HIVE CASE STUDY

(DSC30)- Sahil Thakare

PROBLEM STATEMENT:

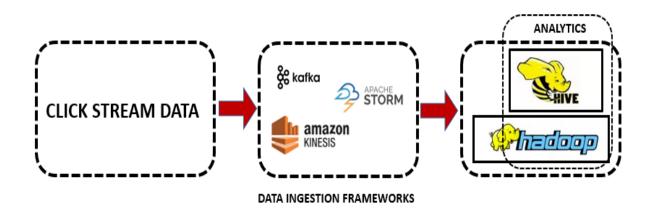
With online sales gaining popularity, tech companies are exploring ways to improve their sales by analyzing customer behavior and gaining insights about product trends. Furthermore, the websites make it easier for customers to find the products they require without much scavenging. In this case study, we are working with clickstream data by getting insights and making decisions upon how the E-commerce websites canimprove their sales.

OBJECTIVE:

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

DATA:

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



You will find the data in the link given below.

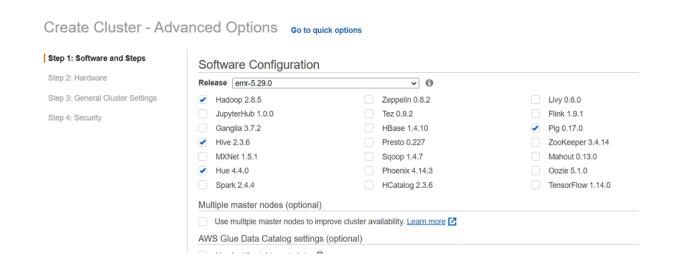
https://e-commerce-events-ml.s3.amazonaws.com/2019-Oct.csv https://e-commerce-events-ml.s3.amazonaws.com/2019-Nov.csv

The implementation phase:

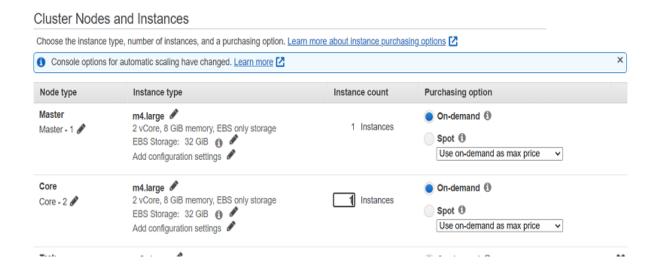
- Copying the data set into the HDFS:
 - Launch an EMR cluster that utilizes the Hive services
 - Move the data from the S3 bucket into the HDFS
- Creating the database and launching Hive queries on your EMR cluster:
 - Create the structure of your database,
 - Use optimized techniques to run your queries as efficiently as possible
 - Show the improvement of the performance after using optimization on any single query.
 - Run Hive queries to answer the questions given below.
- Cleaning up
 - Drop your database, and
 - Terminate your cluster

*** EMR CLUSTER CREATION**

EMR Cluster Landing page > Create Cluster > Advanced Options > Selecting the release emr-5.29 and the required services.



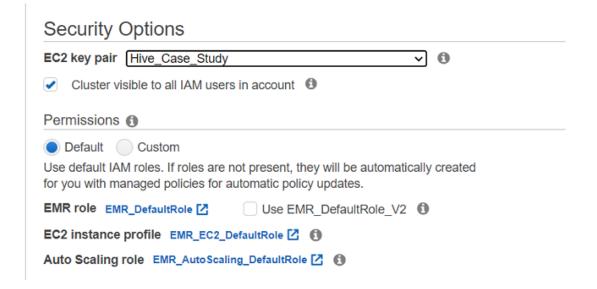
Hardware Configuration Page > To define the cluster & nodes: Instance type for both master & core nodes are M4. Large



Naming the Cluster:



Selecting the Key-pair (Created before creating the cluster):



Cluster "Hive Case Study 01" is successfully created and launched.

Using the "Hive_Case_Study" key-pair we enter the terminal.

