

HIVE CASE STUDY

By Rakhee Kumari & Steven

Problem Statement

With online sales gaining popularity, tech companies are exploring ways to improve their sales by analysing customer behaviour 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.

Data for the case study is in the link given below.

<https://e-commerce-events-ml>.

<https://e-commerce-events-ml>.

We used a **2-node** EMR cluster with both the master and core nodes as **M4.large**.

IMPORTING THE DATA INTO HDFS

We login to Nuvepro dashboard, go to the console and then to EMR home page → Click on Create Cluster → select release EMR 5.29.0 and select required service for the case study.

1 Launching an EMR cluster that utilizes Hive services

The screenshot displays the AWS Management Console interface for an Amazon EMR cluster. The top navigation bar includes the AWS logo, a search bar, and the user's account information (N. Virginia, upgradstevencheriyil @ 0794-8777-7402). The left sidebar shows the navigation menu with categories like Amazon EMR, EMR Studio, EMR on EC2, Clusters, Notebooks, Git repositories, Security configurations, Block public access, VPC subnets, Events, EMR on EKS, Virtual clusters, Help, and What's new. The main content area shows the details for the cluster 'Hive_project', which is in the 'Waiting' state. A red warning box at the top indicates that auto-termination is not available for this account when using this release of EMR. The cluster details are organized into several sections: Summary, Configuration details, Application user interfaces, Network and hardware, and Security and access. The Summary section provides key information such as the cluster ID (j-OANF5ZAOP954), creation date (2022-04-27 16:25 UTC+5:30), elapsed time (55 minutes), and the master public DNS (ec2-100-25-163-224.compute-1.amazonaws.com). The Configuration details section shows the release label (emr-5.29.0), Hadoop distribution (Amazon 2.8.5), applications (Hive 2.3.6, Pig 0.17.0, Hue 4.4.0), and the log URI. The Application user interfaces section shows that persistent user interfaces are not enabled. The Network and hardware section shows the availability zone (us-east-1e), subnet ID (subnet-110ea620), and the number of master and core nodes (1 m4.large each). The Security and access section is partially visible at the bottom.

Amazon EMR

EMR Studio

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

What's new

Clone Terminate AWS CLI export

Cluster: Hive_project **Waiting** Cluster ready after last step completed.

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-OANF5ZAOP954

Creation date: 2022-04-27 16:25 (UTC+5:30)

Elapsed time: 55 minutes

After last step completes: Cluster waits

Termination protection: On [Change](#)

Tags: -- [View All / Edit](#)

Master public DNS: ec2-100-25-163-224.compute-1.amazonaws.com [Connect to the Master Node Using SSH](#)

Configuration details

Release label: emr-5.29.0

Hadoop distribution: Amazon 2.8.5

Applications: [Hive 2.3.6](#), [Pig 0.17.0](#), [Hue 4.4.0](#)

Log URI: [s3://aws-logs-079487777402-us-east-1/elasticmapreduce/](#)

EMRFS consistent view: Disabled

Custom AMI ID: --

Application user interfaces

Persistent user interfaces [🔗](#): --

On-cluster user interfaces [🔗](#): Not Enabled [Enable an SSH Connection](#)

Network and hardware

Availability zone: us-east-1e

Subnet ID: [subnet-110ea620](#) [🔗](#)

Master: **Running** 1 m4.large

Core: **Running** 1 m4.large

Task: --

Cluster scaling: Not enabled

Security and access

Feedback Looking for language selection? Find it in the new [Unified Settings](#) [🔗](#)

© 2022, Amazon Internet Services Private Ltd. or its affiliates. [Privacy](#) [Terms](#) [Cookie preferences](#)

