# Hive Case Study
# By
# Nishant Raj And Guresh Kumar

We have completed the case study in 2 parts.
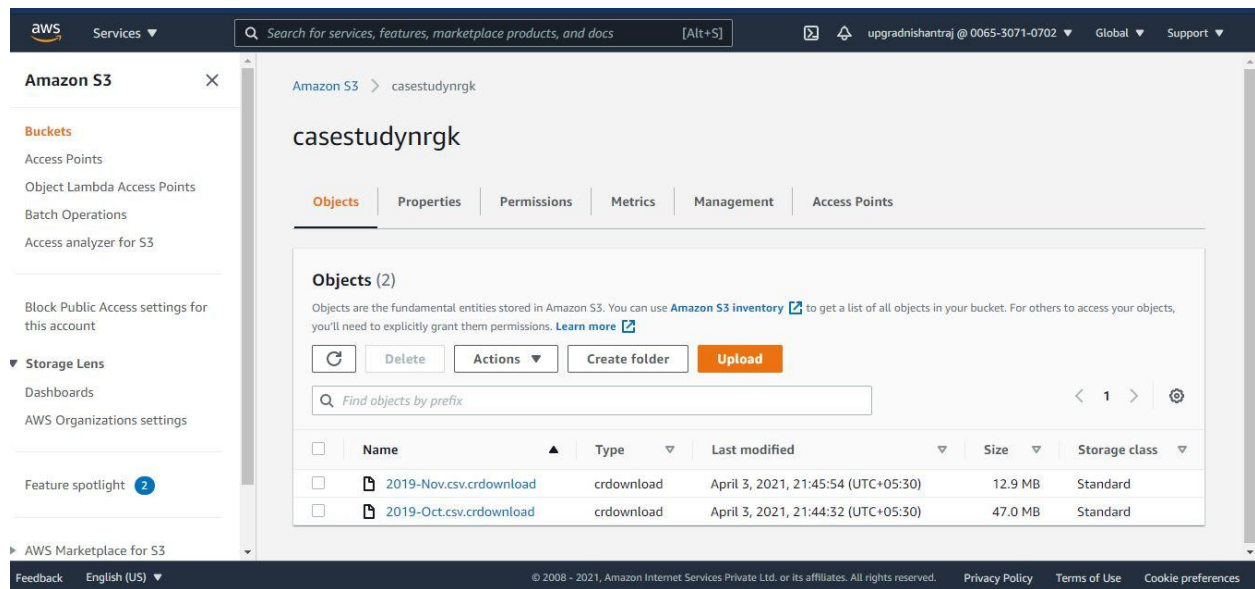1. Working with S3 and EMR clusters.
2. Querying with HIVE (HQL)

## 1. Working with S3 and EMR clusters:
    (a) Load data into S3 bucket

The first step is to create a S3 bucket and then load Data in it. We are provided with links through which we can get the data.

