

Hive Case Study

- Srilathaa Vasu, Shaik Athaullah

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. 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.

This is done by tracking their clicks on the website and searching for patterns within them. The clickstream data contains all the logs as to how the customer navigated through the website. It also contains other details such as time spent on every page, etc. From this, tech companies make use of data ingesting frameworks such as Apache Kafka or AWS Kinesis in order to store it in frameworks such as Hadoop. From there, machine learning engineers or business analysts use this data to derive valuable insights.

Data for this case study was available in:

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

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

Objective:

To extract data and gather insights from an e-commerce company dataset, using AWS EMR and S3, Hadoop and Hive systems.

Steps Involved:

1. Launching an EMR Cluster
2. Loading the data into S3 and ingesting the data to HDFS

3. Launching and creating Hive Schema
4. Optimization of tables
5. Querying – Assignment Questions
6. Improvement of performance after using optimization on tables
7. Cleaning-up

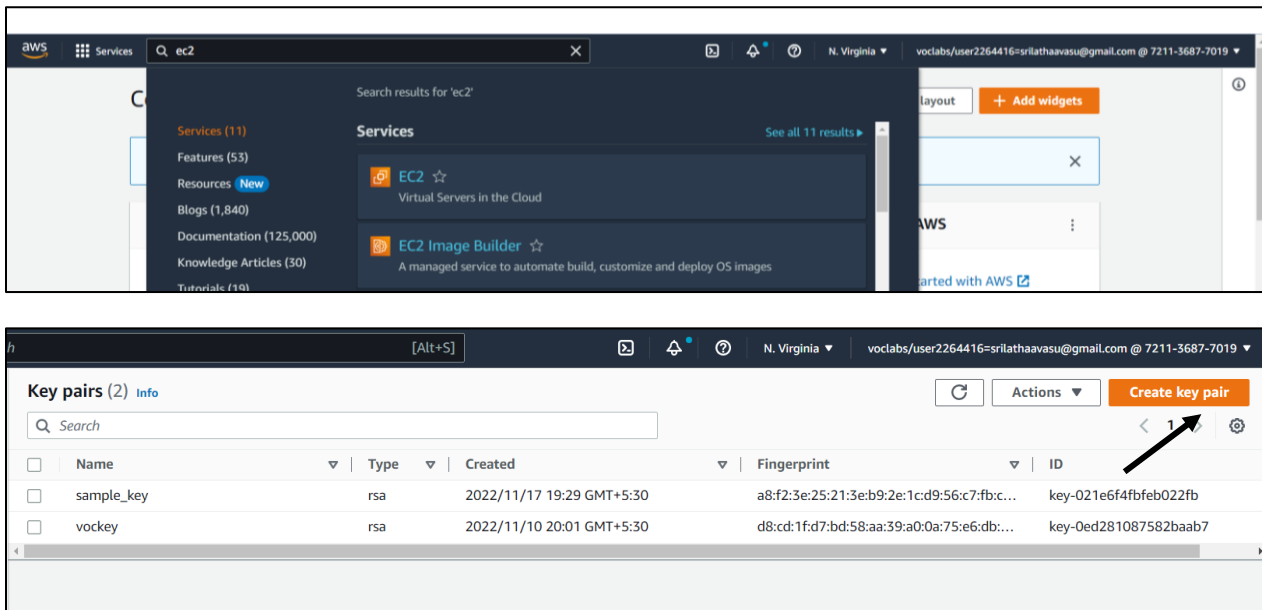
1. Launching an EMR Cluster:

The following steps are followed while launching an EMR cluster,

- i. Key-pair creation
- ii. Configuring & Launching an EMR cluster
- iii. Connecting to master node using SSH
- iv. Launching of Hive shell

i. Key-pair creation:

We need to create an EC2 key-pair file and download it as .pem file.



Key pair
A key pair, consisting of a private key and a public key, is a set of security credentials that you use to prove your identity when connecting to an instance.

Name

The name can include up to 255 ASCII characters. It can't include leading or trailing spaces.

Key pair type [info](#)
☒ RSA
☐ ED25519

Private key file format
☒ .pem
 For use with OpenSSH
☐ .ppk
 For use with PuTTY

Tags - *optional*
No tags associated with the resource.

[Add new tag](#)
You can add up to 50 more tags.

[Cancel](#) [Create key pair](#)

Successfully created key pair

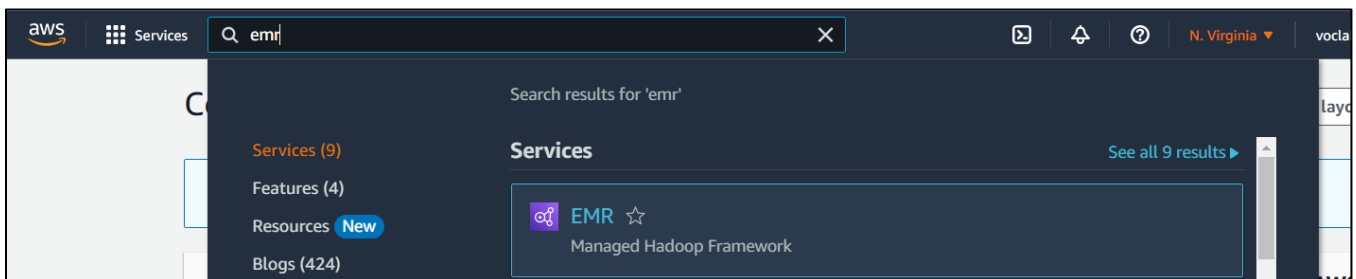
Key pairs (3) [Info](#)

