# HIVE CASE STUDY (Retail Store)
## (DSC41)
### Submitted by : SaiTeja , Kailash

## PROBLEM STATEMENT:

With online sales gaining popularity, tech companies are exploring ways to improve their sales by analyzing customer behavior and gaining insights about product trends. Furthermore, the websites make it easier for customers to find the products they require without much scavenging. Needless to say, the role of big data analysts is among the most sought-after job profiles of this decade. Therefore, as part of this assignment, we will be challenging you, as a big data analyst, to extract data and gather insights from a real-life data set of an e-commerce company.

## OBJECTIVE:

The aim is to extract data and gather insights from a real-life data set of an e-commerce Company.

## DATA:

The data used for this assignment is a public clickstream dataset of a cosmetic store. The clickstream data contains all the logs as to how one navigated through the e-commerce website. It also contains other data such as customer time spent on every page, a number of
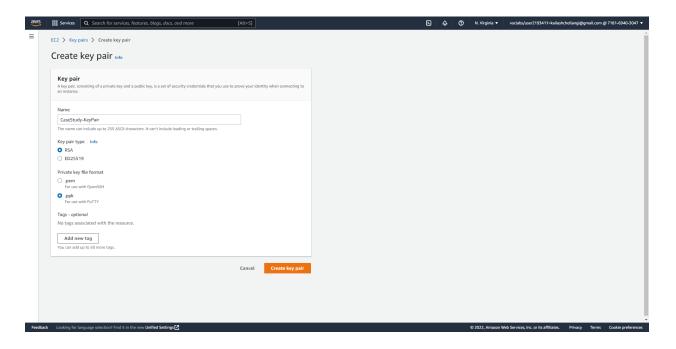clicks made, adding items to the cart, customer id, etc.



DATA INGESTION FRAMEWORKS

## OVERVIEW OF STEPS:

➢ Copying the data set into HDFS:
  ● Launch an EMR cluster that utilizes the hive services, and
  ● Move the data from S3 bucket into the HDFS
➢ Creating the database and launching hive queries on your EMR cluster:
  ● Create the structure of your database,
  ● Use optimized techniques to run your queries as effectively as possible
  ● Show the improvement in performance after optimizing
  ● Run hive queries to answer the given questions.
➢ Cleaning up:
  ● Drop your database and
  ● Terminate your cluster

## KEY-PAIR CREATION:
Creating Key Pair with ppk file format directly to use with putty





**CaseStudy-KeyPair** successfully created and downloaded (refer above screenshot)

## S3 BUCKET:

To Store the data – Click on "Create Bucket"



Creating "**hivecasestudy-teja-kailash**" with all default options



Bucket Successfully created.

Successfully uploaded the 2019 October and 2019 November csv files to S3 bucket

# EMR CLUSTER CREATION:

## Click on "Create cluster" button to create the EMR cluster



## Creating cluster with advanced options



**Software and Steps page**: Changed the Release from emr-5.36.0 to "**emr-5.29.0**"

**Hardware page**: Changing the Master and Core nodes from m5.xlarge to "**m4.large**"

*Specify target capacity and how Amazon EMR fulfills it for each node type. Mix instance types and purchasing options. Learn more ⧉*

## Networking

Use a Virtual Private Cloud (VPC) to process sensitive data or connect to a private network. Launch the cluster into a VPC with a public, private, or shared subnet. Subnets may be associated with and AWS Outpost or AWS Local Zone.

Launch the cluster into a VPC with a public, private, or shared subnet. Subnets may be associated with an AWS Outpost or AWS Local Zone.

| | |
|---|---|
| Network | vpc-023aa35c21ef643b6 (172.31.0.0/16) (default) ▾    Create a VPC ⧉ ⓘ |
| EC2 Subnet | subnet-00b8049369443263f \| Default in us-east-1c ▾ |

## Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. Learn more about instance purchasing options ⧉

ⓘ Console options for automatic scaling have changed. Learn more ⧉                                              ✕

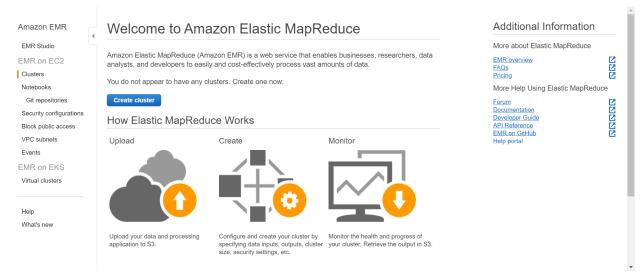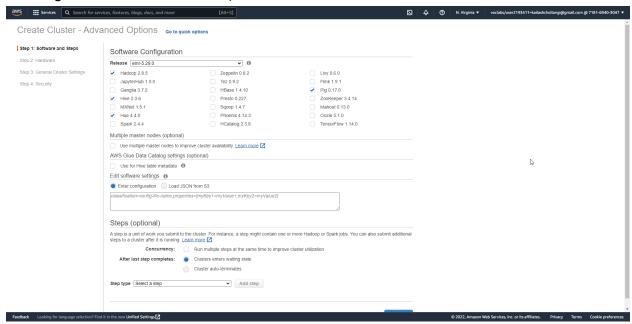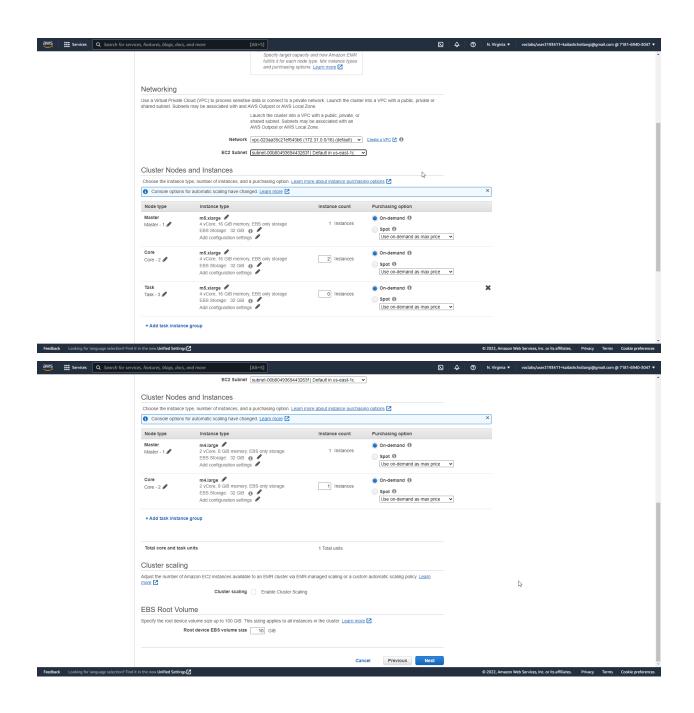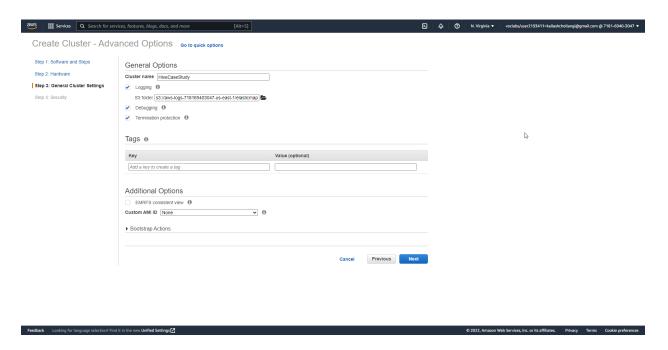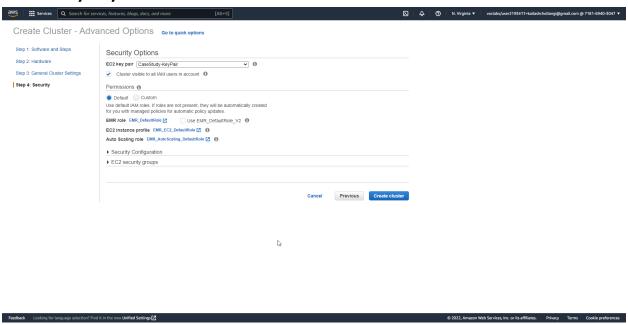| Node type | Instance type | Instance count | Purchasing option |
|---|---|---|---|
| **Master**<br>Master - 1 ✎ | **m5.xlarge** ✎<br>4 vCore, 16 GiB memory, EBS only storage<br>EBS Storage:  32 GiB  ⓘ ✎<br>Add configuration settings ✎ | 1   Instances | ⦿ On-demand ⓘ<br>○ Spot ⓘ<br>Use on-demand as max price ▾ |
| **Core**<br>Core - 2 ✎ | **m5.xlarge** ✎<br>4 vCore, 16 GiB memory, EBS only storage<br>EBS Storage:  32 GiB  ⓘ ✎<br>Add configuration settings ✎ | 2   Instances | ⦿ On-demand ⓘ<br>○ Spot ⓘ<br>Use on-demand as max price ▾ |
| **Task**<br>Task - 3 ✎ | **m5.xlarge** ✎<br>4 vCore, 16 GiB memory, EBS only storage<br>EBS Storage:  32 GiB  ⓘ ✎<br>Add configuration settings ✎ | 0   Instances | ⦿ On-demand ⓘ<br>○ Spot ⓘ<br>Use on-demand as max price ▾ | ✖ |