**https://e-commerce-events-ml.s3.amazonaws.com/2019-Oct.csv**

**https://e-commerce-events-ml.s3.amazonaws.com/2019-Nov.csv**



We have loaded the S3 bucket with 2 given data sets that we need to work on,
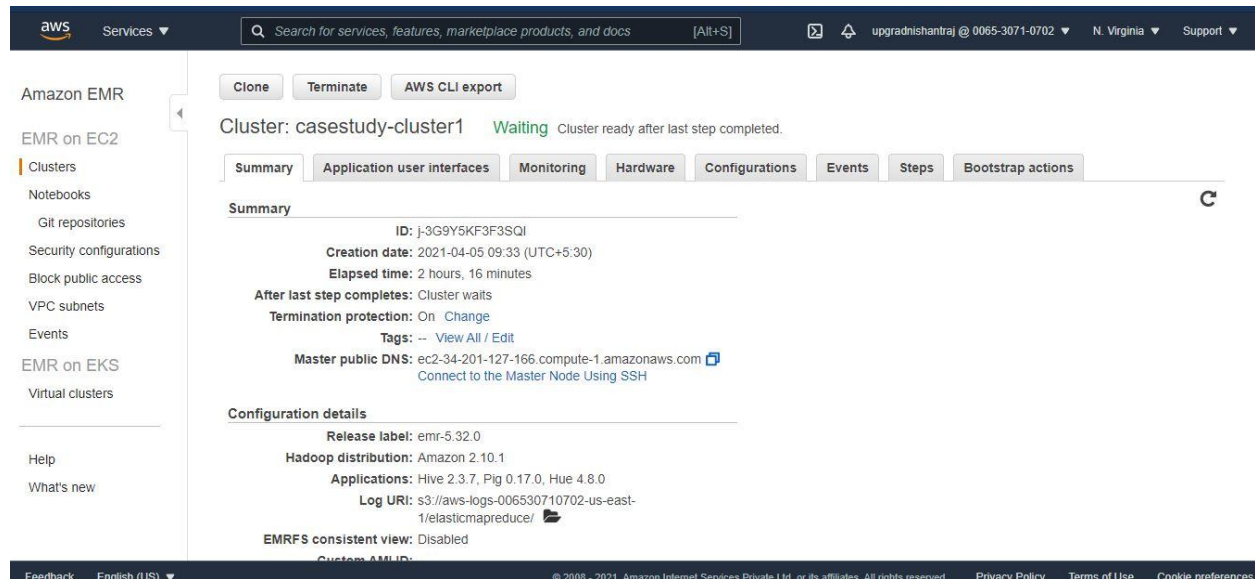
Data set 1: 2019-Nov.csv.crdownload and
Data set 2: 2019-Oct.csv.crdownload.

We need to extract data and gather insights from these data sets. In the next step will see how the data in S3 bucket is moved into the HDFS.

    (b)  Launch an EMR cluster

We have created an Elastic MapReduce cluster with the suggested configuration. We have named our cluster as casestudy-cluster1. We have created the cluster with release version 5.32.0 and used two-node cluster [1 Master M4.large and 1 Core M4.large]



We'll utilize Hive services with this EMR cluster and give instructions through commands using PuTTy.

    (c)  Copy Data from EMR cluster to HDFS

We have used **Hadoop distcp** to copy files in the S3 bucket to the Hadoop file system. This job runs like a MapReduce job in the background to copy the files from EMR cluster to HDFS.

load data local inpath "/home/hadoop/2019-Oct.csv.crdownload" into table oct_info;

load data local inpath "/home/hadoop/2019-Nov.csv.crdownload" into table nov_info;
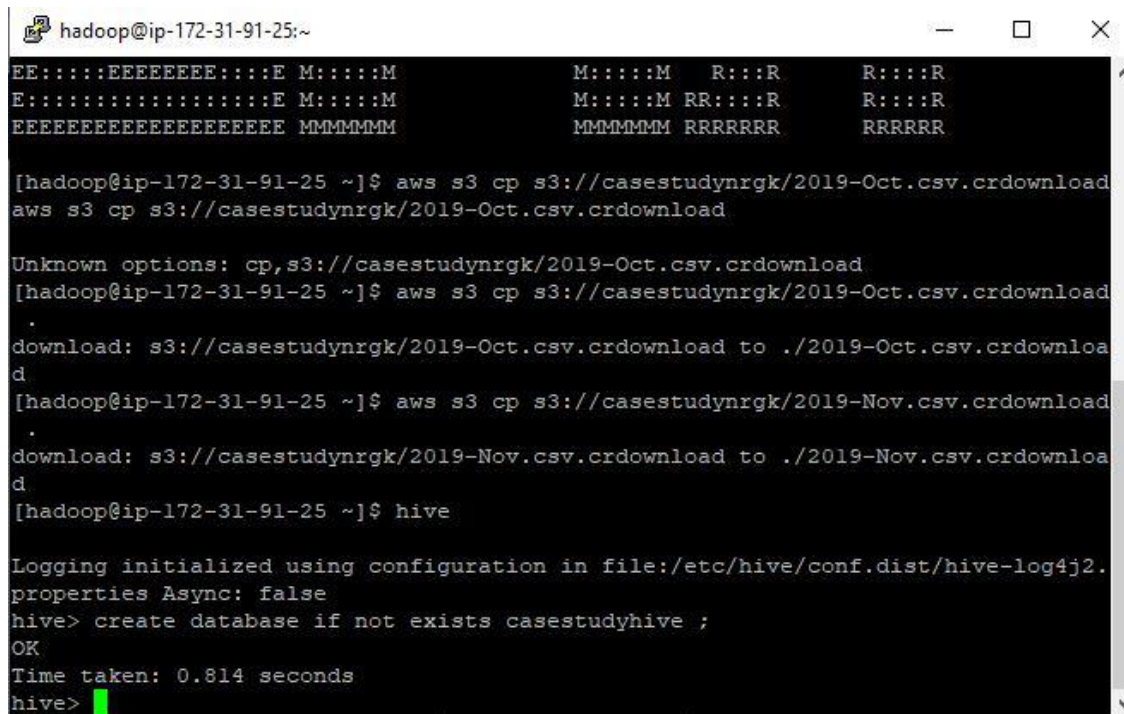
We have used the above commands to achieve the task.