<input type="checkbox"/>	Name	Type	Created	Fingerprint	ID
<input checked="" type="checkbox"/>	case_study	rsa	2022/11/28 22:51 GMT+5:30	21:9b:3e:c4:76:21:e0:ee:41:52:bc:35:b2:...	key-052eae7183ca88876
<input type="checkbox"/>	sample_key	rsa	2022/11/17 19:29 GMT+5:30	a8:f2:3e:25:21:3e:b9:2e:1c:d9:56:c7:fb:c...	key-021e6f4fbfeb022fb
<input type="checkbox"/>	vockey	rsa	2022/11/10 20:01 GMT+5:30	d8:cd:1f:d7:bd:58:aa:39:a0:0a:75:e6:db:...	key-0ed281087582baab7

We will use this private key to securely connect to the EMR cluster through SSH.

ii. Configuring & Launching an EMR cluster:

Amazon EMR (previously called Amazon Elastic MapReduce) is a managed cluster platform that simplifies running big data frameworks, such as Hadoop and Apache Spark on AWS to process and analyze vast amounts of data.



EMR Serverless is now GA.
With EMR Serverless, get the benefits of Amazon EMR such as open source compatibility, latest versions and performance optimized runtime for job startup, automatic capacity management, and simple cost controls. [Get Started with EMR Serverless.](#)

Create cluster

View details

Clone

Terminate

Filter: All clusters Filter clusters ... 5 clusters (all loaded)

	Name	ID	Status	Creation time (UTC+5:30)
<input type="checkbox"/>	emr_quering	j-EWMTE7M4RSQZ	Terminated Auto-terminate	2022-11-18 16:03 (UTC+5:30)
<input type="checkbox"/>	Hive_cluster	j-L4Y9Q71Z05TH	Terminated User request	2022-11-17 19:44 (UTC+5:30)
<input type="checkbox"/>	hive_emr	j-1ACS7IF8XNITY	Terminated User request	2022-11-17 19:42 (UTC+5:30)

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

☒ Logging

S3 folder

Launch mode ☒ Cluster ☐ Step execution

Software Configuration

Release

- | | | |
|--|---|--|
| <input checked="" type="checkbox"/> Hadoop 2.8.5 | <input type="checkbox"/> Zeppelin 0.8.2 | <input type="checkbox"/> Livy 0.6.0 |
| <input type="checkbox"/> JupyterHub 1.0.0 | <input type="checkbox"/> Tez 0.9.2 | <input type="checkbox"/> Flink 1.9.1 |
| <input type="checkbox"/> Ganglia 3.7.2 | <input type="checkbox"/> HBase 1.4.10 | <input checked="" type="checkbox"/> Pig 0.17.0 |
| <input checked="" type="checkbox"/> Hive 2.3.6 | <input type="checkbox"/> Presto 0.227 | <input type="checkbox"/> ZooKeeper 3.4.14 |
| <input type="checkbox"/> MXNet 1.5.1 | <input type="checkbox"/> Sqoop 1.4.7 | <input type="checkbox"/> Mahout 0.13.0 |
| <input checked="" type="checkbox"/> Hue 4.4.0 | <input type="checkbox"/> Phoenix 4.14.3 | <input type="checkbox"/> Oozie 5.1.0 |
| <input type="checkbox"/> Spark 2.4.4 | <input type="checkbox"/> HCatalog 2.3.6 | <input type="checkbox"/> TensorFlow 1.14.0 |

Multiple master nodes (optional)

☐ Use multiple master nodes to improve cluster availability. [Learn more](#)

AWS Glue Data Catalog settings (optional)

☐ Use for Hive table metadata

Edit software settings

☒ Enter configuration ☐ Load JSON from S3

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 Master - 1	m4.large 2 vCore, 8 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Core Core - 2	m4.large 2 vCore, 8 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	<input type="text" value="1"/> Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

+ Add task instance group

We have chosen emr-5.29.0 release for this case study.

General Options

Cluster name

☒ Logging ⓘ

S3 folder

☒ Debugging ⓘ

☒ Termination protection ⓘ

Tags ⓘ

Key	Value (optional)
<input type="text" value="Add a key to create a tag"/>	<input type="text"/>

Additional Options

☐ EMRFS consistent view ⓘ

Custom AMI ID ⓘ

Security Options

EC2 key pair ⓘ

☒ Cluster visible to all IAM users in account ⓘ

Permissions ⓘ

☒ Default ☐ Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

EMR role [EMR_DefaultRole](#) ⓘ ☐ Use EMR_DefaultRole_V2 ⓘ

EC2 instance profile [EMR_EC2_DefaultRole](#) ⓘ

Auto Scaling role [EMR_AutoScaling_DefaultRole](#) ⓘ

▼ Security Configuration

Security configuration ⓘ

▶ EC2 security groups

We have selected the “case_study” private key which we have created earlier.