+ Add task instance group

---

| | |
|---|---|
| EC2 Subnet | subnet-00b8049369443263f \| Default in us-east-1c ▾ |

## Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. Learn more about instance purchasing options ⧉

ⓘ Console options for automatic scaling have changed. Learn more ⧉                                              ✕

| Node type | Instance type | Instance count | Purchasing option |
|---|---|---|---|
| **Master**<br>Master - 1 ✎ | **m4.large** ✎<br>2 vCore, 8 GiB memory, EBS only storage<br>EBS Storage:  32 GiB  ⓘ ✎<br>Add configuration settings ✎ | 1   Instances | ⦿ On-demand ⓘ<br>○ Spot ⓘ<br>Use on-demand as max price ▾ |
| **Core**<br>Core - 2 ✎ | **m4.large** ✎<br>2 vCore, 8 GiB memory, EBS only storage<br>EBS Storage:  32 GiB  ⓘ ✎<br>Add configuration settings ✎ | 1   Instances | ⦿ On-demand ⓘ<br>○ Spot ⓘ<br>Use on-demand as max price ▾ |

+ Add task instance group

| | |
|---|---|
| **Total core and task units** | 1 Total units |

## Cluster scaling

Adjust the number of Amazon EC2 instances available to an EMR cluster via EMR-managed scaling or a custom automatic scaling policy. Learn more ⧉

**Cluster scaling**   ☐ Enable Cluster Scaling

## EBS Root Volume

Specify the root device volume size up to 100 GiB. This sizing applies to all instances in the cluster. Learn more ⧉

**Root device EBS volume size**   10  GiB

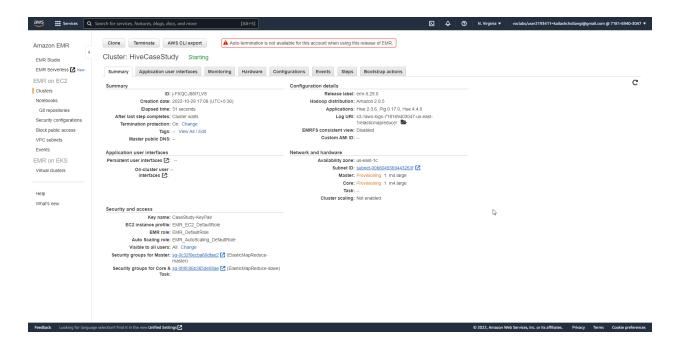Cancel   Previous   **Next**

**General Cluster Settings page**: Giving the name to cluster "HiveCaseStudy"
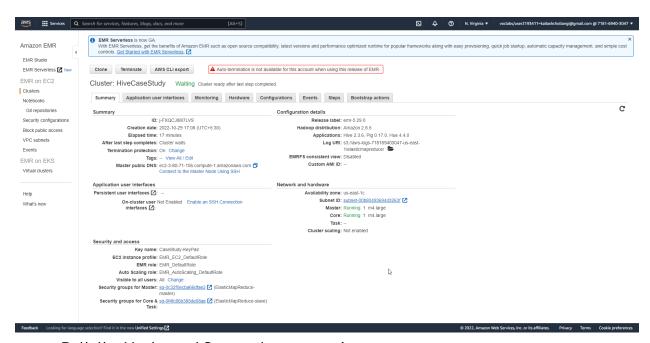


**Security page**: changing the EC2 key pair option to our created key pair – "**CaseStudy-KeyPair**"
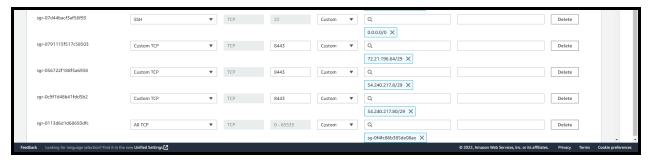


Click on "Create Cluster" button

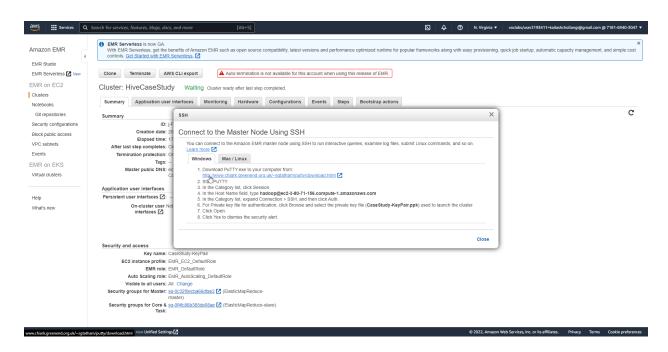Cluster is ready with status "Waiting"



- Both the Master and Core nodes are running
- We need to make sure before connecting to SSH, ensure that the port is open to establish a connection. For this, click on Security groups for Master node.
- Click on edit Inbound rules.
- Add a new rule by selecting SSH and change the IP address to Anywhere

- Then save the SSH rule to the inbound rules
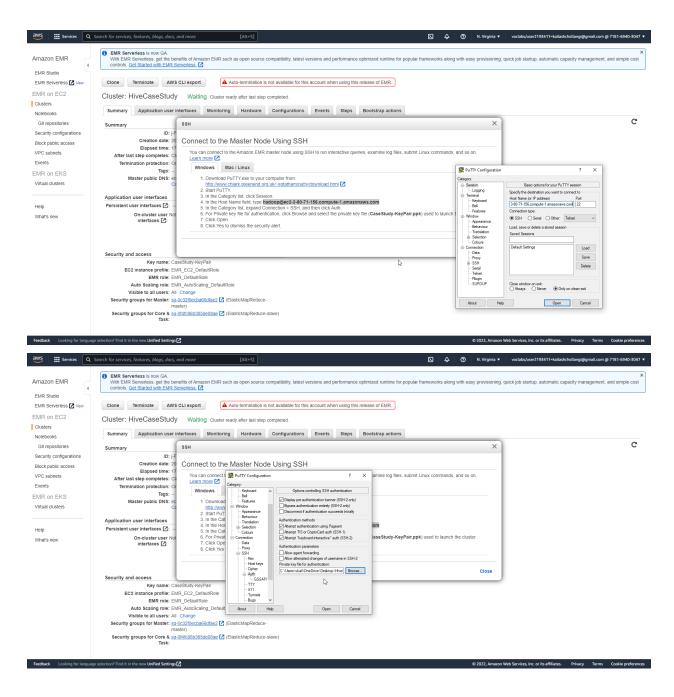
## CONNECT TO MASTER NODE:

Open the putty and enter the Host Name as hadoop@ec2-3-80-71-156.compute-1.amazonaws.com and navigate to Connection > SSH > Auth then browse and select the private key, which we created initially.