2 Creating a folder on the Hadoop file system

```
hadoop@ip-172-31-61-99:~  
Using username "hadoop".  
Authenticating with public key "imported-openssh-key"  
Last login: Wed Apr 27 11:07:43 2022  
  
    _|_  _|_  )  
    _|_ (  _|_ /  Amazon Linux AMI  
    _|\_  _|_  |  
  
https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/  
72 package(s) needed for security, out of 102 available  
Run "sudo yum update" to apply all updates.  
  
EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRRRR  
E::::::::::::::::::::E M::::::::M          M::::::::M R::::::::::::R  
EE::::::::EEEEEEEEEE E M::::::::M          M::::::::M R::::RRRRRR::::R  
  E::::E          EEEEE M::::::::M          M::::::::M RR::::R      R::::R  
  E::::E          M::::M M::::M M::::M M::::M R:::R      R::::R  
  E::::EEEEEEEEEE M::::M M::::M M::::M M::::M R::RRRRRR::::R  
  E::::::::::E M::::M M::::M M::::M M::::M R::::::::::::RR  
  E::::EEEEEEEEEE M::::M M::::M M::::M M::::M R::RRRRRR::::R  
  E::::E          M::::M M::::M M::::M M::::M R:::R      R::::R  
  E::::E          EEEEE M::::M          MMM M::::M R:::R      R::::R  
EE::::::::EEEEEEEEEE E M::::M          M::::M R:::R      R::::R  
E::::::::::::::::::::E M::::M          M::::M RR::::R      R::::R  
EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRR      RRRRRR  
  
[hadoop@ip-172-31-61-99 ~]$ ls  
[hadoop@ip-172-31-61-99 ~]$ hadoop fs -mkdir /test-folder  
[hadoop@ip-172-31-61-99 ~]$ hadoop fs -ls /  
Found 5 items  
drwxr-xr-x - hdfs hadoop 0 2022-04-27 11:02 /apps  
drwxr-xr-x - hadoop hadoop 0 2022-04-27 11:11 /test-folder  
drwxrwxrwt - hdfs hadoop 0 2022-04-27 11:05 /tmp  
drwxr-xr-x - hdfs hadoop 0 2022-04-27 11:02 /user  
drwxr-xr-x - hdfs hadoop 0 2022-04-27 11:02 /var
```

3 Moving the data from S3 bucket into the HDFS

- Importing 2019-Nov.csv file from S3 to HDFS


hadoop@ip-172-31-61-99:~

```
-bash: distcp: command not found
[hadoop@ip-172-31-61-99 ~]$ hadoop distcp s3n://upgradprojectbucket/2019-Nov.csv
/tmp/test-folder/2019-Nov.csv
22/04/27 11:19:37 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='uniformsize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3n://upgradprojectbucket/2019-Nov.csv], targetPath=/tmp/test-folder/2019-Nov.csv, targetPathExists=false, filtersFile='null'}
22/04/27 11:19:38 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-61-99.ec2.internal/172.31.61.99:8032
22/04/27 11:19:44 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
22/04/27 11:19:44 INFO tools.SimpleCopyListing: Build file listing completed.
22/04/27 11:19:44 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
22/04/27 11:19:44 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
22/04/27 11:19:44 INFO tools.DistCp: Number of paths in the copy list: 1
22/04/27 11:19:44 INFO tools.DistCp: Number of paths in the copy list: 1
22/04/27 11:19:44 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-61-99.ec2.internal/172.31.61.99:8032
22/04/27 11:19:45 INFO mapreduce.JobSubmitter: number of splits:1
22/04/27 11:19:45 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1651057437121_0001
22/04/27 11:19:46 INFO impl.YarnClientImpl: Submitted application application_1651057437121_0001
22/04/27 11:19:47 INFO mapreduce.Job: The url to track the job: http://ip-172-31-61-99.ec2.internal:20888/proxy/application_1651057437121_0001/
22/04/27 11:19:47 INFO tools.DistCp: DistCp job-id: job_1651057437121_0001
22/04/27 11:19:47 INFO mapreduce.Job: Running job: job_1651057437121_0001
22/04/27 11:19:58 INFO mapreduce.Job: Job job_1651057437121_0001 running in uber mode : false
22/04/27 11:19:58 INFO mapreduce.Job: map 0% reduce 0%
22/04/27 11:20:16 INFO mapreduce.Job: map 100% reduce 0%
22/04/27 11:20:19 INFO mapreduce.Job: Job job_1651057437121_0001 completed successfully
22/04/27 11:20:19 INFO mapreduce.Job: Counters: 38
File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=172471
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=359
    HDFS: Number of bytes written=545839412
    HDFS: Number of read operations=12
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=4
```

hadoop@ip-172-31-61-99:~

```
22/04/27 11:19:58 INFO mapreduce.Job: map 0% reduce 0%
22/04/27 11:20:16 INFO mapreduce.Job: map 100% reduce 0%
22/04/27 11:20:19 INFO mapreduce.Job: Job job_1651057437121_0001 completed successfully
22/04/27 11:20:19 INFO mapreduce.Job: Counters: 38
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=172471
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=359
    HDFS: Number of bytes written=545839412
    HDFS: Number of read operations=12
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=4
    S3N: Number of bytes read=545839412
    S3N: Number of bytes written=0
    S3N: Number of read operations=0
    S3N: Number of large read operations=0
    S3N: Number of write operations=0
  Job Counters
    Launched map tasks=1
    Other local map tasks=1
    Total time spent by all maps in occupied slots (ms)=590112
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=18441
    Total vcore-milliseconds taken by all map tasks=18441
    Total megabyte-milliseconds taken by all map tasks=18883584
  Map-Reduce Framework
    Map input records=1
    Map output records=0
    Input split bytes=136
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=370
    CPU time spent (ms)=20990
    Physical memory (bytes) snapshot=563040256
    Virtual memory (bytes) snapshot=3299827712
    Total committed heap usage (bytes)=460324864
  File Input Format Counters
    Bytes Read=223
  File Output Format Counters
    Bytes Written=0
  DistCp Counters
    Bytes Copied=545839412
    Bytes Expected=545839412
    Files Copied=1
```