Cluster: cosmetics_cluster **Starting**

[Summary](#) [Application user interfaces](#) [Monitoring](#) [Hardware](#) [Configurations](#) [Events](#) [Steps](#) [Bootstrap actions](#)

Summary

ID: j-28HOL3ZDHQ8SH
Creation date: 2022-11-29 19:56 (UTC+5:30)
Elapsed time: 1 second
After last step completes: Cluster waits
Termination protection: On [Change](#)
Tags: -- [View All / Edit](#)
Master public DNS: --

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-721136877019-us-east-1/elasticmapreduce/
EMRFS consistent view: Disabled
Custom AMI ID: --

Application user interfaces

Persistent user interfaces : --
On-cluster user -- interfaces :

Network and hardware

Availability zone: --
Subnet ID: [subnet-005f018cb791095a6](#)
Master: Provisioning 1 m4.large
Core: Provisioning 1 m4.large
Task: --
Cluster scaling: Not enabled

Security and access

Key name: case_study
EC2 instance profile: EMR_EC2_DefaultRole
EMR role: EMR_DefaultRole
Auto Scaling role: EMR_AutoScaling_DefaultRole
Visible to all users: All [Change](#)

Our cluster Hive Assignment is created and launched successfully and is now in “Waiting” state.

Clone
Terminate
AWS CLI export
⚠ Auto-termination is not available for this account when using this release of EMR.

Cluster: cosmetics_cluster
Waiting
Cluster ready after last step completed.

Summary
Application user interfaces
Monitoring
Hardware
Configurations
Events
Steps
Bootstrap actions

Summary

ID: j-28HOL3ZDHQ8SH

Creation date: 2022-11-29 19:56 (UTC+5:30)

Elapsed time: 12 minutes

After last step completes: Cluster waits

Termination protection: On [Change](#)

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

Master public DNS: ec2-54-208-163-205.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-721136877019-us-east-1/elasticmapreduce/ [View Log](#)

EMRFS consistent view: Disabled

Custom AMI ID: --

Application user interfaces

Persistent user interfaces [View](#): --

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

Network and hardware

Availability zone: us-east-1d

Subnet ID: subnet-005f018cb791095a6 [View](#)

Master: Running 1 m4.large

Core: Running 1 m4.large

Task: --

Cluster scaling: Not enabled

Security and access

Key name: case_study

EC2 instance profile: EMR_EC2_DefaultRole

EMR role: EMR_DefaultRole

Auto Scaling role: EMR_AutoScaling_DefaultRole

We need to ensure that the port is open to establish a connection. Under the cluster information page, click on the security groups of the master node. Here, we add an additional inbound rule as below.

