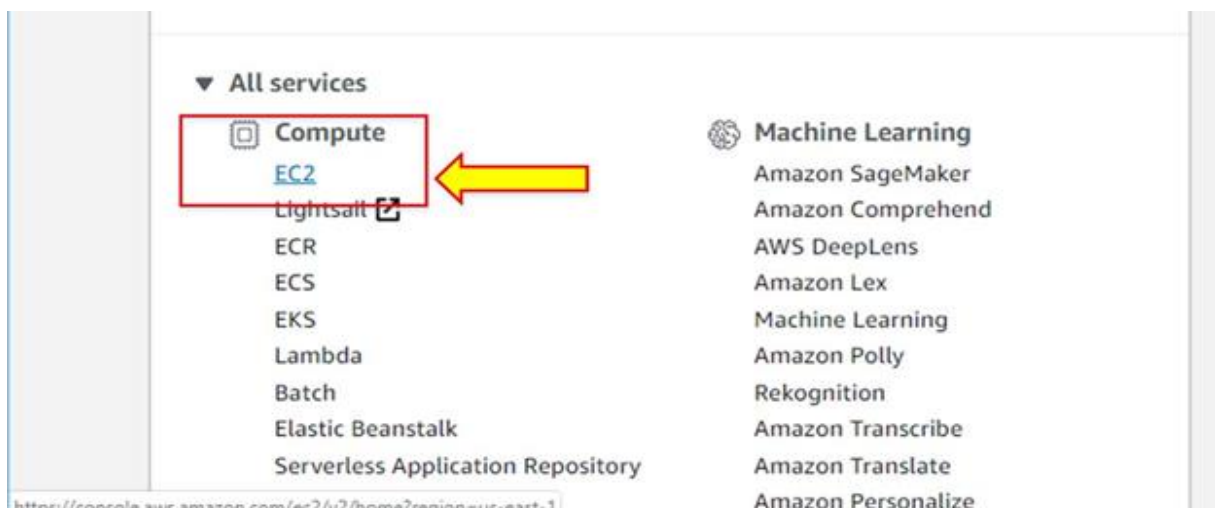**SOLUTION:**

# BIG DATA ANALYTICS FOR ANALYZING LARGE DATA SETS

**Steps:**

For the requirement, they use Elastic MapReduce tool for analyzing the process and in that use Hive for querying.

**Step 1:** First, create a Security group with appropriate SSH rule for EMR master node.

**Edit inbound rules**

| Type ⓘ | Protocol ⓘ | Port Range ⓘ | Source ⓘ | |
|---|---|---|---|---|
| SSH ▾ | TCP | 22 | Anywhere ▾ | 0.0.0.0/0, ::/0 |

Add Rule

NOTE: Any edits made on existing rules will result in the edited rule being deleted and a new rule created with the new
on that rule to be dropped for a very brief period of time until the new rule can be created.



Create

**Step 2:** Create an S3 bucket which is public to store the table.





**Step 3:** Go to EMR under Analytics on AWS Console

**Step 4:** Click on Create Cluster and go to advanced options.

**Step 5:** Go to the latest EMR release. Check Hadoop, Hive and Spark as we are going to use these tools. Click Next.



**Step 6:** In Hardware configuration, select the VPC which has public subnet. Here we are using default VPC.

**Step 7:** Now go to the node and select the type for master node and its purchasing option.

There is no need for a core node and task node.



**Step 8:** Click Next

**Step 9:** Enter the cluster name and leave everything as it is. Click Next.

**Step 10:** In security option, select EC2 key pair as we want that to SSH the master node. Here we have used the key which already created in EC2 service.



**Step 11:** Leave everything as it is and in the security group of master node add additional security group. Here we add that Security group that we created in the start for SSH. Click on Assign Security group.

**Step 12:** Click on Create Cluster.

**Step 13:** Now cluster has started creating. Once it is created, take public DNS name of the cluster for SSH.



The cluster status will turn to waiting.

**Step 14:** For Windows Users use PuttyKeyGen to generate a PPK file out of the PEM file and login using putty.

Click on the SSH link as shown in the picture:



Go to terminal here using Linux. Go to SSH via pem key and DNS name to connect to master node.

For MAC and LInux users:

 "chmod 400 new.pem "

 "ssh -i new.pem hadoop@<DNS name>"