- **Importing 2019-Oct.csv file from S3 to HDFS**

 hadoop@ip-172-31-61-99:~

```
ov.csv
[hadoop@ip-172-31-61-99 ~]$ hadoop distcp s3n://upgradprojectbucket/2019-Oct.csv
/tmp/test-folder/2019-Oct.csv
22/04/27 11:28:28 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='uniformsize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3n://upgradprojectbucket/2019-Oct.csv], targetPath=/tmp/test-folder/2019-Oct.csv, targetPathExists=false, filtersFile='null'}
22/04/27 11:28:29 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-61-99.ec2.internal/172.31.61.99:8032
22/04/27 11:28:34 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
22/04/27 11:28:34 INFO tools.SimpleCopyListing: Build file listing completed.
22/04/27 11:28:34 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
22/04/27 11:28:34 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
22/04/27 11:28:34 INFO tools.DistCp: Number of paths in the copy list: 1
22/04/27 11:28:34 INFO tools.DistCp: Number of paths in the copy list: 1
22/04/27 11:28:34 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-61-99.ec2.internal/172.31.61.99:8032
22/04/27 11:28:34 INFO mapreduce.JobSubmitter: number of splits:1
22/04/27 11:28:35 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1651057437121_0003
22/04/27 11:28:35 INFO impl.YarnClientImpl: Submitted application application_1651057437121_0003
22/04/27 11:28:35 INFO mapreduce.Job: The url to track the job: http://ip-172-31-61-99.ec2.internal:20888/proxy/application_1651057437121_0003/
22/04/27 11:28:35 INFO tools.DistCp: DistCp job-id: job_1651057437121_0003
22/04/27 11:28:35 INFO mapreduce.Job: Running job: job_1651057437121_0003
22/04/27 11:28:45 INFO mapreduce.Job: Job job_1651057437121_0003 running in uber mode : false
22/04/27 11:28:45 INFO mapreduce.Job:  map 0% reduce 0%
22/04/27 11:29:02 INFO mapreduce.Job:  map 100% reduce 0%
22/04/27 11:29:05 INFO mapreduce.Job: Job job_1651057437121_0003 completed successfully
22/04/27 11:29:06 INFO mapreduce.Job: Counters: 38
    File System Counters
        FILE: Number of bytes read=0
        FILE: Number of bytes written=172471
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=359
        HDFS: Number of bytes written=482542278
        HDFS: Number of read operations=12
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=4
```

hadoop@ip-172-31-61-99:~

```
FILE: Number of bytes read=0
FILE: Number of bytes written=172471
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=359
HDFS: Number of bytes written=482542278
HDFS: Number of read operations=12
HDFS: Number of large read operations=0
HDFS: Number of write operations=4
S3N: Number of bytes read=482542278
S3N: Number of bytes written=0
S3N: Number of read operations=0
S3N: Number of large read operations=0
S3N: Number of write operations=0
Job Counters
  Launched map tasks=1
  Other local map tasks=1
  Total time spent by all maps in occupied slots (ms)=551072
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=17221
  Total vcore-milliseconds taken by all map tasks=17221
  Total megabyte-milliseconds taken by all map tasks=17634304
Map-Reduce Framework
  Map input records=1
  Map output records=0
  Input split bytes=136
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=410
  CPU time spent (ms)=18980
  Physical memory (bytes) snapshot=543019008
  Virtual memory (bytes) snapshot=3295203328
  Total committed heap usage (bytes)=445644800
File Input Format Counters
  Bytes Read=223
File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=482542278
  Bytes Expected=482542278
  Files Copied=1
```

4 Checking if the files are correctly imported to HDFS

hadoop fs -ls /hiveassignment

```
[hadoop@ip-172-31-49-96 ~]$ hadoop fs -ls /hiveassignment
Found 2 items
-rw-r--r-- 1 hadoop hadoop 545839412 2022-01-02 06:21 /hiveassignment/2019-Nov.csv
-rw-r--r-- 1 hadoop hadoop 482542278 2022-01-02 06:17 /hiveassignment/2019-Oct.csv
```