sgr-03d97385932325906
Custom TCP
TCP
8443
Custom
Q
Delete

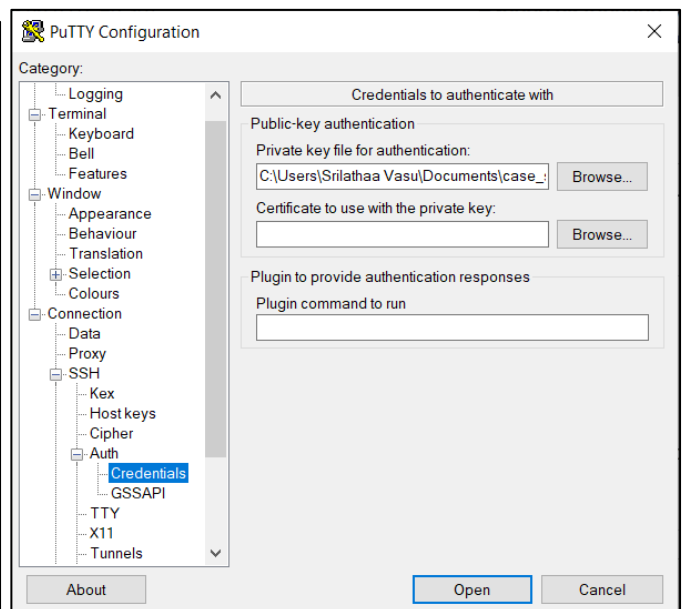
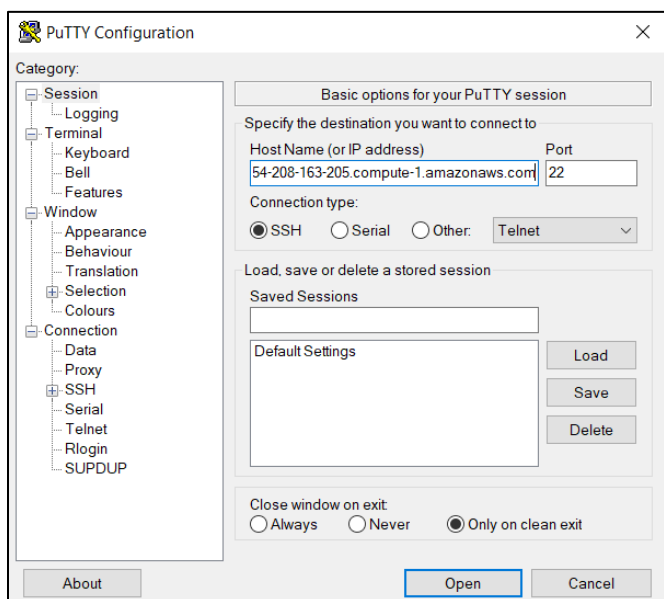
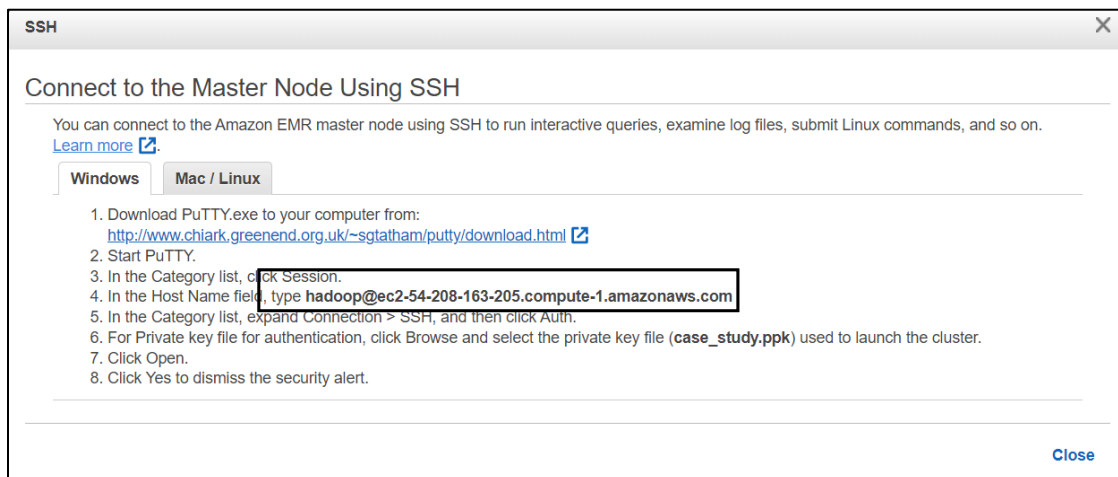
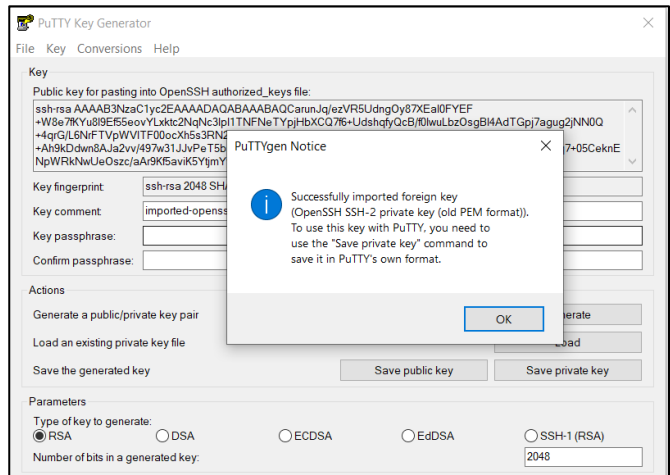
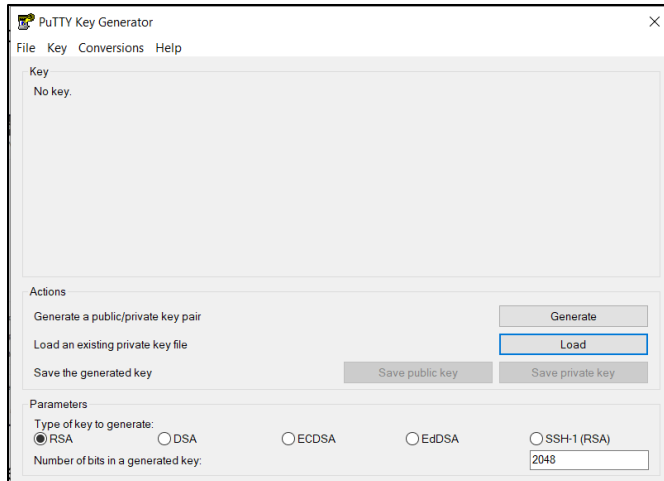
-
SSH
TCP
22
Anywh...
72.21.198.64/29 X
Delete

Add rule

Cancel
Preview changes
Save rules

iii. Connecting to master node using SSH:

For connecting the master node using SSH, we need to use two programs namely, PuTTY and PuTTYgen. We first need to open PuTTYgen, then load the private key(case_study) and save it as .ppk file at a desired location. Next, we need to open PuTTY to connect to master node. We need to copy Master DNS address from the EMR cluster summary page.



iv. Launching of Hive Shell:

```
hadoop@ip-172-31-87-236:~$
Using username "hadoop".
Authenticating with public key "imported-openssh-key"
Last login: Tue Nov 29 15:36:18 2022

  _ _ | _ _ | _ _ |
  _ _ | _ _ | _ _ |

Amazon Linux AMI

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
69 package(s) needed for security, out of 97 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMM MMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::::E M::::::::M M::::::::M R::::::::::::R
EE::::::::EEEEEEEEEE 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 R::::RRRRR::::R
E::::EEEEEEEEEE M::::M M::::M M::::M R::::RRRRR::::R
E::::E M::::M M::::M M::::M R::::R R::::R
E::::E EEEEE M::::M M M 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 MMMMMMM MMMMMMM RRRRRRR RRRRRR

[hadoop@ip-172-31-87-236 ~]$
```

2. Loading the data into S3 and ingesting the data to HDFS:

Uploading databases to S3 buckets:

Create bucket [Info](#)

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

Bucket name must be globally unique and must not contain spaces or uppercase letters. [See rules for bucket naming](#)

AWS Region

US East (N. Virginia) us-east-1

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

Choose bucket

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

cosmetics-case-study [Info](#)

Objects Properties Permissions Metrics Management Access Points

Objects (2)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	2019-Nov.csv	csv	November 29, 2022, 21:52:31 (UTC+05:30)	520.6 MB	Standard
<input type="checkbox"/>	2019-Oct.csv	csv	November 29, 2022, 21:52:31 (UTC+05:30)	460.2 MB	Standard

Creating a temporary directory in HDFS:

Checking the directories already present in HDFS.