## 2. Querying with HIVE (HQL):
   (a) Structure of your database
Step 1: Create a Database in Hadoop

As we can see from the screenshot the data from S3 has been loaded to Hadoop File System and the database with the name of casestudyhive is created.



Step 2: Create required tables in the Database

After creating the database we have created the tables so that data could be loaded into the tables. Here we have created two tables namely oct_info and nov_info which will be used in the subsequent screenshots for querying .

**Step 3: Load Data into tables**

After creating the tables, now we have added data into the tables we can see clearly that the data is perfectly loaded into the tables.

(b) Select an optimized technique to run your queries as efficiently as possible

We have 2 optimization techniques to run a query. One is Partitioning and the other one is Bucketing.
In Partitioning we'll divide the data present in tables into several parts based on columns and the conditions that we apply using the partition keys. In Bucketing, the partitioned data is further subdivided into buckets based on Hash function. Thus, Bucketting helps us with an optimized solution for our queries.

(c) Applying both the techniques for a query to see the performance
**\*\* We haven't captured the screenshot for the query and terminated the cluster. We have mentioned the code that we have used to check the performance.**

```
create table if not exists part_oct_info( event_time string , event_type
string , product_id string ,category_id string , category_code string ,
price float , user_id bigint , user_session string) partitioned by (brand
string) row farmat delimited by "," lines terminated by "\n";


insert into table part_oct_info partition(category_code) select event_time
string , event_type string , product_id string ,category_id string , brand
string , price float , user_id bigint , user_session string;

exit;

hadoop fs -ls "user/hive/warehouse/part_oct_info"

SELECT user_id , sum(price) from oct_info where brand="xyz" group by
category_code;

SELECT user_id , sum(price) from part_oct_info where brand="xyz" group by
category_code;

=============================================

BUCKETING

create table if not exists buck_oct_info( event_time string , event_type
string , product_id string ,category_id string , category_code string ,
price float , user_id bigint , user_session string) partitioned by (brand
string) clustered by (category_code) into 10 buckets row farmat delimited
by "," lines terminated by "\n";
```

```
insert into table buck_oct_info partition(category_code) select event_time
string , event_type string , product_id string ,category_id string , brand
string , price float , user_id bigint , user_session string;

hadoop fs -ls /user/hive/warehouse/buck_oct_info;
```

We can clearly see that bucketing will significantly reduce the query
time.

### (d) Run Hive queries to answer the Case study questions

## 1. Find the total revenue generated due to purchases made in October.
We have used SELECT SUM(PRICE) FROM oct_info; command to solve the query.

**2. Write a query to yield the total sum of purchases per month in a single output.**

The below query produces the desired output.



**3. Write a query to find the change in revenue generated due to purchases from October to November.**

The below query produces the desired output.