We can confirm the databases were loaded successfully.

```
[hadoop@ip-172-31-49-96 ~]$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive>
```

Creating database “upgrad_assignment”:

Create database if not exists upgrad_assignment;

use upgrad_assignment;

```
hive> Create database if not exists upgrad_assignment;
OK
Time taken: 0.065 seconds
```

```
hive> use upgrad_assignment;
OK
Time taken: 0.052 seconds
```

Creating an External Table, Sales:

create External table if not exists sales(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' WITH SERDEPROPERTIES ("separatorChar"=",", "quoteChar"="\\"", "escapeChar"="\") stored as textfile

Location '/hiveassignment' TBLPROPERTIES("skip.header.line.count"="1");

```
hive> create External Table if not exists sales(event_time timestamp,event_type
> string,product_id string,category_id string,category_code string,brand string,price f
> user_id bigint,user session string) ROW FORMAT SERDE
> 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
> WITH SERDEPROPERTIES ("separatorChar"=",", "quoteChar"="\\"", "escapeChar"="\")
> stored as textfile
> Location '/hiveassignment' TBLPROPERTIES("skip.header.line.count"="1");
OK
Time taken: 0.336 seconds
```

desc sales;

```
hive> desc sales;
OK
event_time          string              from deserializer
event_type           string              from deserializer
product_id           string              from deserializer
category_id          string              from deserializer
category_code        string              from deserializer
brand                string              from deserializer
price                string              from deserializer
user_id              string              from deserializer
user_session         string              from deserializer
Time taken: 0.418 seconds, Fetched: 9 row(s)
```

Loading the Data into the table:

hive> load data inpath '/hiveassignment/2019-Oct.csv' into table sales;

hive> load data inpath '/hiveassignment/2019-Nov.csv' into table sales;

```
hive> load data inpath '/hiveassignment/2019-Oct.csv' into table sales;
Loading data to table upgrad_assignment.sales
OK
Time taken: 3.264 seconds
```

```
hive> load data inpath '/hiveassignment/2019-Nov.csv' into table sales;
Loading data to table upgrad_assignment.sales
OK
Time taken: 1.149 seconds
```


We are required to provide answers to the questions given below:

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

```
hive> set hive.cli.print.header=true;
```

```
hive> select sum(price) from sales where Month(event_time)=10 and event_type='purchase';
```

```
hive> set hive.cli.print.header=true;
hive> select sum(price) from sales where Month(event_time)=10 and event_type='purchase';
Query ID = hadoop_20220102064210_59e412ed-d3e3-4769-b292-a6c50c8e59cb
Total jobs = 1
Launching Job 1 out of 1
Ter session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1641103761278_0004)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    1         1         0         0         0         0
-----
VERTICES: 02/02 [=====>>] 100% ELAPSED TIME: 40.41 s
-----
OK
_c0
1211530.4299997438
Time taken: 56.901 seconds, Fetched: 1 row(s)
hive>
```

Here the query takes 56.90 seconds which can be optimized by creating dynamic partition and then compare the execution time.

Dynamic Partitioning and Bucketing:

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

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

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

Creating a table by name sales_dp to store the dataset which we partitioned by using 'event_type' and clustered by 'user_id'.

```
Desc sales_dp;
```

```
hive> create External table if not exists sales_dp(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(user_id) into 5
> buckets ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' stored as
> TextFile;
OK
Time taken: 0.089 seconds
hive> desc sales_dp;
OK
col_name      data_type      comment
event_time     string          from deserializer
product_id     string          from deserializer
category_id    string          from deserializer
category_code  string          from deserializer
brand          string          from deserializer
price         string          from deserializer
user_id        string          from deserializer
user_session   string          from deserializer
event_type     string
# Partition Information
# col_name      data_type      comment
event_type     string
```

Loading the data into the new table:

insert into sales_dp partition(event_type) select event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type from sales;

```
hive> insert into sales_dp partition(event_type) select event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type from sales;
Query ID = hadoop_20220102065934_89d5fa6-8c44-41b8-ac8f-b495018e87e8
Total jobs = 1
Launching Job 1 out of 1
782 session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1641103761278_0005)

-----
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: 92/92  [=====] 100% ELAPSED TIME: 113.54 s
-----
Loading data to table upgrad_assignment.sales_dp partition (event_type=null)
Loaded : 4/4 partitions.
Time taken to load dynamic partitions: 0.617 seconds
Time taken for adding to write entity : 0.007 seconds
OK
event_time    product_id    category_id    category_code    brand    price    user_id    user_session    event_type
Time taken: 126.946 seconds
hive>
```