Click on "open" and then Accept the connection

EMR CLI is launched

Verifying the services that are running on Hadoop cluster with command "**sudo initctl list**"

We can see that Hive services are running

Verifying the Hadoop file system with command "**hadoop fs -ls /**"



All the above are inbuilt directories in HDFS.

CREATING A NEW DIRECTORY FOR HIVE CASE STUDY:

Creating a new directory under user>hive for Hive case study to store the data files and directory name creating is "**hive-casestudy**" and verifying whether the new directory is listed in Hadoop file system>user>hive

hadoop fs -mkdir /user/hive/hive-casestudy
hadoop fs -ls /user/hive/



New directory is successfully created

## LOADING THE DATA FROM S3 BUCKET to HDFS:
Copying the file path from S3



Distributed copy command is using to copy the data from S3 to HDFS –

For 2019 October:

**hadoop distcp s3n://hivecasestudy-teja-kailash/2019-Oct.csv
/user/hive/hive-casestudy/2019-Oct.csv**

For 2019 November:

**hadoop distcp s3n://hivecasestudy-teja-kailash/2019-Nov.csv
/user/hive/hive-casestudy/2019-Nov.csv**

Below are the screenshots for copying October 2019 and November 2019 data individually

October 2019:

## November 2019





Verifying whether the data is successfully copied into HDFS from S3 buckets

Command: **hadoop fs -ls /user/hive/hive-casestudy**

Inspecting the table data to know which columns are available before creating the hive table with command **"hadoop fs -cat /user/hive/hive-casestudy/2019-Oct.csv |head" and "hadoop fs -cat /user/hive/hive-casestudy/2019-Nov.csv |head"**

```
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -cat /user/hive/hive-casestudy/2019-Oct.csv |head
event_time,event_type,product_id,category_id,category_code,brand,price,user_id,user_session
2019-10-01 00:00:00 UTC,cart,5773203,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:03 UTC,cart,5773353,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:07 UTC,cart,5881589,2151191071051219817,,lovely,13.48,429681830,49e8d843-adf3-428b-a2c3-fe8bc6a307c9
2019-10-01 00:00:07 UTC,cart,5723490,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:15 UTC,cart,5881449,1487580013522845895,,lovely,0.56,429681830,49e8d843-adf3-428b-a2c3-fe8bc6a307c9
2019-10-01 00:00:16 UTC,cart,5857269,1487580005134238553,,runail,2.62,430174032,73deale7-664e-43f4-8b30-d32b9d5af04f
2019-10-01 00:00:19 UTC,cart,5739055,1487580008246412266,,kapous,4.75,377667011,81326ac6-daa4-4f0a-b488-fd0956a78733
2019-10-01 00:00:24 UTC,cart,5825598,1487580009445982239,,,0.56,467916806,2f5b5546-b8cb-9ee7-7ecd-84276f8ef486
2019-10-01 00:00:25 UTC,cart,5698989,1487580006317032337,,,1.27,385985999,d30965e8-1101-44ab-b45d-cclbb9fae694
cat: Unable to write to output stream.
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -cat /user/hive/hive-casestudy/2019-Nov.csv |head
event_time,event_type,product_id,category_id,category_code,brand,price,user_id,user_session
2019-11-01 00:00:02 UTC,view,5802432,1487580009286598681,,,0.32,562076640,09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:09 UTC,cart,5844397,1487580006317032337,,,2.38,553329724,2067216c-31b5-455d-alcc-af0575a34ffb
2019-11-01 00:00:10 UTC,view,5837166,1783999064103190764,,pnb,22.22,556138645,57ed222e-a54a-4907-9944-5a875c2d7f4f
2019-11-01 00:00:11 UTC,cart,5876812,1487580010100293687,,jessnail,3.16,564506666,186c1951-8052-4b37-adce-dd9644bld5f7
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,,3.33,553329724,2067216c-31b5-455d-alcc-af0575a34ffb
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,,3.33,553329724,2067216c-31b5-455d-alcc-af0575a34ffb
2019-11-01 00:00:25 UTC,view,5856189,1487580009026551821,,runail,15.71,562076640,09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:32 UTC,view,5837835,1933472286753424063,,,3.49,514649199,432a4e95-375c-4b40-bd36-0fc039e77580
2019-11-01 00:00:34 UTC,remove_from_cart,5870838,1487580007675986893,,milv,0.79,429913900,2f0bff3c-252f-4fe6-afcd-5d8a6a92839a
cat: Unable to write to output stream.
[hadoop@ip-172-31-20-212 ~]$
```

Both the tables are having same columns of data

Moving to hive:

```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive>
```

## CREATING AN EXTERNAL TABLE IN HIVE:

CREATE EXTERNAL TABLE IF NOT EXISTS retailsstore (event_time timestamp, event_type string, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS TEXTFILE LOCATION '/user/hive/hive-casestudy' tblproperties("skip.header.line.count"="1");

```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> CREATE EXTERNAL TABLE IF NOT EXISTS retailsstore (event_time timestamp, event_type string, product_id string, cate
E 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS TEXTFILE LOCATION '/user/hive/hive-casestudy' tblproperties("sk
OK
Time taken: 1.823 seconds
hive>
```

```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> CREATE EXTERNAL TABLE IF NOT EXISTS retailsstore (event_time timestamp, event_type string, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string) ROW FORMAT SERD
E 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS TEXTFILE LOCATION '/user/hive/hive-casestudy' tblproperties("skip.header.line.count"="1");
OK
Time taken: 1.823 seconds
hive>
```

Below command is used to set the display the header columns

**set hive.cli.print.header = true;**

## APPLYING OPTIMIZATION TECHNIQUES - PARTITIONING AND BUCKETING:

Below commands are to enable the dynamic partitioning and bucketing
hive> set hive.exec.dynamic.partition.mode = nonstrict;
hive> set hive.exec.dynamic.partition = true;
hive> set hive.enforce.bucketing = true;

```
hive> set hive.exec.dynamic.partition.mode = nonstrict;
hive> set hive.exec.dynamic.partition = true;
hive> set hive.enforce.bucketing = true;
hive>
```

Creating an optimized table by applying partitioning on "event_type" and bucketing on "price"

CREATE TABLE IF NOT EXISTS dynpart_buck_retailsstore(event_time timestamp, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string)
PARTITIONED BY (event_type string)
CLUSTERED BY (price) INTO 10 BUCKETS
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
STORED AS TEXTFILE
LOCATION '/user/hive/hive-casestudy'
tblproperties('skip.header.line.count' = '1');

```
    >
    > CREATE TABLE IF NOT EXISTS dynpart_buck_retailsstore(event_time timestamp, product_id string, category_id string, category_code string, brand string,price float, user_id bigint, user_session string)
    > PARTITIONED BY (event_type string)
    > CLUSTERED BY (price) INTO 10 BUCKETS
    > ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
    > STORED AS TEXTFILE
    > LOCATION'/user/hive/hive-casestudy'
    > tblproperties('skip.header.line.count' = '1');
OK
Time taken: 0.109 seconds
```