Command: *hadoop fs -ls /*

Output:

```
[hadoop@ip-172-31-87-236 ~]$ hadoop fs -ls /
Found 4 items
drwxr-xr-x  - hdfs hadoop      0 2022-11-29 14:33 /apps
drwxrwxrwt  - hdfs hadoop      0 2022-11-29 14:35 /tmp
drwxr-xr-x  - hdfs hadoop      0 2022-11-29 14:33 /user
drwxr-xr-x  - hdfs hadoop      0 2022-11-29 14:33 /var
```

Creating a temporary directory.

Command: *hadoop fs -mkdir /user/hivecasestudy*
hadoop fs -ls /user/

Output:

```
[hadoop@ip-172-31-87-236 ~]$ hadoop fs -mkdir /user/hivecasestudy
[hadoop@ip-172-31-87-236 ~]$ hadoop fs -ls /user/
Found 7 items
drwxrwxrwx  - hadoop hadoop      0 2022-11-29 14:33 /user/hadoop
drwxr-xr-x  - mapred mapred      0 2022-11-29 14:33 /user/history
drwxrwxrwx  - hdfs  hadoop      0 2022-11-29 14:33 /user/hive
drwxr-xr-x  - hadoop hadoop      0 2022-11-29 16:50 /user/hivecasestudy
drwxrwxrwx  - hue   hue        0 2022-11-29 14:33 /user/hue
drwxrwxrwx  - oozie oozie      0 2022-11-29 14:33 /user/oozie
drwxrwxrwx  - root  hadoop      0 2022-11-29 14:33 /user/root
[hadoop@ip-172-31-87-236 ~]$
```

Loading the data to HDFS:

Since the size of the data is large we'll load the data into HDFS from S3 and into the local storage.

Command:

```
hadoop distcp 's3://cosmetics-case-study/2019-Oct.csv' /user/hivecasestudy/2019-Oct.csv
```

```
hadoop distcp 's3://cosmetics-case-study/2019-Nov.csv' /user/hivecasestudy/2019-Nov.csv
```

Output:

```
hadoop@ip-172-31-95-196:~$ hadoop distcp 's3://cosmetics-case-study/2019-Oct.csv' /user/hivecasestudy/2019-Oct.csv
22/11/30 11:17:30 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=[s3://cosmetics-case-study/2019-Oct.csv], targetPath=/user/hivecasestudy/2019-Oct.csv, targetPathExists=false, filtersFile='null')
22/11/30 11:17:30 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-95-196.ec2.internal/172.31.95.196:8032
22/11/30 11:17:35 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
22/11/30 11:17:35 INFO tools.SimpleCopyListing: Build file listing completed.
22/11/30 11:17:35 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
22/11/30 11:17:35 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
22/11/30 11:17:35 INFO tools.DistCp: Number of paths in the copy list: 1
22/11/30 11:17:35 INFO tools.DistCp: Number of paths in the copy list: 1
22/11/30 11:17:35 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-95-196.ec2.internal/172.31.95.196:8032
22/11/30 11:17:35 INFO mapreduce.JobSubmitter: number of splits:1
22/11/30 11:17:35 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1669806251379_0003
22/11/30 11:17:36 INFO impl.YarnClientImpl: Submitted application application_1669806251379_0003
22/11/30 11:17:36 INFO mapreduce.Job: The url to track the job: http://ip-172-31-95-196.ec2.internal:20888/proxy/application_1669806251379_0003/
22/11/30 11:17:36 INFO tools.DistCp: DistCp job-id: job_1669806251379_0003
22/11/30 11:17:36 INFO mapreduce.Job: Running job: job_1669806251379_0003
22/11/30 11:17:44 INFO mapreduce.Job: Job job_1669806251379_0003 running in uber mode : false
22/11/30 11:17:44 INFO mapreduce.Job: map 0% reduce 0%
22/11/30 11:18:03 INFO mapreduce.Job: map 100% reduce 0%
22/11/30 11:18:06 INFO mapreduce.Job: Job job_1669806251379_0003 completed successfully
22/11/30 11:18:06 INFO mapreduce.Job: Counters: 38

File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=172497
  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
  S3: Number of bytes read=482542278
  S3: Number of bytes written=0
  S3: Number of read operations=0
  S3: Number of large read operations=0
  S3: 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)=606848
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=18964
  Total vcore-milliseconds taken by all map tasks=18964
  Total megabyte-milliseconds taken by all map tasks=19419136

Map-Reduce Framework
  Map input records=1
  Map output records=0
  Input split bytes=136
  Spilled Records=0

hadoop@ip-172-31-95-196:~$ hadoop distcp 's3://cosmetics-case-study/2019-Nov.csv' /user/hivecasestudy/2019-Nov.csv
22/11/30 11:08:00 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=[s3://cosmetics-case-study/2019-Nov.csv], targetPath=/user/hivecasestudy/2019-Nov.csv, targetPathExists=false, filtersFile='null')
22/11/30 11:08:00 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-95-196.ec2.internal/172.31.95.196:8032
22/11/30 11:08:06 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
22/11/30 11:08:06 INFO tools.SimpleCopyListing: Build file listing completed.
22/11/30 11:08:06 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
22/11/30 11:08:06 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
22/11/30 11:08:06 WARN hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
  at java.lang.Object.wait(Native Method)
  at java.lang.Thread.join(Thread.java:1257)
  at java.lang.Thread.join(Thread.java:1331)
  at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:973)
  at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:624)
  at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:801)
22/11/30 11:08:06 INFO tools.DistCp: Number of paths in the copy list: 1
22/11/30 11:08:06 INFO tools.DistCp: Number of paths in the copy list: 1
22/11/30 11:08:06 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-95-196.ec2.internal/172.31.95.196:8032
22/11/30 11:08:06 INFO mapreduce.JobSubmitter: number of splits:1
22/11/30 11:08:06 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1669806251379_0002
22/11/30 11:08:07 INFO impl.YarnClientImpl: Submitted application application_1669806251379_0002
22/11/30 11:08:07 INFO mapreduce.Job: map 0% reduce 0%
22/11/30 11:08:07 INFO mapreduce.Job: Counters: 38

File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=172497
  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
  S3: Number of bytes read=482542278
  S3: Number of bytes written=0
  S3: Number of read operations=0
  S3: Number of large read operations=0
  S3: 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)=601856
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=18808
  Total vcore-milliseconds taken by all map tasks=18808
  Total megabyte-milliseconds taken by all map tasks=19259392

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)=274
  CPU time spent (ms)=20510
  Physical memory (bytes) snapshot=614346752
  Virtual memory (bytes) snapshot=3299844096
  Total committed heap usage (bytes)=502792192

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
```