***** HADOOP & HIVE QUERIES:

Terminal > Connecting to EMR Cluster(Putty).

```
Using username
Authenticating with public key "imported-openssh-key"
Last login: Tue Nov 30 12:50:56 2021
                    Amazon Linux AMI
nttps://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
68 package(s) needed for security, out of 106 available Run "sudo yum update" to apply all updates.
EEEEEEEEEEEEEEEEE MMMMMMMM
                                      M::::::: M R:::::::::::::R
EE:::::EEEEEEEEE:::E M:::::::M
                                    M::::::M R:::::RRRRRR:::::R
          EEEEE M::::::M
                                   M:::::::: M RR::::R
 E::::E
                   M::::::::M
                                                 R:::R
                                                           R::::R
                   M:::::M M:::M M::::M
 E::::EEEEEEEEE
                                                 R:::RRRRRR::::R
 E::::EEEEEEEEE
                                                 R:::RRRRRR::::R
 E::::E
             EEEEE M:::::M
                               MMM
                                       M:::::M
                                                 R:::R
                                                           R::::R
E:::::EEEEEEEE::::E M:::::M
                                       M:::::M
                                                 R:::R
                                                           R::::R
M:::::M RR::::R
                                                           R::::R
EEEEEEEEEEEEEEEEE MMMMMMM
                                       MMMMMMM RRRRRRR
                                                           RRRRRR
[hadoop@ip-10-0-0-13 ~]$
```

Creating a directory "casestudy"

hadoop fs -mkdir /casestudy

hadoop fs -ls /

```
[hadoop@ip-10-0-0-61 ~]$ hadoop fs -mkdir /casestudy
[hadoop@ip-10-0-0-61 \sim]$ hadoop fs -ls /
Found 5 items
drwxr-xr-x
                                       0 2021-12-07 11:27 /apps
            - hdfs
                      hadoop
                                       0 2021-12-07 11:36 /casestudy
drwxr-xr-x

    hadoop hadoop

                                       0 2021-12-07 11:29 /tmp
drwxrwxrwt
             - hdfs
                      hadoop
                                       0 2021-12-07 11:27 /user
                      hadoop
drwxr-xr-x
             - hdfs
                                       0 2021-12-07 11:27 /var
drwxr-xr-x
             - hdfs
                      hadoop
```

Loading the datasets into HDFS from S3:

hadoopdistcp 's3://e-commerce-events-ml/2019-Oct.csv' /hive02/2019_Oct.csv

```
[hadoop@ip-10-0-0-61 ~]$ hadoop distcp 's3://e-commerce-events-ml/2019-Oct.csv' /hive02/2019-Oct.cs/
21/12/07 11:41:07 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, skipCRC=false, blocking=true, numListstatusThreads=0, maxMaps=20, mapBandwidth=100, sslConfigurationFile='null', copyStrategy='unifor nsize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://e-commerce-events-ml/2019-Oct.csv], targetPath=/hive02/2019-Oct.csv, targetPathExists=false, filtersFile='null'}
21/12/07 11:41:07 INFO client.RMProxy: Connecting to ResourceManager at ip-10-0-0-61.ec2.internal/10.0.0.61:8032
21/12/07 11:41:12 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
21/12/07 11:41:12 INFO tools.SimpleCopyListing: Build file listing completed.
```

```
21/12/07 11:41:45 INFO mapreduce.Job: Job job_1638876533904_0001 completed successfully
21/12/07 11:41:45 INFO mapreduce.Job: Counters: 38

File System Counters

FILE: Number of bytes read=0

FILE: Number of bytes written=172396

FILE: Number of read operations=0

FILE: Number of large read operations=0

FILE: Number of write operations=0

HDFS: Number of bytes written=482542278

HDFS: Number of read operations=12

HDFS: Number of large read operations=0

HDFS: Number of large read operations=0

HDFS: Number of read operations=12

HDFS: Number of write operations=4
```