## 4. Find distinct categories of products. Categories with null category code can be ignored.

The below query produces the desired output.



## 5. Find the total number of products available under each category.

The below query produces the desired output.

## 6. Which brand had the maximum sales in October and November combined?



## 7. Which brands increased their sales from October to November?

The below query produces the desired output.

```
hadoop@ip-172-31-91-25:~                                                    —  □  ×
Time taken: 23.915 seconds, Fetched: 224 row(s)
hive> SELECT SUM(price), brand  FROM oct_info GROUP BY brand;
Query ID = hadoop_20210405094910_e9fde009-ace4-4639-9ebc-dea54af23557
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1617595861253_0010)

--------------------------------------------------------------------------------
        VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------
Map 1 .......... container    SUCCEEDED    3      3        0        0        0       0
Reducer 2 ...... container    SUCCEEDED    2      2        0        0        0       0
--------------------------------------------------------------------------------
VERTICES: 02/02  [==========================>>] 100%  ELAPSED TIME: 15.66 s
--------------------------------------------------------------------------------
OK
4412.549980640411       almea
3613.059979915619       ardell
4120.210026979446       art-visage
4643.859920978546       artex
56.03999996185303       australis
20.65000057220459       barbie
963.2300157546997       batiste
30217.169900417328      beautix
1424.4499959945679      beautyblender
6668.800020098686       bioaqua
695.1200037002563       biore
60.31999969482422       blise
121.29999768733978      blixz
NULL    brand
31954.33073091507       browxenna
16.649999618530273      busch
382.99999713897705      carmex
24536.669979840517      concept
352.1999931335449       consly
1873.8199858665466      coocla
```

## 8. Your company wants to reward the top 10 users of its website with a Golden Customer plan. Write a query to generate a list of top 10 users who spend the most.

The below query produces the desired output.



```
hive> WITH new_table AS( SELECT * FROM oct_info UNION SELECT * FROM nov_info) SE
LECT SUM(price) AS Highest ,user_id FROM new_table GROUP BY user_id ORDER BY Hig
hest DESC LIMIT 10;
Query ID = hadoop_20210405062431_c3552e8e-9736-4f88-ade0-0e4a8231ace9
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1617595861253
_0005)

Map 1: 0/3      Map 5: 0/1      Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 0/3      Map 5: 0/1      Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 0(+1)/3  Map 5: 0/1      Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 0(+3)/3  Map 5: 0(+1)/1  Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 1(+2)/3  Map 5: 0(+1)/1  Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 1(+2)/3  Map 5: 0(+1)/1  Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 1(+2)/3  Map 5: 0(+1)/1  Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 1(+2)/3  Map 5: 0(+1)/1  Reducer 3: 0/2  Reducer 4: 0/1
Map 1: 1(+2)/3  Map 5: 1/1      Reducer 3: 0(+1)/2      Reducer 4: 0/1
Map 1: 2(+1)/3  Map 5: 1/1      Reducer 3: 0(+1)/2      Reducer 4: 0/1
Map 1: 2(+1)/3  Map 5: 1/1      Reducer 3: 0(+2)/2      Reducer 4: 0/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 0(+2)/2      Reducer 4: 0/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 0(+2)/2      Reducer 4: 0/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 0(+2)/2      Reducer 4: 0/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 1(+1)/2      Reducer 4: 0(+1)/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 2/2  Reducer 4: 0(+1)/1
Map 1: 3/3      Map 5: 1/1      Reducer 3: 2/2  Reducer 4: 1/1
OK
5741.189976811409       546619064
4927.870005518198       553074123
4393.639978170395       556045950
4368.31982421875        463764281
4161.389979839325       475358605
3986.2999315559864      476180402
3819.020003736019       541475200
3784.1000034809113      501203199
3580.310007929802       551985809
3405.979974269867       512161964
Time taken: 30.911 seconds, Fetched: 10 row(s)
hive> █
```

## * Cleaning up:
   (a) Dropping Database
       The below query Drops off the Database.

(b) Terminating the Cluster

The below query terminates the Database.

# Thank you!