Verifying the created table

```
hive> show tables;
OK
tab_name
dynpart_buck_retailsstore
retailsstore
Time taken: 0.374 seconds, Fetched: 2 row(s)
hive>
```

**INSERTING THE DATA INTO NEWLY CREATED OPTIMIZED TABLE (dynpart_buck_retailsstore) FROM EXISTING TABLE(retailsstore):**

INSERT INTO TABLE dynpart_buck_retailsstore
PARTITION (event_type)
SELECT event_time,
product_id, category_id, category_code, brand, price, user_id, user_session, event_type
FROM retailsstore;

```
hive> INSERT INTO TABLE dynpart_buck_retailsstore
    > PARTITION (event_type)
    > SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type
    > FROM retailsstore;
Query ID = hadoop_20221029131305_27ad2a73-5aac-4c79-b971-d7a4864957ee
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1667043937063_0004)

----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      2         2        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      5         5        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 02/02  [==========================>>] 100%  ELAPSED TIME: 170.96 s
----------------------------------------------------------------------------------------------
Loading data to table default.dynpart_buck_retailsstore partition (event_type=null)

Loaded : 4/4 partitions.
        Time taken to load dynamic partitions: 0.957 seconds
        Time taken for adding to write entity : 0.004 seconds
OK
event_time      product_id      category_id     category_code   brand   price   user_id user_session    event_type
Time taken: 185.845 seconds
hive>
```

Output: Based on the above results, it partitioned into 4
Verifying the partitioned in the Hadoop file system

```
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -ls /user/hive/warehouse/
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -ls /user/hive/hive-casestudy/
Found 6 items
-rw-r--r--   1 hadoop hadoop  545839412 2022-10-29 12:40 /user/hive/hive-casestudy/2019-Nov.csv
-rw-r--r--   1 hadoop hadoop  482542278 2022-10-29 12:35 /user/hive/hive-casestudy/2019-Oct.csv
drwxr-xr-x   - hadoop hadoop          0 2022-10-29 13:16 /user/hive/hive-casestudy/event_type=cart
drwxr-xr-x   - hadoop hadoop          0 2022-10-29 13:16 /user/hive/hive-casestudy/event_type=purchase
drwxr-xr-x   - hadoop hadoop          0 2022-10-29 13:16 /user/hive/hive-casestudy/event_type=remove_from_cart
drwxr-xr-x   - hadoop hadoop          0 2022-10-29 13:16 /user/hive/hive-casestudy/event_type=view
[hadoop@ip-172-31-20-212 ~]$
```

Randomly verifying the partitioned data in hadoop

```
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -ls /user/hive/hive-casestudy/event_type=cart
Found 10 items
-rwxr-xr-x   1 hadoop hadoop   27512579 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000000_0
-rwxr-xr-x   1 hadoop hadoop   32190447 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000001_0
-rwxr-xr-x   1 hadoop hadoop   33302805 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000002_0
-rwxr-xr-x   1 hadoop hadoop   32602023 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000003_0
-rwxr-xr-x   1 hadoop hadoop   34104132 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000004_0
-rwxr-xr-x   1 hadoop hadoop   32538513 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000005_0
-rwxr-xr-x   1 hadoop hadoop   39257340 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000006_0
-rwxr-xr-x   1 hadoop hadoop   24825787 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000007_0
-rwxr-xr-x   1 hadoop hadoop   28504487 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000008_0
-rwxr-xr-x   1 hadoop hadoop   35410315 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=cart/000009_0
```

```
[hadoop@ip-172-31-20-212 ~]$ hadoop fs -ls /user/hive/hive-casestudy/event_type=purchase
Found 10 items
-rwxr-xr-x   1 hadoop hadoop    6241877 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000000_0
-rwxr-xr-x   1 hadoop hadoop    7235640 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000001_0
-rwxr-xr-x   1 hadoop hadoop    7231471 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000002_0
-rwxr-xr-x   1 hadoop hadoop    7526313 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000003_0
-rwxr-xr-x   1 hadoop hadoop    7227979 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000004_0
-rwxr-xr-x   1 hadoop hadoop    7310389 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000005_0
-rwxr-xr-x   1 hadoop hadoop    8915123 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000006_0
-rwxr-xr-x   1 hadoop hadoop    5366094 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000007_0
-rwxr-xr-x   1 hadoop hadoop    6469070 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000008_0
-rwxr-xr-x   1 hadoop hadoop    8004214 2022-10-29 13:15 /user/hive/hive-casestudy/event_type=purchase/000009_0
[hadoop@ip-172-31-20-212 ~]$
```

## VERIFYING THE PERFORMANCE OF BOTH THE TABLES – BEFORE AND AFTER OPTIMIZED TECHNIQUES:

**select * from retailsstore limit 5;**



Time taken to retrieve first 5 rows of data **before optimization is 4.876 seconds (above)**

**select * from dynpart_buck_retailsstore limit 5;**



Time taken to retrieve the first 5 rows of data **after** optimization is **0.273 seconds (above screenshot)**

## ANSWERING GIVEN QUESTIONS:
### 1. Find the total revenue generated due to purchases made in October

**Base table:**
SELECT SUM(price) AS tot_revenue_oct FROM retailsstore WHERE MONTH(event_time) = '10' AND event_type = 'purchase';



Time taken is **134.492 seconds**

**Optimized table:**
SELECT SUM(price) AS tot_revenue_oct FROM dynpart_buck_retailsstore WHERE MONTH(event_time) = 10 AND event_type = 'purchase';

```
hive> SELECT SUM(price) AS tot_revenue_oct FROM dynpart_buck_retailsstore WHERE MONTH(event_time) = 10 AND event_type = 'purchase';
Query ID = hadoop_20221029134223_32d89022-1a79-43e2-9a54-cd84d111be42
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      3          3        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      1          1        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 02/02 [==========================>>] 100%  ELAPSED TIME: 24.18 s
----------------------------------------------------------------------------------------------
OK
tot_revenue_oct
1211532.4500002791
Time taken: 25.342 seconds, Fetched: 1 row(s)
hive>
```

Time taken with optimized table is **25.342 seconds**

## Insights:

1. The total revenue generated based on Purchase made in the month of October is 1,211,538.43 /-
2. Non-optimized table query took the execution time of 134.492 seconds whereas optimized table query took execution time of 25.342 seconds. We can see there is a significant drop in the execution time of the same query.
3. Hence, optimized table gives better performance in execution time.