hadoopdistcp 's3://e-commerce-events-ml/2019-Oct.csv' /hive02/2019_Nov.csv

```
[higoop@ip-10-0-0-61 ~]$ hadoop distcp 's3://e-commerce-events-ml/2019-Nov.csv' /hive02/2019-Nov.cs / 21/12/07 11:42:44 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, skipCRC=false, blocking=true, numListstatusThreads=0, maxMaps=20, mapBandwidth=100, sslConfigurationFile='null', copyStrategy='unifor msize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://e-commerce-events-ml/2019-Nov.csv], targetPath=/hive02/2019-Nov.csv, targetPathExists=false, filtersFile='null'} 21/12/07 11:42:44 INFO client.RMProxy: Connecting to ResourceManager at ip-10-0-0-61.ec2.internal/1 0.0.0.61:8032 21/12/07 11:42:48 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0 21/12/07 11:42:48 INFO tools.SimpleCopyListing: Build file listing completed. 21/12/07 11:42:48 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapred uce.task.io.sort.mb 21/12/07 11:42:48 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapred uce.task.io.sort.factor 11:42:48 INFO tools.DistCp: Number of paths in the copy list: 1 21/12/07 11:42:48 INFO tools.DistCp: Number of paths in the copy list: 1 21/12/07 11:42:48 INFO client.RMProxy: Connecting to ResourceManager at ip-10-0-0-61.ec2.internal/1 21/12/07 11:42:48 INFO client.RMProxy: Connecting to ResourceManager at ip-10-0-0-61.ec2.internal/1
```

Viewing the data

hadoop fs -cat /hive02/2019-Oct.csv | head

```
[hadoop@ip-10-0-0-61 \sim]$ hadoop fs -cat /hive02/2019-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-
3d2c-92e149dab885
2019-10-01 00:00:03 UTC,cart,5773353,1487580005134238553,,runai1,2.62,463240011,26dd6e6e-4dac-4778-
d2c-92e149dab885
019-10-01 00:00:07 UTC,cart,5881589,2151191071051219817,,lovely,13.48,429681830,49e8d843-adf3-428b
a2c3-fe8bc6a307c9
019-10-01 00:00:07 UTC,cart,5723490,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-
d2c-92e149dab885
019-10-01 00:00:15 UTC,cart,5881449,1487580013522845895,,lovely,0.56,429681830,49e8d843-adf3-428b-
2c3-fe8bc6a307c9
2019-10-01 00:00:16 UTC,cart,5857269,1487580005134238553,,runail,2.62,430174032,73dea1e7-664e-43f4-
8b30-d32b9d5af04f
019-10-01 00:00:19 UTC,cart,5739055,1487580008246412266,,kapous,4.75,377667011,81326ac6-daa4-4f0a-
b488-fd0956a78733
019-10-01 00:00:24 UTC,cart,5825598,1487580009445982239,,,0.56,467916806,2f5b5546-b8cb-9ee7-7ecd-8
```

hadoop fs -cat /hive02/2019-Nov.csv | head

```
[hadoop@ip-10-0-0-61 ~]$ hadoop fs -cat /hive02/2019-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-3
3527f4de241
2019-11-01 00:00:09 UTC,cart,5844397,1487580006317032337,,2.38,553329724,2067216c-31b5-455d-a1cc-a
f0575a34ffb
2019-11-01 00:00:10 UTC,view,5837166,1783999064103190764,pnb,22.22,556138645,57ed222e-a54a-4907-99
44-5a875c2d7f4f
2019-11-01 00:00:11 UTC,cart,5876812,1487580010100293687,jessnail,3.16,564506666,186c1951-8052-4b3
7-adce-dd9644b1d5f7
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,3.33,553329724,2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,3.33,553329724,2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,3.33,553329724,2067216c-31b5-455d-a1cc-af0575a34ffb
```

Data are successfully loaded.

Launch Hive

Creating new database

hive> CREATE DATABASE IF NOT EXISTS hive_assignment;

hive> SHOW DATABASES;

hive> DESCRIBE DATABASE hive_assignment;

```
hive> CREATE DATABASE IF NOT EXISTS hive_assignment;
OK
Time taken: 0.319 seconds
hive> SHOW DATABASES;
OK
default
hive_assignment
Time taken: 0.021 seconds, Fetched: 2 row(s)
hive> DESCRIBE DATABASE hive_assignment;
OK
```

Creating retail

hive >CREATE EXTERNAL TABLE IF NOT EXISTS retail (event_time timestamp, event_type string,product_id string, category_id string, category_code string, brand string, price decimal(10,3user_session string) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' WITHSERDEPROPERTIES ("separatorChar" = "," , "quoteChar" = "\"", "escapeChar" = "\\") stored as tLOCATION '/casestudy' TBLPROPERTIES ("skip.header.line.count"="1");

hive > DESCRIBE retail;