Now executing the same Q1:

select sum(price) from sales_dp where Month(event_time)=10 and event_type='purchase';

```
hive> select sum(price) from sales_dp where Month(event_time)=10 and
> event_type='purchase';
Query ID = hadoop_20220102070525_e507749a-b0fa-4f5b-afc8-ce22feeffbd2
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_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: 17.53 s
-----
OK
_c0
1211538.4299999105
Time taken: 18.169 seconds, Fetched: 1 row(s)
```

We can notice how the time taken reduced drastically due to partitioning and bucketing. Now it took only 18.17 sec.

The total sales in the month of October is 1211538.42.

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

select Month(event_time) as Month, sum(price) as sum, count(event_type) as cnt from sales where event_type='purchase' group by Month(event_time);

```
hive> select Month(event_time) as Month, sum(price) as sum, COUNT(event_type) as cnt
> from sales where event_type='purchase' group by Month(event_time);
Query ID = hadoop_20220102070705_5d14246c-98bc-4352-bc60-454a103166c0
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_0005)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    3         3         0         0         0         0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 37.67 s
-----
OK
month  sum      cnt
10     1211538.4265997438    245624
11     1531016.900000122    322417
Time taken: 38.284 seconds, Fetched: 2 row(s)
hive>
```

In the month of October, the total purchase is 245624 and sales is 1211538.42

In the month of November, the total purchase is 322417 and sales is 1531016.90

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

with CTE1 as (select sum(case when Month(event_time)=10 then price else 0 end) as Oct, sum(case when Month(event_time)=11 then price else 0 end) as Nov from sales_dp where event_type='purchase' and Month(event_time) in (10,11)) select Oct,Nov,(Nov-Oct) as diff from CTE1;

```
hive> with CTE1 as (select sum(case when Month(event_time)=10 then price else 0 end) as Oct, sum(case when Month(event_time)=11 then price else 0 end) as Nov from sales_dp where event_type='purchase' and Month(event_time) in (10,11) ) select Oct,Nov,(Nov-Oct) as diff from CTE1;
Query ID = hadoop_20220102070857_d1393f88-239a-4fbd-888c-3351b447ea16
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_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: 19.94 s
-----
OK
oct  nov  diff
1211538.4265999105    1531016.9000000267    319478.47000001613
Time taken: 20.694 seconds, Fetched: 1 row(s)
hive>
```

We can see the difference in the revenue is 319478.47

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

select distinct split(category_code,'\\\.')[0] as Categories from sales_dp where split(category_code,'\\\.')[0]<>'';

```
hive> select distinct split(category_code,'\\\.')[0] as Categories from sales_dp where
> split(category_code,'\\\.')[0]<>'';
Query ID = hadoop_20220102071052_af73873d-0f35-414f-8bf1-a299f35b0985
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_0005)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED   14      14          0         0         0         0
Reducer 2 ..... container  SUCCEEDED    5         5          0         0         0         0
-----
VERTICES: 02/02 [=====>>>] 100% ELAPSED TIME: 57.04 s
-----
OK
categories
furniture
appliances
accessories
apparel
sport
stationery
Time taken: 57.803 seconds, Fetched: 6 row(s)
```

We can see the 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 Prodcoun from sales group by split(category_code,'\\\.')[0] order by Prodcoun desc;

```
hive> select split(category_code,'\\\.')[0] as category, count(product_id) as Prodcoun from sales group by split(category_code,'\\\.')[0] order by Prodcoun desc;
Query ID = hadoop_20220102071551_600baf7-0eb8-45e9-9a7d-b5bb28f6c043
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_0005)

-----
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
Reducer 3 ..... container  SUCCEEDED    1         1          0         0         0         0
-----
VERTICES: 03/03 [=====>>>] 100% ELAPSED TIME: 52.54 s
-----
OK
category      prodcoun
8594895
appliances    61736
stationery    26722
furniture     23604
apparel 18232
accessories   12929
sport         2
Time taken: 53.145 seconds, Fetched: 7 row(s)
hive>
```

The total number of products under each category is as follows: Appliances 61736, Stationery 26722, Furniture 23604, Apparel 18232, Accessories 12929, Sport 2