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

**Base Query:**

SELECT MONTH(event_time) AS month, COUNT(event_type) AS sum_of_purchases FROM retailsstore WHERE event_type = 'purchase' GROUP BY MONTH(event_time);

```
hive> SELECT MONTH(event_time) AS month, COUNT(event_type) AS sum_of_purchases FROM retailsstore WHERE event_type = 'purchase' GROUP BY MONTH(event_time);
Query ID = hadoop_20221029134351_b57375c5-b507-42e8-808e-9a7b633f4599
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      5          5        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      3          3        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 02/02 [==========================>>] 100%  ELAPSED TIME: 105.71 s
----------------------------------------------------------------------------------------------
OK
month    sum_of_purchases
10       245624
11       322417
Time taken: 106.436 seconds, Fetched: 2 row(s)
```

Time taken is **106.436 seconds**

**Optimized table:**

SELECT MONTH(event_time) AS month, COUNT(event_type) AS sum_of_purchases FROM dynpart_buck_retailsstore WHERE event_type = 'purchase' GROUP BY MONTH(event_time);

```
hive> SELECT MONTH(event_time) AS month, COUNT(event_type) AS sum_of_purchases FROM dynpart_buck_retailsstore WHERE event_type = 'purchase' GROUP BY MONTH(event_time);
Query ID = hadoop_20221029134635_47690273-13e0-4048-8da4-08e50a2adef7
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------------------
Map 1 ......... container     SUCCEEDED      3         3        0        0       0       0
Reducer 2 ..... container     SUCCEEDED      1         1        0        0       0       0
--------------------------------------------------------------------------------------------
VERTICES: 02/02 [==========================>>] 100%  ELAPSED TIME: 30.18 s
--------------------------------------------------------------------------------------------
OK
month   sum_of_purchases
10      245619
11      322412
Time taken: 31.114 seconds, Fetched: 2 row(s)
hive>
```

Time taken is **31.114 seconds**

## Insights:

- Sum of purchases made in the month of October is 245624 and in the month of November 322417, which means number of purchases are increased in November month
- Non-optimized table query took the execution time of 106.436 seconds whereas optimized table query took execution time of 31.114 seconds. We can see there is a significant drop in the execution time of the same query.
- Hence, with proper partitioning and bucketing on table we can reduce execution time.

**Using Optimized table from below questions onwards:**

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

SELECT (SUM(CASE WHEN MONTH(event_time)=11 THEN price ELSE 0 END) - SUM(CASE WHEN MONTH(event_time)=10 THEN PRICE ELSE 0 END)) AS change_in_rev FROM dynpart_buck_retailsstore WHERE event_type = 'purchase' AND MONTH(event_time) in ('10','11');

```
hive> SELECT (SUM(CASE WHEN MONTH(event_time)=11 THEN price ELSE 0 END) - SUM(CASE WHEN MONTH(event_time)=10 THEN PRICE ELSE 0 END)) AS change_in_rev FROM dynpart_buck_retailsstore WHERE event_type = 'purchase'
AND MONTH(event_time) in ('10','11');
Query ID = hadoop_20221029134926_74d61db5-5491-4049-a53d-04070f4b06a7
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------------------
Map 1 ......... container     SUCCEEDED      3         3        0        0       0       0
Reducer 2 ..... container     SUCCEEDED      1         1        0        0       0       0
--------------------------------------------------------------------------------------------
VERTICES: 02/02 [==========================>>] 100%  ELAPSED TIME: 27.11 s
--------------------------------------------------------------------------------------------
OK
change_in_rev
319437.7899997565
Time taken: 27.765 seconds, Fetched: 1 row(s)
hive>
```

## Insights:

1. Time taken to execute the query is 27.765 seconds
2. Revenue increased in November by 319437.789 from October

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

```
SELECT DISTINCT SPLIT(category_code,'\\.')[0] AS Category
FROM dynpart_buck_retailsstore
WHERE category_code != '';
```

```
hive> SELECT DISTINCT SPLIT(category_code,'\\.')[0] AS Category FROM dynpart_buck_retailsstore WHERE category_code != '';
Query ID = hadoop_20221029135114_0f5f4ecd-af92-4016-8e1d-01de293bc5f5
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      6          6        0        0       0       0
Reducer 2 ..... container     SUCCEEDED      5          5        0        0       0       0
--------------------------------------------------------------------------------
VERTICES: 02/02  [========================>>] 100%  ELAPSED TIME: 68.03 s
--------------------------------------------------------------------------------
OK
category
furniture
appliances
accessories
apparel
sport
stationery
Time taken: 68.715 seconds, Fetched: 6 row(s)
hive>
```

## Insights:

1. Time taken to execute the query is 68.715 seconds
2. Total we got 6 distinct categories are – furniture, appliances, accessories, apparel, sport, stationery.

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

```
SELECT SPLIT(category_code,'\\.')[0] AS Category, COUNT(product_id) AS num_of_prod
FROM dynpart_buck_retailsstore
WHERE category_code != ''
GROUP BY SPLIT(category_code,'\\.')[0]
ORDER BY num_of_prod DESC;
```

```
hive> SELECT SPLIT(category_code,'\\.')[0] AS Category, COUNT(product_id) AS num_of_prod
    > FROM dynpart_buck_retailsstore
    > WHERE category_code != ''
    > GROUP BY SPLIT(category_code,'\\.')[0]
    > ORDER BY num_of_prod DESC;
Query ID = hadoop_20221029135418_4685ce12-95b8-4b84-8f67-341f5cb2971e
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      6          6        0        0       0       0
Reducer 2 ..... container     SUCCEEDED      5          5        0        0       0       0
Reducer 3 ..... container     SUCCEEDED      1          1        0        0       0       0
--------------------------------------------------------------------------------
VERTICES: 03/03  [========================>>] 100%  ELAPSED TIME: 68.55 s
--------------------------------------------------------------------------------
OK
category        num_of_prod
appliances      61736
stationery      26722
furniture       23604
apparel 18232
accessories     12928
sport   2
Time taken: 69.331 seconds, Fetched: 6 row(s)
hive>
```

## Insights:

1. Time taken to execute the query is 69.331 seconds
2. Appliances are having highest number of products available with 61736 compared to other categories.
3. Stationary and Furniture categories are almost equally registered with available ranges from 23000 to 27000.
4. Sports category is least available with 2 products.

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

WITH tot_sales AS(
SELECT brand, (SUM(CASE WHEN MONTH(event_time)=10 THEN price ELSE 0 END) +
SUM(CASE WHEN MONTH(event_time)=11 THEN PRICE ELSE 0 END)) AS total_sales
FROM dynpart_buck_retailsstore
WHERE event_type = 'purchase' AND MONTH(event_time) in
('10','11') AND brand != ''
GROUP BY brand)
SELECT brand, total_sales
FROM tot_sales
ORDER BY total_sales DESC
LIMIT 1;

```
hive> WITH tot_sales AS(
    > SELECT brand, (SUM(CASE WHEN MONTH(event_time)=10 THEN price ELSE 0 END) + SUM(CASE WHEN MONTH(event_time)=11 THEN PRICE ELSE 0 END)) AS total_sales
    > FROM
    > dynpart_buck_retailsstore
    > WHERE event_type = 'purchase' AND MONTH(event_time) in ('10','11') AND brand != ''
    > GROUP BY brand)
    > SELECT brand, total_sales
    > FROM tot_sales
    > ORDER BY total_sales DESC
    > LIMIT 1;
Query ID = hadoop_20221029135911_8caff0d4-cf8a-485e-9257-c70f5d67c27b
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------
Map 1 ......... container     SUCCEEDED      3         3        0        0       0       0
Reducer 2 ..... container     SUCCEEDED      1         1        0        0       0       0
Reducer 3 ..... container     SUCCEEDED      1         1        0        0       0       0
--------------------------------------------------------------------------------
VERTICES: 03/03  [=========================>>] 100%  ELAPSED TIME: 27.50 s
--------------------------------------------------------------------------------
OK
brand   total_sales
runail  148292.46000001638
Time taken: 28.184 seconds, Fetched: 1 row(s)
hive>
```

### Insights:
1. Runail is the brand that has the highest sales in total of both the months October and November.
2. It seems that Runail brand has high popularity among cosmetic lovers and bringing in more.
3. Products related to Runail brand could help in increasing their profit.