Loading data into table "retail";

hive> LOAD DATA INPATH '/hive02/2019-Oct.csv' INTO TABLE retail;

hive> LOAD DATA INPATH '/hive02/2019-Nov.csv' INTO TABLE retail;

```
hive> LOAD DATA INPATH '/hive02/2019-Oct.csv' INTO TABLE retail;
Loading data to table default.retail
OK
Time taken: 1.116 seconds
hive> LOAD DATA INPATH '/hive02/2019-Nov.csv' INTO TABLE retail;
Loading data to table default.retail
OK
```

Performing data check:

```
hive> SELECT * FROM retail WHERE MONTH(event_time)=11 limit 5;
```

hive> SELECT * FROM retail WHERE MONTH(event_time)=10 limit 5;

```
nive> SELECT * FROM retail WHERE MONTH(event_time)=11 limit 5
                                 5802432 1487580009286598681
                                                                                    0.32
                                                                                             562076640 0
9fafd6c-6c99-46b1-834f-33527f4de241
                                5844397 1487580006317032337
2019-11-01 00:00:09 UTC cart
                                                                                    2.38
                                                                                             553329724 2
067216c-31b5-455d-a1cc-af0575a34ffb
                                 5837166 1783999064103190764
                                                                                             556138645 5
                                                                                    22.22
                                                                            pnb
ed222e-a54a-4907-9944-5a875c2d7f4f
2019-11-01 00:00:11 UTC cart
                                5876812 1487580010100293687
506666 186c1951-8052-4b37-adce-dd9644b1d5f7
2019-11-01 00:00:24 UTC remove_from_cart
                                                  5826182 1487580007483048900
                        2067216c-31b5-455d-a1cc-af0575a34ffb
       553329724
nive> SELECT * FROM retail WHERE MONTH(event time)=10 limit 5;
2019-10-01 00:00:00 UTC cart 5773
6dd6e6e-4dac-4778-8d2c-92e149dab885
                                 5773203 1487580005134238553
                                                                            runail 2.62
                                                                                             463240011 2
                                5773353 1487580005134238553
                                                                                             463240011 2
5dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:07 UTC cart 5881589 2151191071051219817
                                                                            lovely 13.48
                                                                                            429681830 4
e8d843-adf3-428b-a2c3-fe8bc6a307c9
2019-10-01 00:00:07 UTC cart 572:
6dd6e6e-4dac-4778-8d2c-92e149dab885
                                5723490 1487580005134238553
                                                                            runail 2.62
                                                                                             463240011 2
```

DYNAMIC PARTITIONING

hive> set hive.exec.dynamic.partition=true;

hive> set hive.exec.dynamic.partition.mode=nonstrict;

PARTITION TABLE 1: retail_part_1

hive> CREATE EXTERNAL TABLE IF NOT EXISTS retail_part_1 (event_time timestamp, product_id string, category_id string, category_code string, brand string, price decimal(10,3), user_idbigint, user_sessionstring) PARTITIONED BY(event_type string) CLUSTERED BY (user_id) INTO 5 buckets ROW FORMATSERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS textfile;

hive> DESCRIBE retail_part_1;

```
hive> DESCRIBE retail part 1 ;
OK
event_time
                       string
                                                from deserializer
product id
                       string
                                                from deserializer
category id
                                                from deserializer
                       string
category code
                                                from deserializer
                       string
brand
                       string
                                                from deserializer
price
                                                from deserializer
                       string
user_id
                                                from deserializer
                       string
                                                from deserializer
user_session
                       string
event_type
                        string
```

hive> INSERT INTO TABLE retail_part_1 PARTITION (event_type) SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type FROM retail;

hive> SELECT SUM(price) FROM retail_part_1 WHERE MONTH(event_time)=10 AND event_type='purchase';

Time Taken to execute the below query is 25.038 sec.