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

select brand, sum(price) as totalsales from sales_dp where brand <>' ' and event_type='purchase' group by brand order by totalsales desc limit 1;

```
hive> select brand, sum(price) as totalsales from sales_dp where brand <>' ' and
> event_type='purchase' group by brand order by totalsales desc limit 1;
Query ID = hadoop_20220102071804_975b0c72-6115-47de-86e3-12979ff0bd27
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_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: 16.75 s
OK
brand totalsales
runail 148297.93999999954
Time taken: 17.426 seconds, Fetched: 1 row(s)
```

We can see that Runail is the brand with the maximum sales for oct and nov. Total sales is 148297.94

7 Which brands increased their sales from October to November?

with CTE2 as (select brand, sum(case when Month(event_time)=10 then price else 0 end) as Oct, sum(case when Month(event_time)=11 then price else 0 end) as Nov from sales_dp where event_type='purchase' and group by brand) select brand, Oct,Nov,(Nov-Oct) as diff from CTE2 where (Nov-Oct)>0 ORDER BY diff;

```
hive> with CTE2 as(select brand, sum(case when month(event_time)=10 then price else 0
> end) as Oct,sum(case when month(event_time)=11 then price else 0 end) as Nov from
> sales_dp where event_type='purchase' group by brand) select brand , Oct,Nov,(Nov-Oct) as
> diff from CTE2 where (Nov-Oct)>0 ORDER BY diff;
Query ID = hadoop_20220102071955_2f02bd7f-96ec-46ef-a5fb-29217a6cf34f
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_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: 10.92 s
OK
brand oct nov diff
ovale 2.54 3.1 0.56
cosima 20.229999999999997 20.929999999999996 0.6999999999999991
grace 100.91999999999997 102.61 1.69000000000000241
helloganic 0.0 3.1 3.1
skinity 8.88 12.440000000000001 3.5600000000000005
bodytton 1276.24000000000004 1280.6400000000002 4.2999999999999727
moyoa 5.71 10.280000000000001 4.5700000000000001
neolcor 43.41 51.7 8.2900000000000006
solao 204.19999999999997 212.52999999999996 8.3300000000000155
jaguar 1102.11000000000001 1110.65 8.5399999999999964
tertio 238.16000000000003 245.79999999999998 7.6399999999999958
fly 17.14 27.17 10.030000000000001
rasyan 18.799999999999997 28.939999999999998 10.14
deoproce 316.84 329.16999999999996 12.329999999999994
barbie 0.0 12.39 12.39
supertan 50.36999999999999 66.51 16.140000000000015
treaclemoon 163.37 181.49000000000004 18.120000000000033
kamill 63.009999999999998 81.49000000000001 18.480000000000018
tuan 0.0 31.08 31.08
```


juno	0.0	21.08	21.08		
veraclara		50.11	71.21000000000001	21.10000000000001	
glysolid		69.72999999999999	91.59	21.860000000000014	
godefroy		401.22	425.12	23.899999999999977	
binacil 0.0		24.259999999999998		24.259999999999998	
blixx 30.95		63.400000000000006		24.450000000000003	
profepil		93.36	118.02000000000002	24.660000000000025	
estelara		444.81000000000006	471.87000000000002	27.060000000000116	
orly	902.38	931.09	28.710000000000036		
biore	60.650000000000006		90.31	29.659999999999997	
besutyblender	78.740000000000001		109.41	30.669999999999987	
villanta	107.50000000000007		231.20000000000003	31.600000000000006	
mavala	409.04000000000001		446.32000000000005	37.27999999999997	
likato	296.06	340.97	44.910000000000025		
ladykin	125.65	170.57	44.91999999999999		
foamie	35.04	80.49	45.449999999999996		
elskin	251.08999999999997		307.65	56.56	
balbcara	155.330000000000004		212.380000000000005	57.050000000000001	
koelcia 55.5		112.75000000000001	57.250000000000014		
profhenna		679.22999999999998	736.8499999999997	57.619999999999989	
kares 0.0		59.449999999999996	59.449999999999996		
marutaka-foot		49.22	109.33	60.11	
dewal 0.0		61.289999999999999	61.289999999999999		
lm	288.02	351.21	63.19		
laboratorium	246.49999999999997		312.52	66.020000000000001	
cutrin	499.37	361.62	68.49		
egomania		77.47	146.04	68.57	
konad	735.8299999999999		810.67000000000002	70.840000000000026	
hirwel	163.040000000000002		234.33	71.289999999999999	
koelf	422.7299999999999		507.28999999999996	84.560000000000006	
plazan	101.37	194.01	92.63999999999999		
aura	83.95	177.51	93.55999999999999		
kerasys	430.91000000000001		525.1999999999999	94.289999999999985	
enjoy	41.349999999999994		136.570000000000002	95.220000000000003	
depillflax	2707.0699999999974		2803.7799999999998	96.710000000000049	
eos	54.339999999999999		152.60999999999999	98.27	
carmex	145.07999999999998		243.36	98.280000000000003	
batiste	772.40000000000001		874.17	101.76999999999997	
osmo	645.580000000000002		762.31000000000001	116.72999999999999	
dizao	819.13	945.51000000000003	126.380000000000034		
igrochauty		513.860000000000002	645.07	131.409999999999985	
finish 90.38		230.30	132.0		
nefertiti		233.520000000000004	366.64	133.11999999999995	
elizabeth		70.53	204.3	133.77	
miskin	158.040000000000002		293.06999999999994	135.02999999999992	
latinoil		249.52	384.590000000000003	135.070000000000002	
famona	1692.460000000000005		1843.43	150.96999999999997	
cristalinas		427.63	584.95	157.320000000000005	
chi	358.93999999999994		538.61000000000001	179.67000000000002	