## 7. Which brands increased their sales from October to November?
WITH brand_sales AS(
SELECT brand, SUM(CASE WHEN MONTH(event_time)=10 THEN price ELSE 0 END) AS
Oct_sales, SUM(CASE WHEN MONTH(event_time)=11 THEN PRICE ELSE 0 END) AS
Nov_sales FROM dynpart_buck_retailsstore WHERE event_type = 'purchase' AND
MONTH(event_time) in ('10','11') AND brand != '' GROUP BY brand)
SELECT brand, Oct_sales, Nov_sales, Nov_sales-Oct_sales AS sale_diff
FROM brand_sales
WHERE Nov_sales-Oct_sales > 0
ORDER BY sale_diff DESC;

```
Time taken: 28.184 seconds, Fetched: 1 row(s)
hive> WITH brand_sales AS(
    > SELECT brand, SUM(CASE WHEN MONTH(event_time)=10 THEN price ELSE 0 END) AS Oct_sales, SUM(CASE WHEN MONTH(event_time)=11 THEN PRICE ELSE 0 END) AS Nov_sales FROM dynpart_buck_retailsstore WHERE event_type
= 'purchase' AND MONTH(event_time) in ('10','11') AND brand != '' GROUP BY brand)
    > SELECT brand, Oct_sales, Nov_sales, Nov_sales-Oct_sales AS sale_diff
    > FROM brand_sales
    > WHERE Nov_sales-Oct_sales > 0
    > ORDER BY sale_diff DESC;
Query ID = hadoop_20221029140234_047efa69-bdc5-4d36-a2c4-ec36cb16faf3
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

--------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS   TOTAL   COMPLETED   RUNNING   PENDING   FAILED   KILLED
--------------------------------------------------------------------------------------------
Map 1 ......... container   SUCCEEDED     3        3          0         0         0        0
Reducer 2 ..... container   SUCCEEDED     1        1          0         0         0        0
Reducer 3 ..... container   SUCCEEDED     1        1          0         0         0        0
--------------------------------------------------------------------------------------------
VERTICES: 03/03 [==========================>>] 100%  ELAPSED TIME: 27.49 s
--------------------------------------------------------------------------------------------
OK
brand    oct_sales      nov_sales       sale_diff
grattol 35445.54000000205    71472.7100000044       36027.1700000042
uno     35302.03000000019    51039.75000000047      15737.720000000285
lianail 5892.840000000024    16394.240000000544     10501.400000000052
ingarden        23161.3900000002      33566.210000000625      10404.820000000425
strong 29196.63000000005    38671.269999999975      9474.63999999997
jessnail        26287.840000000084    33345.22999999992       7057.389999999839
cosmoprofi      8322.809999999901     14536.989999999909      6214.180000000008
polarus 6013.719999999998    11371.930000000004     5358.210000000055
runail 71537.77000001163    76754.69000000473       5216.919999993101
freedecor       3421.7800000000097    7671.800000000062       4250.020000000052
staleks 8519.730000000014    11875.610000000015     3355.880000000001
bpw.style       11572.15000000046     14837.440000000377      3265.289999999917
lovely 8704.379999999985    11939.059999999967     3234.679999999982
marathon        7280.749999999996     10273.099999999999      2992.350000000002
haruyama        9390.690000000088     12352.910000000145      2962.2200000000576
yoko   8756.909999999983    11707.879999999976     2950.969999999992
italwax 21940.23999999968    24799.369999999864     2859.130000000183
benovy 409.6199999999999    3259.970000000002      2850.350000000002
kaypro 881.3399999999999    3268.6999999999985     2387.3599999999988
estel  21756.75000000025    24142.670000000027     2385.920000000002
concept 11032.13999999978    13380.399999999978     2348.26
kapous 11927.15999999738    14093.079999999864     2165.9200000001256
f.o.x  6624.230000000007    8577.279999999979      1953.049999999972
masura 31266.07999999318    33058.46999999992      1792.3900000006033
milv   3904.940000000026    5642.0100000001285     1737.0700000001025
beautix 10493.94999999997    12222.95       1729.000000000031
artex  2730.6400000000002   4327.249999999996      1596.6099999999942
domix  10472.04999999992    12009.169999999851     1537.1199999999317
shik   3341.2000000000035   4839.720000000001      1498.5199999999977
```

```
OK
brand    oct_sales      nov_sales       sale_diff
grattol 35445.54000000205    71472.7100000044       36027.1700000042
uno     35302.03000000019    51039.75000000047      15737.720000000285
lianail 5892.840000000024    16394.240000000544     10501.400000000052
ingarden        23161.3900000002      33566.210000000625      10404.820000000425
strong 29196.63000000005    38671.269999999975      9474.63999999997
jessnail        26287.840000000084    33345.22999999992       7057.389999999839
cosmoprofi      8322.809999999901     14536.989999999909      6214.180000000008
polarus 6013.719999999998    11371.930000000004     5358.210000000055
runail 71537.77000001163    76754.69000000473       5216.919999993101
freedecor       3421.7800000000097    7671.800000000062       4250.020000000052
staleks 8519.730000000014    11875.610000000015     3355.880000000001
bpw.style       11572.15000000046     14837.440000000377      3265.289999999917
lovely 8704.379999999985    11939.059999999967     3234.679999999982
marathon        7280.749999999996     10273.099999999999      2992.350000000002
haruyama        9390.690000000088     12352.910000000145      2962.2200000000576
yoko   8756.909999999983    11707.879999999976     2950.969999999992
italwax 21940.23999999968    24799.369999999864     2859.130000000183
benovy 409.6199999999999    3259.970000000002      2850.350000000002
kaypro 881.3399999999999    3268.6999999999985     2387.3599999999988
estel  21756.75000000025    24142.670000000027     2385.920000000002
concept 11032.13999999978    13380.399999999978     2348.26
kapous 11927.15999999738    14093.079999999864     2165.9200000001256
f.o.x  6624.230000000007    8577.279999999979      1953.049999999972
masura 31266.07999999318    33058.46999999992      1792.3900000006033
milv   3904.940000000026    5642.0100000001285     1737.0700000001025
beautix 10493.94999999997    12222.95       1729.000000000031
artex  2730.6400000000002   4327.249999999996      1596.6099999999942
domix  10472.04999999992    12009.169999999851     1537.1199999999317
shik   3341.2000000000035   4839.720000000001      1498.5199999999977
smart  4457.2599999999875   5902.139999999909      1444.88000000001
roubloff        3491.3600000000006    4913.770000000003       1422.4100000000026
levrana 2243.560000000002    3664.100000000002      1420.5400000000018
oniq   8425.409999999947    9841.64999999998       1416.2400000000325
irisk  45591.95999999969    46946.03999999916      1354.0799999994633
severina        4775.879999999955     6120.479999999956       1344.6000000000013
joico  705.52    2015.1000000000006      1309.580000000006
zeitun 708.6600000000067    2009.63        1300.9699999999998
beauty-free     554.1700000000014     1782.8599999999983      1228.6899999999969
swarovski       1887.9299999999978    3043.1599999999794      1155.2299999999896
de.lux 1659.6999999999784    2775.509999999973      1115.8099999999945
metzger 5373.45000000001     6457.160000000005       1083.7099999999955
markell 1768.7499999999998    2834.429999999994       1065.6799999999942
sanoto 157.14    1209.6799999999993      1052.54
nagaraku        4369.740000000042     5327.6800000000285      957.9399999999869
ecolab 262.85    1214.3000000000009      951.4500000000008
art-visage      2092.709999999978     2997.800000000056       905.0900000000079
levissime       2227.5000000000064    3085.309999999854       857.809999999979
missha 1293.83 2150.2799999999997      856.4499999999998
solomeya        1899.7000000000012    2685.799999999996       786.0999999999949
rosi   3077.040000000002    3841.560000000001      764.5199999999991
```

```
rosi   3077.040000000002    3841.560000000001      764.5199999999991
refectocil      2716.180000000003     3475.580000000036       759.4000000000005
kaaral 4412.429999999999    5086.07 673.6400000000003
kosmekka        1181.4399999999993    1813.37 631.9300000000003
kinetics        6334.250000000006     6945.259999999998       611.009999999992
browxenna       14331.370000000057    14916.73000000007       585.3600000000133
airnails        5118.900000000015     5691.520000000021       572.6200000000063
uskusi 5142.269999999997    5690.310000000087      548.0400000000509
coifin 903.0000000000002    1428.4899999999998      525.4899999999996
s.care 412.68   913.07 500.39000000000004
limoni 1308.9000000000005    1796.600000000004      487.6999999999998
matrix 3243.250000000001    3726.740000000016      483.4900000000007
gehwol 1089.0700000000002    1557.68 468.6099999999999
greymy 29.21   489.49 460.28000000000003
bioaqua 942.8900000000004    1398.1200000000001      455.23
farmavita       837.37  1291.97 454.6
sophin 1067.8600000000003    1515.520000000002      447.66000000000145
yu-r   271.41   673.7099999999999      402.2999999999999
kiss   421.55   817.3300000000001      395.78000000000037
naomi  0.0      389.0   389.0
lador  2083.610000000002    2471.5300000000016     387.9199999999996
ellips 245.84999999999997    606.04  360.19