PARTITION TABLE 2: retail_part_3

Partition on: month

hive> CREATE EXTERNAL TABLE IF NOT EXISTS retail_part_3 (event_time timestamp, event_type string, product_id string, category_id string, category_code string, brand string, price decimal(10,3), user_idbigint, user_session string) PARTITIONED BY(month int) CLUSTERED BY (brand) INTO 5 buckets ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS textfile; hive> DESCRIBE retail_part_3

hive > DESCRIBE retail_part_3;

```
ive> SELECT SUM(price) FROM retail_part_1 WHERE MONTH(event_time)=10 AND
> event_type='purchase';
Query ID = hadoop_20211207121646_afff7385-d1f1-474c-b112-20587aade2c5
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application 1638876533904 0004)
        VERTICES
                     MODE
                                   STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
Map 1 ..... container
                                 SUCCEEDED
Reducer 2 ..... container
1211538.4299998982
Time taken: 25.038 seconds, Fetched: 1 row(s)
hive> CREATE EXTERNAL TABLE IF NOT EXISTS retail_part_3 (event_time timestamp, event_type string,
   > product_id string, category_id string, category_code string, brand string, price decimal(10,3
    > user session string) PARTITIONED BY (month int) CLUSTERED BY (brand) INTO 5 buckets ROW FORMAT
    > SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS textfile;
Time taken: 0.098 seconds
hive> DESCRIBE retail_part_3;
OK
event_time
event_type
                                                   from deserializer
                                                    from deserializer
product_id
                                                   from deserializer
category_id
category_code
                                                   from deserializer
                                                   from deserializer
brand
                         string
                                                   from deserializer
                                                   from deserializer
                                                   from deserializer
user_session
                          string
```

hive> INSERT INTO TABLE retail_part_3 PARTITION (month) SELECT event_time, event_type, product_id, category_id, category_code, brand, price, user_id, user_session, MONTH(CAST(REPLACE(event_time,'UTC','') AS timestamp)) FROM retail;

```
hive> INSERT INTO TABLE retail_part_3 PARTITION (month) SELECT event_time, event_type, product_id, category_id, category_code, brand, price, user_id, user_session, MONTH(CAST(REPLACE(event_time,'UTC','') As timestamp)) FROM retail;
Query ID = hadoop_20211207122626_f63b8503-fd65-49f9-8675-0e986fdfe5aa
Total jobs = 1
Launching Job 1 out of 1
Tex session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1638876533904_0005)

VERTICES MODE STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED

Map 1 ...... container SUCCEEDED 2 2 0 0 0 0 0
Reducer 2 ..... container SUCCEEDED 5 5 5 0 0 0 0 0
VERTICES: 02/02 [==============>>>] 100% ELAPSED TIME: 188.50 s

Loading data to table default.retail_part_3 partition (month=null)

Loaded: 2/2 partitions.

Time taken to load dynamic partitions: 0.177 seconds

Time taken for adding to write entity: 0.001 seconds

OK
Time taken: 198.216 seconds
```

Executing the same query with the new table "retail part 3" to check the time.

hive> SELECT SUM(price) FROM retail_part_3 WHERE MONTH(event_time)=10 AND event_type='purchase';

Time Taken to execute the above query is 81.213 sec.

QUESTIONS

1. Find the total revenue generated due to purchases made in October.

hive> SELECT SUM(price) FROM retail_part_1 WHERE MONTH(event_time)=10 AND event_type='purchase';

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

hive> SELECT MONTH(event_time), SUM(price) as sum_purchase, COUNT(event_type) as cnt FROM retail_part_1 WHERE event_type='purchase' GROUP BY MONTH(event_time);

```
Query ID = hadoop_20211207123359_28681afd-e1df-40de-be0b-1b574ad5be2e
otal jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application 1638876533904 0005)
      VERTICES
                  MODE
                            STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
                          SUCCEEDED
Map 1 ..... container
                          SUCCEEDED
Reducer 2 ..... container
      1211538.4299998982
                          245624
      1531016.8999999384
                           322417
```

In October month, 245624 purchases generated revenue of 1211538.4299 Similarly in November month, 322417 purchases generated revenue of 1531016.8999.

Time Taken to execute the above query is 24.765 sec.