chi	358.93999999999994		538.61000000000001	179.67000000000002	
matreshka	0.0	182.670000000000002	182.670000000000002		
freshbubble		318.7	502.340000000000003	183.640000000000004	
mane	66.789999999999999		260.26	193.47	
keen	236.34999999999997		435.620000000000006	199.270000000000001	
ecocraft		41.1600000000000004	241.95	200.79	
fedus	52.38	263.81	211.43		
provoc	827.98999999999994		1063.8199999999993	235.82999999999998	
skinlite		651.940000000000003	890.45000000000004	238.51000000000001	
entity	479.710000000000009		719.25999999999995	239.54999999999985	
trind	298.070000000000005		342.96	244.89	
pronokeratin	201.25	456.74	255.540000000000002		
beauugreen		511.510000000000005	768.3499999999999	256.83999999999986	
bluesky	10307.2399999999982		10565.5299999999948	258.2900000000006636	
candy	534.95999999999996		799.37999999999999	264.41999999999996	
insight	1443.700000000000003		1721.960000000000005	278.260000000000002	
kocostar		310.84999999999997	594.93	284.08	
happyfens		801.920000000000003	1091.590000000000006	289.670000000000003	
kims	330.03999999999996		632.040000000000001	302.000000000000001	
shary	871.95999999999996		1176.48999999999999	304.52999999999993	
nitrite	847.27999999999999		1162.67999999999998	315.4	
lowence	242.840000000000003		367.75	324.90999999999997	
jas		3318.960000000000002	3657.43000000000017	338.46999999999998	
ellips	245.850000000000002		606.03999999999998	360.18999999999998	
lador	2083.610000000000024		2471.530000000000025	387.920000000000001	
naomi 0.0	344.0	484.0			
kiss	421.55	817.32999999999997	395.77999999999997		
yu-r	271.40999999999997		673.71	402.300000000000007	
sophin	1067.860000000000004		1515.520000000000002	447.659999999999965	
farmavita		837.370000000000001	1291.96999999999998	454.59999999999997	
bioaqua	942.890000000000001		1398.120000000000003	455.230000000000025	
greymy	29.21	488.49	460.280000000000003		
qehwol	1089.070000000000002		1557.680000000000003	468.610000000000001	
matrix	3243.250000000000001		3726.740000000000007	483.48999999999998	
limoni	1308.9	1796.600000000000006	487.700000000000005		
s.care	412.68	913.07	500.390000000000004		
coifin	903.0	1426.46999999999998	525.48999999999998		
uskusi	5142.27000000000011		5690.30999999999995	548.03999999999981	
alnails		5118.860000000000006	5601.520000000000004	572.620000000000053	
browenna		14331.3699999999982	14916.730000000000005	585.360000000000133	
kinetics		6334.25000000000018	6945.36000000000015	611.009999999999966	
kosmekka		1181.44	1813.37	631.929999999999998	
kaatal	4412.430000000000003		5086.070000000000004	673.640000000000012	
refectocil		2716.180000000000002	3475.580000000000005	759.400000000000028	
rosi	3077.040000000000001		3841.560000000000002	764.520000000000013	
solomeya		1899.7	2685.79999999999984	786.099999999999963	
missha	1293.83	2150.27999999999997	856.44999999999998		
levissime		2227.500000000000004	3085.310000000000104	857.810000000000063	
art-visage		2092.710000000000002	2997.800000000000007	905.090000000000051	

```

ecolab 262.6499999999997 1214.3000000000002 951.4500000000003
nagaraku 4369.7400000000011 5327.6600000000004 957.9399999999932
vanoto 157.14 1209.6799999999998 1052.54
markail 1768.7499999999982 2834.4799999999994 1065.6800000000012
metzger 5373.4500000000001 6457.1599999999996 1083.7099999999955
de.lux 1659.6999999999966 2775.5099999999988 1115.8099999999913
swarovski 1887.9299999999955 3043.1599999999917 1155.2299999999962
beauty-free 554.1700000000001 1782.8600000000022 1228.6900000000002
reitun 708.6600000000002 2009.6299999999997 1300.9699999999993
joico 705.52 2015.1000000000001 1309.5800000000002