jas    3318.960000000002    3657.4300000000026     338.4700000000007
lowence 242.83999999999997    567.7499999999999       324.9099999999999
nitrile 847.279999999999     1162.679999999999       315.4
shary  871.9599999999997    1176.4899999999996     304.5300000000002
kims   330.04   632.0400000000001      302.00000000000006
happyfons       801.9200000000004     1091.5900000000008      289.6700000000004
kocostar        310.8499999999999     594.9299999999998       284.0799999999999
insight 1443.7000000000005    1721.9600000000001     278.2600000000003
candy  534.9599999999999    799.3799999999994      264.41999999999995
bluesky 10307.240000000156    10565.529999999784      258.2899999999628
beauugreen      511.51000000000016    768.3499999999999       256.8399999999975
protokeratin    201.25  456.79  255.54000000000002
trind  298.07000000000005    542.96  244.89
entity 479.71000000000157    719.2599999999991      239.5499999999975
skinlite        651.9399999999997     890.4499999999998       238.5100000000001
provoc 827.9900000000004    1063.8200000000024     235.8300000000021
fedua  52.38    263.81  211.43
ecocraft        41.16000000000004     241.9499999999996       200.78999999999996
keen   236.35   435.6199999999995      199.26999999999995
mane   66.78999999999999     260.26  193.47
freshbubble     318.6999999999999     502.33999999999975      183.63999999999987
matreshka       0.0      182.67000000000004      182.67000000000004
chi    358.9400000000001    538.6100000000001      179.67000000000002
cristalinas     427.6299999999999     584.9499999999999       157.32000000000005
farmona 1692.4600000000003    1843.4299999999998      150.96999999999957
latinoil        249.52  384.59000000000015      135.07000000000014
miskin 158.04000000000002    293.0699999999999      135.0299999999994
elizavecca      70.53   204.30000000000004      133.77000000000004
nefertiti       233.52000000000007    366.64  133.11999999999992
finish 98.38    230.38  132.0
igrobeauty      513.6600000000005     645.0700000000006       131.41000000000008
```

```
igrobeauty     513.6600000000005          645.0700000000006          131.41000000000008
dizao   819.1300000000003          945.5100000000014          126.38000000000102
osmo    645.58  762.31  116.7299999999999
batiste 772.4000000000001          874.1699999999998          101.76999999999975
carmex  145.08000000000004          243.36  98.27999999999997
eos     54.339999999999996          152.61  98.27000000000001
depilflax       2707.0699999999956          2803.779999999998          96.71000000000231
enjoy   41.35   136.57000000000002          95.22000000000003
kerasys 430.91000000000014          525.2   94.2899999999999
aura    83.95   177.58999999999996          93.55999999999996
plazan  101.36999999999999          194.01000000000005          92.64000000000006
koelf   422.72999999999996          507.28999999999985          84.55999999999989
nirvel  163.04  234.32999999999987          71.28999999999988
konad   739.8299999999997          810.6699999999992          70.83999999999946
egomania        77.47   146.04000000000002          68.57000000000002
cutrin  299.37  367.62  68.25
laboratorium    246.5   312.52  66.01999999999998
inm     288.0199999999999          351.2100000000001          63.19000000000017
dewal   0.0     61.28999999999999          61.28999999999999
marutaka-foot   49.22   109.33000000000001          60.110000000000014
kares   0.0     59.45   59.45
profhenna       679.2300000000002          736.8499999999999          57.61999999999966
koelcia 55.5    112.75  57.25
balbcare        155.32999999999996          212.37999999999997          57.05000000000001
elskin  251.0900000000001          307.65000000000015          56.56000000000006
foamie  35.04   80.49   45.44999999999996
ladykin 125.64999999999999          170.57  44.92
likato  296.05999999999983          340.96999999999997          44.91000000000014
mavala  409.0400000000001          446.3200000000001          37.28000000000003
vilenta 197.59999999999997          231.21000000000004          33.61000000000007
beautyblender   78.74000000000001          109.40999999999998          30.669999999999973
biore   60.65000000000006          90.31   29.659999999999997
orly    902.3800000000002          931.0900000000004          28.71000000000015
estelare        444.81000000000005          471.87000000000023          27.059999999999718
profepil        93.36000000000001          118.02000000000001          24.659999999999997
blixz   38.95   63.4    24.449999999999996
binacil 0.0     24.259999999999998          24.259999999999998
godefroy        401.22  425.12  23.899999999999977
glysolid        69.72999999999998          91.58999999999999          21.86000000000014
veraclara       50.11000000000001          71.21000000000001          21.1
juno    0.0     21.08   21.08
kamill  63.010000000000005          81.49000000000001          18.480000000000004
treaclemoon     163.37000000000003          181.49000000000004          18.120000000000005
supertan        50.37000000000001          66.51000000000002          16.140000000000008
barbie  0.0     12.39   12.39
deoproce        316.84000000000003          329.1700000000001          12.330000000000041
rasyan  18.799999999999997          28.93999999999998          10.14
fly     17.14   27.16999999999998          10.029999999999998
tertio  236.16  245.8   9.640000000000015
jaguar  1102.1100000000004          1110.6500000000003          8.539999999999964
soleo   204.2   212.5299999999998          8.329999999999814
neoleor 43.41   51.7    8.290000000000006
moyou   5.71    10.280000000000001          4.570000000000001
```

```
kerasys 430.91000000000014          525.2   94.2899999999999
aura    83.95   177.58999999999996          93.55999999999996
plazan  101.36999999999999          194.01000000000005          92.64000000000006
koelf   422.72999999999996          507.28999999999985          84.55999999999989
nirvel  163.04  234.32999999999987          71.28999999999988
konad   739.8299999999997          810.6699999999992          70.83999999999946
egomania        77.47   146.04000000000002          68.57000000000002
cutrin  299.37  367.62  68.25
laboratorium    246.5   312.52  66.01999999999998
inm     288.0199999999999          351.2100000000001          63.19000000000017
dewal   0.0     61.28999999999999          61.28999999999999
marutaka-foot   49.22   109.33000000000001          60.110000000000014
kares   0.0     59.45   59.45
profhenna       679.2300000000002          736.8499999999999          57.61999999999966
koelcia 55.5    112.75  57.25
balbcare        155.32999999999996          212.37999999999997          57.05000000000001
elskin  251.0900000000001          307.65000000000015          56.56000000000006
foamie  35.04   80.49   45.44999999999996
ladykin 125.64999999999999          170.57  44.92
likato  296.05999999999983          340.96999999999997          44.91000000000014
mavala  409.0400000000001          446.3200000000001          37.28000000000003
vilenta 197.59999999999997          231.21000000000004          33.61000000000007
beautyblender   78.74000000000001          109.40999999999998          30.669999999999973
biore   60.65000000000006          90.31   29.659999999999997
orly    902.3800000000002          931.0900000000004          28.71000000000015
estelare        444.81000000000005          471.87000000000023          27.059999999999718
profepil        93.36000000000001          118.02000000000001          24.659999999999997
blixz   38.95   63.4    24.449999999999996
binacil 0.0     24.259999999999998          24.259999999999998
godefroy        401.22  425.12  23.899999999999977
glysolid        69.72999999999998          91.58999999999999          21.86000000000014
veraclara       50.11000000000001          71.21000000000001          21.1
juno    0.0     21.08   21.08
kamill  63.010000000000005          81.49000000000001          18.480000000000004
treaclemoon     163.37000000000003          181.49000000000004          18.120000000000005
supertan        50.37000000000001          66.51000000000002          16.140000000000008
barbie  0.0     12.39   12.39
deoproce        316.84000000000003          329.1700000000001          12.330000000000041
rasyan  18.799999999999997          28.93999999999998          10.14
fly     17.14   27.16999999999998          10.029999999999998
tertio  236.16  245.8   9.640000000000015
jaguar  1102.1100000000004          1110.6500000000003          8.539999999999964
soleo   204.2   212.5299999999998          8.329999999999814
neoleor 43.41   51.7    8.290000000000006
moyou   5.71    10.280000000000001          4.570000000000001
bodyton 1376.3399999999983          1380.639999999999          4.300000000000637
skinity 8.88    12.44000000000001          3.560000000000005
helloganic      0.0     3.1     3.1
grace   100.92000000000002          102.60999999999999          1.6899999999999693
cosima  20.23   20.929999999999993          0.6999999999999922
ovale   2.54    3.1     0.56
Time taken: 28.051 seconds, Fetched: 160 row(s)
hive>
```