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

hive>WITH diff AS (SELECT SUM(CASE WHEN date_format(event_time,'MM')=10 THEN price ELSE 0 END) AS October, SUM(CASE WHEN date_format(event_time,'MM')=11 THEN price ELSE 0 END) AS November FROM retail_part_1 WHERE date_format(event_time,'MM') IN (10,11) AND event_type='purchase') SELECT October, November, (November - October) as Difference FROM diff;

```
hive> WITH diff AS ( SELECT SUM(CASE WHEN date_format(event_time, 'MM')=10 THEN price ELSE 0 END) AS
October, SUM(CASE WHEN date_format(event_time, 'MM')=11 THEN price ELSE 0 END) AS
> November FROM retail_part_1 WHERE date_format(event_time, 'MM') IN (10,11) AND event_type='pur chase') SELECT October, November, (November - October) as Difference FROM diff;
Query ID = hadoop_20211207123613_60012341-821f-4b43-a6d8-f200377b349a
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1638876533904_0005)

VERTICES MODE STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED

Map 1 ....... container SUCCEEDED 3 3 3 0 0 0 0 0
Reducer 2 .... container SUCCEEDED 1 1 0 0 0 0
VERTICES: 02/02 [=============>>>] 100% ELAPSED TIME: 44.32 s

OK
1211538.42999898 1531016.899999384 319478.4700000405
Time taken: 44.93 seconds, Fetched: 1 row(s)
```

The change in revenue generated from October to November is 319478.4700

4. Find distinct categories of products. Categories with null category code can beignored.

hive>SELECT DISTINCT split(category_code,'\\.')[0] AS category FROM retail_part_1 WHERE split(category_code,'\\.')[0]<>";

There are 6 distinct categories of products. They are Furniture, appliances, accessories, apparel, sport, and stationary.

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

hive>SELECT split(category_code,'\\.')[0] AS category, COUNT(product_id) AS prd FROM retail_part_1 GROUP BY split(category_code,'\\.')[0] ORDER BY prd DESC;

```
SELECT split(category_code,'\\.')[0] AS category, COUNT(product_id) AS prd FROM retail_part_1
   > GROUP BY split(category_code,'\\.')[0] ORDER BY prd DESC ;
uery ID = hadoop 20211207124025 93dcc1f7-eac7-4694-87d4-0c7f6b21d6d8
Total jobs = 1
Status: Running (Executing on YARN cluster with App id application 1638876533904 0005)
                                STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
       VERTICES
Map 1 ..... container SUCCEEDED
                             SUCCEEDED
Reducer 2 ..... container
Reducer 3 ..... container
                             SUCCEEDED
                                          =>>] 100% ELAPSED TIME: 71.55 s
       8594895
               61736
appliances
stationery
furniture
               23604
apparel 18232
               12929
accessories
sport 2
Time taken: 72.233 seconds, Fetched: 7 row(s)
```

The "Sport" category has the least number of products, whereas "appliances" have 61736 products.

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

SELECT brand, SUM(price) AS Sales FROM retail_part_1 WHERE brand <>" AND event_type='purchase' GROUP BY brand ORDER BY Sales DESC LIMIT 1;

Brand "runail" has the maximum sales for both months combined.

Time Taken to execute the above query is 23.509 sec.