**Step 15:** Now you are connected to master node. For using Hive, type 'hive'. Now in Hive, create a table which for which the storage is S3. Create a folder called "bigdatatest" in the S3 bucket.

```
hive>
    >
    >
    >
    > create external table BigdataSampleTable1
    >  (O_ORDERKEY INT,
    >  O_CUSTKEY INT,
    >  O_ORDERSTATUS STRING,
    >  O_TOTALPRICE DOUBLE,
    >  O_ORDERDATE STRING,
    >  O_ORDERPRIORITY STRING,
    >  O_CLERK STRING,
    >  O_SHIPPRIORITY INT,
    >  O_COMMENT STRING)
    >  ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'
    >  LOCATION 's3://bigdatalabbucketnov13/bigdatatest/';
OK
Time taken: 4.864 seconds
hive> insert into BigdataSampleTable1 values('001', '123', 'complete', '1500', '
14-11-19', 'High', 'Data', '1', 'OrderTable');
Query ID = hadoop_20191113183807_21eeee72-14b9-49ec-9d40-6eaac87d5d92
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1573669394512
_0002)

Map 1: 0/1
Map 1: 0(+1)/1
Map 1: 0(+1)/1
Map 1: 1/1
Loading data to table default.bigdatasampletable1
OK
Time taken: 12.681 seconds
hive>
```
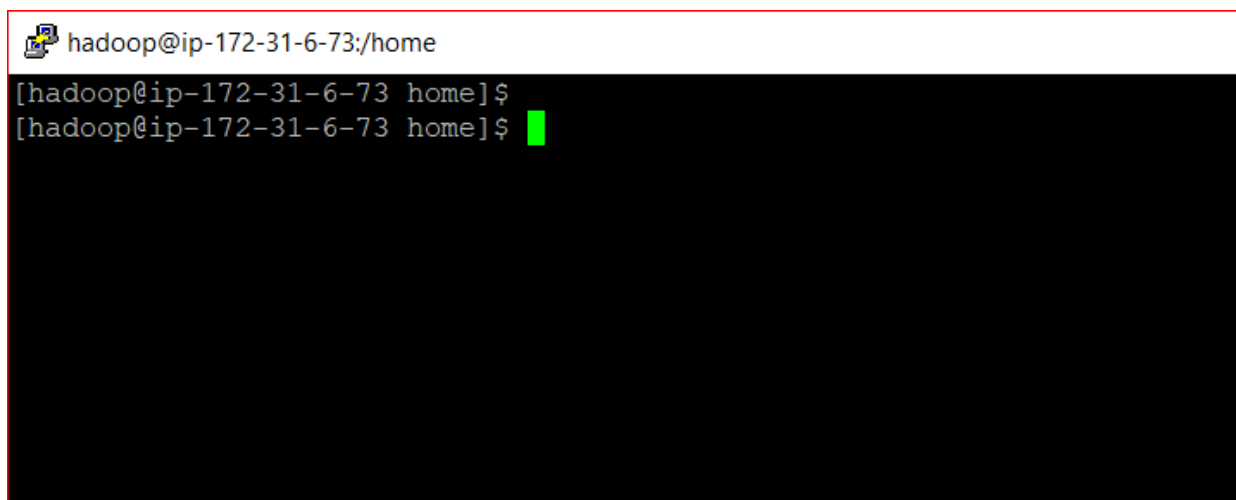
*****

create external table BigdataSampleTable1

 (O_ORDERKEY INT,

 O_CUSTKEY INT,

 O_ORDERSTATUS STRING,

 O_TOTALPRICE DOUBLE,

 O_ORDERDATE STRING,

 O_ORDERPRIORITY STRING,

 O_CLERK STRING,

 O_SHIPPRIORITY INT,

 O_COMMENT STRING)

ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'

LOCATION 's3://bigdatalabbucketnov13/bigdatatest/';

****

Load the Data to the first table : BigdataSampleTable1

****

insert into BigdataSampleTable1 values('001', '123', 'complete', '1500', '14-11-19', 'High', 'Data',

'1', 'OrderTable');

*****

```
hive> insert into BigdataSampleTable1 values('001', '123', 'complete', '1500', '
14-11-19', 'High', 'Data', '1', 'OrderTable');
Query ID = hadoop_20191113183807_21eeee72-14b9-49ec-9d40-6eaac87d5d92
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1573669394512
_0002)

Map 1: 0/1
Map 1: 0(+1)/1
Map 1: 0(+1)/1
Map 1: 1/1
Loading data to table default.bigdatasampletable1
OK
Time taken: 12.681 seconds
hive>
```

**Step 16:** Create another table in the same way.

```
hive>
    >
    > create external table BigdataSampleTable2
    > (L_ORDERKEY INT,
    >  L_PARTKEY INT,
    >  L_NAME STRING)
    >  ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'
    >  LOCATION 's3://bigdatalabbucketnov13/bigdatatest/';
OK
Time taken: 0.159 seconds
```

*****

create external table BigdataSampleTable2

(L_ORDERKEY INT,

 L_PARTKEY INT,

 L_NAME STRING)

 ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'

 LOCATION 's3://bigdatalabbucketnov13/bigdatatest/';

 *****

```
hive>  insert into BigdataSampleTable2 values('001', '123', 'complete');
Query ID = hadoop_20191113184245_76d41bb4-096e-4cd0-90f8-61e54c7a2d27
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1573669394512
_0002)

Map 1: 0/1
Map 1: 0(+1)/1
Map 1: 0(+1)/1
Map 1: 1/1
Loading data to table default.bigdatasampletable2
OK
Time taken: 8.137 seconds
```

Load the Data to the first table : BigdataSampleTable2

****

 insert into BigdataSampleTable2 values('001', '890', 'complete');

****

**Step 17:** Now both tables have been created. Run query on tables

***

select * from BigdataSampleTable1 limit 1;

****

Now you will get the output with the time it takes to run the query.

```
hive>
    >
    >
    > select * from BigdataSampleTable1 limit 1;
OK
1       123      complete         1500.0  14-11-19          High     Data    1         OrderTable
Time taken: 1.726 seconds, Fetched: 1 row(s)
hive>
```

Additional tasks:

1. You can check the EMR logs in the s3 bucket and path mentioned in the LogURI.



2. Check the Web Interfaces Hosted on this Cluster
   a. Hadoop, Ganglia, and other applications publish user interfaces as websites hosted on the master node

3. You can protect view access to all users by changing the below settings.



4. Any time you login to the applications installed in the EMR, it can be monitored under "Application history".