## Insights:

1. Here are some 160 brands with increment in the selling from October to November.
2. 'Grattol' brand has the highest total increment i.e., 36,027 /- and 'Ovale' seems to have the least increment of 0.56 /- from October to November.
3. Among all these brands lists, 'Runail' which was the best brand in terms of selling in October and November combined is also in the top 10 brands with high increment for October (71539.28) to November (76758.61) i.e., increment of total 5219.38.

4. This implies that 'Runail' is the best and popular brand among all other brands within people.

**8. Your company wants to reward the top 10 users of its websites with a golden customer plan. Write a query to generate a list of top 10 users who spend the most.**
SELECT user_id, SUM(price) AS tot_amt_spend FROM dynpart_buck_retailsstore
WHERE event_type = 'purchase'
GROUP BY user_id
ORDER BY tot_amt_spend DESC
LIMIT 10;

```
hive> SELECT user_id, SUM(price) AS tot_amt_spend FROM dynpart_buck_retailsstore
    > WHERE event_type = 'purchase'
    > GROUP BY user_id
    > ORDER BY tot_amt_spend DESC
    > LIMIT 10;
Query ID = hadoop_20221029140711_21c88b43-6ef1-40ce-92ce-af782a809053
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667043937063_0005)

----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      3         3        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      1         1        0        0       0       0
Reducer 3 ...... container     SUCCEEDED      1         1        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 03/03  [=========================>>] 100%  ELAPSED TIME: 26.33 s
----------------------------------------------------------------------------------------------
OK
user_id tot_amt_spend
557790271       2715.869999999995
150318419       1645.9699999999998
562167663       1352.8500000000001
531900924       1329.4499999999996
557850743       1295.4800000000007
522130011       1185.3899999999999
561592095       1109.7000000000003
431950134       1097.5899999999997
566576008       1056.3599999999997
521347209       1040.9100000000003
Time taken: 27.074 seconds, Fetched: 10 row(s)
hive>
```

**Insights:**
1. Here is the list of the top 10 users or buyers who have spent the most and could be rewarded with a Golden Customer plan to attract more people in the coming future.
2. With the Optimized table the execution time reduced with proper partitioning and bucketing.
3. Time taken to execute this query on optimized table is 27.874 seconds.
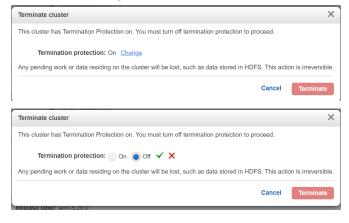
Cleaning up:
Once the analysis is completed, we should drop the tables and databases
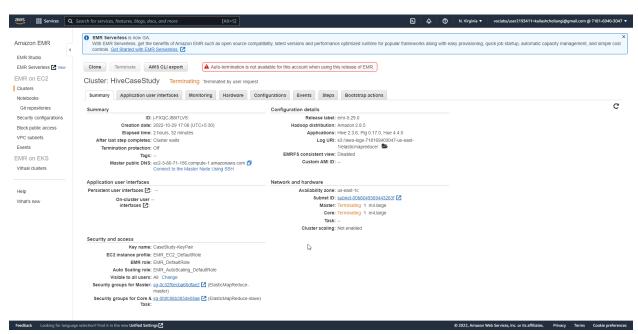
```
hive> show tables;
OK
tab_name
retailsstore
dynpart_buck_retailsstore
Time taken: 0.041 seconds, Fetched: 2 row(s)
hive> drop table retailsstore;
OK
Time taken: 0.113 seconds
hive> drop table dynpart_buck_retailsstore;
ok
Time taken: 0.326 seconds
```

```
hive> drop database casestudy;
OK
Time taken: 0.184 seconds
```
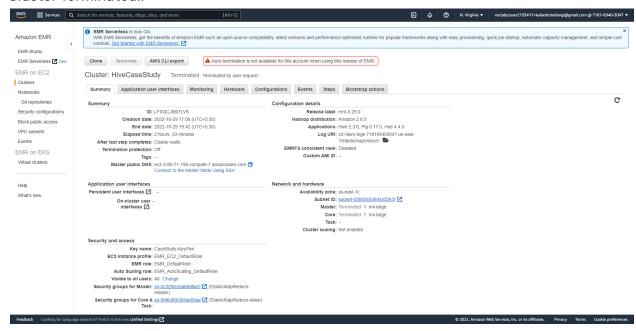
## TERMINATION PROCESS:

After completing our analysis, we should terminate the EMR cluster

**Terminate cluster**　　　　　　　　　　　　　　　　　　　　×

This cluster has Termination Protection on. You must turn off termination protection to proceed.

**Termination protection:** On  Change

Any pending work or data residing on the cluster will be lost, such as data stored in HDFS. This action is irreversible.

Cancel  **Terminate**

---

**Terminate cluster**　　　　　　　　　　　　　　　　　　　　×

This cluster has Termination Protection on. You must turn off termination protection to proceed.

**Termination protection:**  ○ On  ● Off  ✓ ✗

Any pending work or data residing on the cluster will be lost, such as data stored in HDFS. This action is irreversible.

Cancel  **Terminate**

Release label: emr-5.29.0

# Cluster Terminated!!



# Thankyou!!