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

hive>WITH monthly_diff AS (SELECT brand, SUM(CASE WHEN date_format(event_time,'MM')=10 THEN price ELSE 0 END) AS October, SUM(CASE WHEN date_format(event_time,'MM')=11 THEN price ELSE 0 END) AS November FROM retail_part_1 WHERE event_type='purchase' GROUP BY brand) SELECT brand, October, November, (November-October) as Sales_diff FROM monthly_diff WHERE (November-October) >0 ORDER BY Sales_diff;

```
hive> WiTH monthly_diff AS ( SELECT brand, SUM(CASE WHEN date_format(event_time,'MM')=10 THEN price
ELSE 0 END) AS October, SUM(CASE WHEN date_format(event_time,'MM')=11 THEN price ELSE 0 END) AS No
vember FROM retail_part_1 WHERE event_type='purchase' GROUP BY brand) SELECT brand, October, Novemb
er, (November-October) as Sales_diff FROM monthly_diff WHERE (November-October) >0 ORDER BY Sales_d
Query ID = hadoop 20211207125054 6c5dd93c-9708-4a24-9124-2e97b14ab794
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1638876533904_0006)
            VERTICES
                                  MODE
                                                     STATUS TOTAL COMPLETED RUNNING
                                                                                                             PENDING
                                                                                                                            FAILED KILLED
Map 1 ..... container
                                                SUCCEEDED
                                                SUCCEEDED
                                                SUCCEEDED
Reducer 3 ..... container
```

```
ovale
       2.54
               3.1
                       0.56
       20.22999999999997
                               20.93
                                       0.7000000000000028
grace
       100.9199999999999
                               102.610000000000001
                                                       1.69000000000000261
nelloganic
               0.0
                       3.1
                               3.1
skinity 8.88
               12.4400000000000001
                                       3.5600000000000005
bodyton 1376.3400000000006
                               1380.6400000000003
                                                       4.299999999999727
               10.2800000000000001
                                       4.570000000000001
moyou
       5.71
                       8.290000000000006
neoleor 43.41
               51.7
       204.1999999999998
                               212.5299999999999
                                                       8.330000000000098
soleo
                               1110.65 8.539999999999736
jaguar 1102.11000000000004
certio 236.16 245.79999999999999
                                        9.63999999999958
Ely
       17.14
               27.17
                      10.0300000000000001
rasyan 18.79999999999997
                               28.93999999999998
                                                        10.14
               316.84 329.17 12.330000000000041
deoproce
               12.39
barbie 0.0
                       12.39
               50.36999999999999
                                        66.5099999999999
supertan
               163.37 181.49000000000004
                                               18.120000000000033
reaclemoon
       63.00999999999999
                               81.490000000000002
                                                       18.480000000000032
kamill
iuno
       0.0
                       21.08
veraclara
               50.11
                       71.210000000000001
                                               21.10000000000001
glysolid
               69.72999999999999
                                        91.59
                                               21.860000000000014
               401.22 425.1199999999999
                                               23.8999999999992
godefroy
oinacil 0.0
               24.25999999999998
                                       24.25999999999998
       38.95
               63.4000000000000006
                                       24.4500000000000003
lixz
profepil
               93.36 118.02 24.65999999999999
```

```
71539.279999999 76758.65999999984
                                           5219.380000000849
                                                  5358.210000000005
oolarus 6013.719999999999
                           11371.930000000004
                                14536.990000000042
cosmoprofi 8322.809999999994
                                                         6214.180000000048
             26287.840000000127
essnail
                                   33345.23000000014
                                                          7057.390000000014
strong 29196.630000000005 38671.27000000002
                                                  9474.640000000014
ngarden 23161.38999999883 33566.210000000225 10404.820000000342
ianail 5892.839999999865 16394.23999999996 10501.39999999976
      35302.030000000006
                                                  15737.720000000067
                                                  36027.17000000348
grattol 35445.53999999999
                            71472.71000000341
      474679.0600000175
                            619509.2400000119
                                                  144830.17999999435
ime taken: 48.174 seconds, Fetched: 161 row(s)
```

Total of 161 brands have increased their sales from October to November.

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.

hive>SELECT user_id, SUM(price) AS expense FROM retail_part_1 WHERE event_type='purchase' GROUP BY user_id ORDER BY expense DESC LIMIT 10;

```
nive> SELECT user_id, SUM(price) AS expense FROM retail_part_1 WHERE event_type='purchase > GROUP BY user_id ORDER BY expense DESC LIMIT 10;
aunching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application 1638876533904 0006)
                                  STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
        VERTICES
                      MODE
Map 1 ..... container
Reducer 2 .... container
Reducer 3 .... container
                                SUCCEEDED
                                SUCCEEDED
                                SUCCEEDED
557790271
150318419
                1645.9700000000005
562167663
               1352.85
531900924
57850743
522130011
561592095
                1097.5899999999997
431950134
566576008
                1056.3600000000004
21347209
                1040.909999999999
ime taken: 24.447 seconds, Fetched: 10 row(s)
```

Above is the list of the top 10 users.

*** CLEANING UP**

Once the analysis is completed, delete the database &terminate the cluster, and stop EC2 instance.

```
hive> show DATABASES;

OK

default

hive_assignment

Time taken: 0.013 seconds, Fetched: 2 row(s)

hive> DROP DATABASE hive_assignment;

OK

Time taken: 0.187 seconds

hive> show DATABASES;

OK

default

Time taken: 0.009 seconds, Fetched: 1 row(s)
```