severina 4775.8799999999985 6120.479999999998 1344.5999999999993
irick 4444.444444444444 4444.444444444444 1344.0800000000044
onig 8425.41 9841.649999999996 1416.2399999999961
levrana 2243.5599999999995 3664.1000000000002 1420.5400000000027
roubloff 3491.3600000000006 4913.7700000000005 1422.4100000000044
smart 4457.2599999999995 5902.1399999999995 1444.86
shik 3341.2000000000016 4839.7200000000003 1498.5200000000013
domix 10472.05 12009.1699999999982 1537.11999999999826
artex 2730.6399999999994 4327.2500000000001 1596.61000000000015
beautix 10493.9499999999986 12222.949999999999 1729.0000000000036
milv 3904.9399999999978 5642.0099999999994 1737.0700000000152
masura 31266.0799999999012 33058.469999999964 1792.39000000006252
f.o.x 6624.23 8577.2799999999982 1953.0499999999992
kopous 11927.1500000000074 14993.0600000000005 2165.9200000000346
concept 11032.1399999999989 13380.3999999999932 2348.2599999999944
astel 21756.749999999993 24142.670000000002 2385.9200000000089
kaypro 881.34 3266.6999999999994 2387.3599999999999
benovy 409.61999999999985 3259.9700000000003 2850.3500000000004
italwax 21940.2399999999918 24799.3699999999904 2859.1299999999985
yoko 8756.91 11707.8799999999986 2950.9699999999986
haruyama 9390.6899999999988 12352.9100000000014 2962.22000000001267
marathon 7280.7800000000002 10273.1 2982.3499999999995
lovely 8704.3799999999994 11039.0599999999974 3234.67999999999803
bpw.style 11572.1500000000067 14837.4400000000197 3265.290000000013
staleks 8519.7300000000018 11875.6100000000015 3355.8799999999974
freedecor 3421.7799999999986 7671.799999999995 4250.0199999999964
runail 71539.279999999998 76758.659999999987 5219.3800000000082
polarus 6013.72 11371.9300000000004 5358.2100000000004
cosmoprofi 8322.9099999999992 14536.9900000000038 6214.1800000000046
jessnail 26287.840000000127 33345.230000000015 7057.3900000000021
strong 29196.629999999997 38671.270000000002 9474.6400000000021
ingarden 23161.389999999977 33566.210000000009 10404.8200000000316
lianail 5892.8399999999964 16394.239999999972 10501.399999999999
uno 35302.029999999994 51039.750000000009 15727.720000000103
grettol 35445.539999999996 71472.710000000335 36027.1700000003396
474679.0600000162 619509.2400000163 144830.18000000001
Time taken: 19.658 seconds, Fetched: 161 row(s)

```

From the output we can see that 161 brand were able to increase their sales from the month of 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.

select user_id, sum(price) as Totalpurchases from sales_dp where event_type ='purchase' group by user_id order by Totalpurchases DESC limit 10;

```
hive> select user_id, sum(price) as Totalpurchases from sales_dp where event_type='purchase'
> group by user_id order by Totalpurchases DESC limit 10;
Query ID = hadoop_20220102072427_ccd6cc4-2a39-42d0-9acb-e42b30f9bee0
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1641103761278_0005)
```

	VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	3	3	0	0	0	0	0
Reducer 2	container	SUCCEEDED	1	1	0	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0	0

```
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 17.59 s
OK
user_id totalpurchases
557790271 2715.869999999999
150318419 1645.970000000000
562167663 1352.95
531900924 1329.45
557850743 1295.480000000000
522130011 1185.389999999999
561592095 1109.700000000000
421950134 1097.589999999999
566576098 1056.360000000000
521347209 1040.909999999999
Time taken: 18.289 seconds, Fetched: 10 row(s)
hive>
```

We can see the top 10 users in the output.

Finishing Up

Once we are done, we can drop the databases, quit the hive and then terminate the EMR cluster.