Checking the loaded files:

Command: *hadoop fs -ls /user/hivecasestudy/*

Output:

```
[hadoop@ip-172-31-95-196 ~]$ hadoop fs -ls /user/hivecasestudy/
Found 2 items
-rw-r--r--  1 hadoop hadoop  545839412 2022-11-30 11:08 /user/hivecasestudy/2019-Nov.csv
-rw-r--r--  1 hadoop hadoop  482542278 2022-11-30 11:18 /user/hivecasestudy/2019-Oct.csv
[hadoop@ip-172-31-95-196 ~]$
```

We can see that the datasets have been loaded successfully.

3. Launching and creating Hive Schema:

Launching Hive:

Command: *hive*

Output:

```
[hadoop@ip-172-31-95-196 ~]$ hive

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

Creating a new database for the case study:

Command: *create database if not exists hive_casestudy;*
use hive_casestudy;

Output:

```
hive> create database if not exists hive_casestudy;
OK
Time taken: 1.123 seconds
hive> use hive_casestudy;
OK
```

Create an external table to load the data:

Command: *CREATE EXTERNAL TABLE IF NOT EXISTS sales_details(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 '/hivecasestudy'
>TBLPROPERTIES("skip.header.line.count"="1");

Output:

```
hive> CREATE EXTERNAL TABLE IF NOT EXISTS sales_details(event_time timestamp, ev
ent_type string, product_id string, category_id string, category_code string, br
and string, price float, user_id bigint, user_session string) ROW FORMAT SERDE '
org.apache.hadoop.hive.serde2.OpenCSVSerde'
> STORED AS TEXTFILE LOCATION '/user/hivecasestudy/'
> TBLPROPERTIES ("skip.header.line.count"="1");
OK
Time taken: 0.542 seconds
```

Checking the table 'sales_details':

Command: *desc sales_details;*

Output:

```
hive> desc sales_details;
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: 0.18 seconds, Fetched: 9 row(s)
```

Loading the data into the 'sales_details' table and checking:

Command: *load data inpath '/user/hivecasestudy/2019-Oct.csv' into table sales_details;*
load data inpath '/user/hivecasestudy/2019-Nov.csv' into table sales_details;

Output:

```
hive> load data inpath '/user/hivecasestudy/2019-Oct.csv' into table sales_details;
Loading data to table default.sales_details
OK
Time taken: 1.511 seconds
hive> load data inpath '/user/hivecasestudy/2019-Nov.csv' into table sales_details;
Loading data to table default.sales_details
OK
Time taken: 0.842 seconds
```

Command: *select * from sales_details*
>limit 5;

Output:

```
hive> select * from sales_details
> limit 5;
OK
sales_details.event_time      sales_details.event_type      sales_details.product_id
d      sales_details.category_id      sales_details.category_code      sales_details.b
rand      sales_details.price      sales_details.user_id      sales_details.user_session
2019-11-01 00:00:02 UTC view      5802432 1487580009286598681      0.32 5
62076640      09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:09 UTC cart      5844397 1487580006317032337      2.38 5
53329724      2067216c-31b5-455d-alcc-af0575a34ffb
2019-11-01 00:00:10 UTC view      5837166 1783999064103190764      pnb 22.22 5
56138645      57ed222e-a54a-4907-9944-5a875c2d7f4f
2019-11-01 00:00:11 UTC cart      5876812 1487580010100293687      jessnail 3
.16      564506666      186c1951-8052-4b37-adce-dd9644b1d5f7
2019-11-01 00:00:24 UTC remove from cart      5826182 1487580007483048900      3
.33      553329724      2067216c-31b5-455d-alcc-af0575a34ffb
Time taken: 0.632 seconds, Fetched: 5 row(s)
```

Hive Commands for easy view:

Commands: *set hive.cli.print.header=True;*
set hive.resultset.use.unique.column.names=false;
set hive.exec.dynamic.partition=true;
set hive.exec.dynamic.partition.mode=nonstrict;
set hive.enforce.bucketing=true;

Output:

```
hive> set hive.cli.print.header=True;
hive> set hive.resultset.use.unique.column.names=false;
hive> set hive.exec.dynamic.partition=true;
hive> set hive.exec.dynamic.partition.mode=nonstrict;
hive> set hive.enforce.bucketing=true;
```

Explanation:

- *set hive.cli.print.header=True;* - To display the column header in output
- *set hive.resultset.use.unique.column.names=false;* - To remove table name before column name separated by ‘.’
- *set hive.exec.dynamic.partition=true;*
- *set hive.exec.dynamic.partition.mode=nonstrict;* - To enable Dynamic partitioning
- *set hive.enforce.bucketing=true;* - To enable bucketing

4. Optimization of Tables – Partitioning and Bucketing:

Partitioning and Bucketing are query optimisation techniques. Bucketing is a method of segregating the tables or partitions in hive into parts based on the values of a column which in turn reduces the query execution time.

Attribute 'event_type' has unique set of categorical values and it will give us a definite number of partitions (4). So, partitioning will be performed on the attribute 'event_type'. This will also reduce the load to the name node.

Bucketing will be performed on the attribute 'price'. Depending on the data size and the size of HDFS block, the number of buckets will be created. We have calculated to create a minimum of eight buckets.

Command: *create external table cosmetics_par(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 6 buckets
> row format serde 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
> stored as textfile location '/user/hivecasestudy/'
> tblproperties ("skip.header.line.count"="1");

Output:

```
hive> create external table cosmetics_par (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 6 buckets
> row format serde 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
> stored as textfile location '/user/hivecasestudy/'
> tblproperties ("skip.header.line.count"="1");
OK
Time taken: 0.307 seconds
```

Checking the new table cos_opt:

Command: *desc cosmetics_par;*

```
hive> desc cosmetics_par;
OK
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
Time taken: 0.467 seconds, Fetched: 14 row(s)
```

```

Map 1: 4/4      Reducer 2: 4(+3)/9
Map 1: 4/4      Reducer 2: 4(+3)/9
Map 1: 4/4      Reducer 2: 4(+3)/9
Map 1: 4/4      Reducer 2: 4(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 5(+3)/9
Map 1: 4/4      Reducer 2: 6(+2)/9
Map 1: 4/4      Reducer 2: 6(+3)/9
Map 1: 4/4      Reducer 2: 6(+3)/9
Map 1: 4/4      Reducer 2: 6(+3)/9
Map 1: 4/4      Reducer 2: 6(+3)/9
Map 1: 4/4      Reducer 2: 7(+2)/9
Map 1: 4/4      Reducer 2: 7(+2)/9
Map 1: 4/4      Reducer 2: 7(+2)/9
Map 1: 4/4      Reducer 2: 7(+2)/9
Map 1: 4/4      Reducer 2: 8(+1)/9
Map 1: 4/4      Reducer 2: 8(+1)/9
Map 1: 4/4      Reducer 2: 9/9
Loading data to table hive_casestudy.cosmetics_par partition (event_type=null)

Time taken to load dynamic partitions: 0.383 seconds
Time taken for adding to write entity : 0.006 seconds
OK

```

Verifying partitioning in Hive:

Command: *show partitions cosmetics_par;*

Output:

```

hive> show partitions cosmetics_par;
OK
event_type=cart
event_type=purchase
event_type=remove_from_cart
event_type=view
Time taken: 0.144 seconds, Fetched: 4 row(s)

```

Verifying partitioning and bucketing in Hadoop:

Command: *hadoop fs -ls /user/hivecasestudy/*

Output:

```
[hadoop@ip-172-31-83-157 ~]$ hadoop fs -ls /user/hivecasestudy/
Found 4 items
drwxr-xr-x  - hadoop hadoop      0 2022-12-01 12:51 /user/hivecasestudy/event_type=cart
drwxr-xr-x  - hadoop hadoop      0 2022-12-01 12:52 /user/hivecasestudy/event_type=purchase
drwxr-xr-x  - hadoop hadoop      0 2022-12-01 12:52 /user/hivecasestudy/event_type=remove_from_cart
drwxr-xr-x  - hadoop hadoop      0 2022-12-01 12:52 /user/hivecasestudy/event_type=view
```

We can observe that four partitions have been created.

Command: *hadoop fs -ls /user/hivecasestudy/event_type=purchase*

Output:

```
[hadoop@ip-172-31-83-157 ~]$ hadoop fs -ls /user/hivecasestudy/event_type=purchase/
Found 6 items
-rwxr-xr-x  1 hadoop hadoop  23939488 2022-12-01 12:51 /user/hivecasestudy/event_type=purchase/000000_0
-rwxr-xr-x  1 hadoop hadoop  23978398 2022-12-01 12:51 /user/hivecasestudy/event_type=purchase/000001_0
-rwxr-xr-x  1 hadoop hadoop  23228108 2022-12-01 12:50 /user/hivecasestudy/event_type=purchase/000002_0
-rwxr-xr-x  1 hadoop hadoop  23613588 2022-12-01 12:51 /user/hivecasestudy/event_type=purchase/000003_0
-rwxr-xr-x  1 hadoop hadoop  24113362 2022-12-01 12:51 /user/hivecasestudy/event_type=purchase/000004_0
-rwxr-xr-x  1 hadoop hadoop  24183396 2022-12-01 12:50 /user/hivecasestudy/event_type=purchase/000005_0
```

From the output, we can observe that 6 buckets have been created with index starting from 0-5.

Command : *select * from cosmetics_par limit 5;*

Output:

```
hive> select * from cosmetics_par
> limit 5;
OK
2019-10-06 20:46:50 UTC 5853630 1487580006317032337 11.10 557505
362 dfec5efd-eb4a-4d7c-8dc9-0d976fflead8 cart
2019-10-08 09:48:05 UTC 5885214 2069804417665728971 cosmoprofi 5.56 5
58012670 a7a775fa-c671-43ae-928d-667eab5e6e48 cart
2019-10-01 07:00:05 UTC 5847446 1487580009286598681 1.03 555511
769 d29c7b8f-041d-4cbb-acf3-67859a4cc5db cart
2019-10-08 09:48:03 UTC 5885022 2069804417665728971 cosmoprofi 5.56 5
58012694 f0cda670-7691-41d0-b7d8-c445c0e34926 cart
2019-10-06 20:46:49 UTC 5853630 1487580006317032337 11.10 557505
362 dfec5efd-eb4a-4d7c-8dc9-0d976fflead8 cart
Time taken: 3.334 seconds, Fetched: 5 row(s)
```

Now, we have created two tables namely, sales_details and cosmetics_par. We will be using 'cosmetics_par' table for analysing the data which is a partitioned and bucketed table.

```
hive> show tables;
OK
cosmetics_par
sales_details
Time taken: 0.059 seconds, Fetched: 2 row(s)
```

5. Querying – Assignment Questions:

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

Command: *SELECT ROUND(SUM(price),2) AS total_sales*
> FROM cosmetics_par
> WHERE event_type='purchase' AND
> month(event_time)=10;

Output

```
hive> SELECT ROUND(SUM(price),2) AS total_sales
> FROM cosmetics_par
> WHERE event_type='purchase' AND
> month(event_time)=10;
Query ID = hadoop_20221201181907_2c4cda9e-3169-4036-8e6d-4ca117f872fd
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0012)

-----
      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: 24.80 s
-----
OK
total_sales
1211534.31
Time taken: 25.387 seconds, Fetched: 1 row(s)
```

- The total revenue generated in the month of October is Rs.1211534.31

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

Command: *SELECT Month(event_time) as Month, ROUND(SUM(price),2) as total_purchases*
>FROM cosmetics_par
>WHERE event_type='purchase'
>GROUP BY month(event_time);

Output:

```
hive> SELECT month(event_time) as month, ROUND(SUM(price),2) AS total_purchases
> FROM cosmetics_par
> WHERE event_type='purchase'
> GROUP BY month(event_time);
Query ID = hadoop_20221201171315_af38aabd-820a-4376-90fb-eb057953f474
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 25.97 s
-----
OK
month  total_purchases
10     1211534.31
11     1531007.31
Time taken: 26.668 seconds, Fetched: 2 row(s)
```

- In the month of October, the total purchases is Rs.1211534.31
- In the month of November, the total purchases is Rs.1531007.31

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

Command: *SELECT ROUND((Rev_Nov – Rev_Oct),2) AS revenue_change*
>FROM (
>SELECT SUM (CASE WHEN month(event_time)=10 THEN price ELSE 0
END) AS Rev_Oct,
> SUM (CASE WHEN month(event_time)=11 THEN price ELSE 0

```

        END) AS Rev_Nov,
>FROM cosmetics_par,
>WHERE event_type='purchase' AND
>month(event_time) in (10,11)
> )s;

```

Output:

```

hive> SELECT ROUND((Rev_Nov - Rev_Oct),2) AS revenue_change
> FROM (
> SELECT SUM(CASE when month(event_time)=10 then price else 0 END) AS Rev_Oct,
> SUM(CASE when month(event_time)=11 then price else 0 END) AS Rev_Nov
> FROM cosmetics_par
> WHERE event_type='purchase' AND
> month(event_time) in (10,11)
> )s;
Query ID = hadoop_20221201170815_2b6b05f1-dc12-4bb6-9e93-e1b3824a1316
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_1669909606668_0010)

-----
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: 27.85 s
-----
OK
revenue_change
319473.0
Time taken: 37.808 seconds, Fetched: 1 row(s)

```

- From the output, we can infer that the difference in revenue is Rs. 3,19,473.00

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

Command: *SELECT DISTINCT SPLIT (category_code, '\\.')[0] AS Distinct_category_list*
>FROM cosmetics_par;

Output:

```
hive> SELECT DISTINCT SPLIT(category_code,'\\\.')[0] AS Distinct_category_list
> FROM cosmetics_par;
Query ID = hadoop_20221201171658_87fa2462-8b39-4361-ac9f-1fdf5851b5f7
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 74.57 s
-----
OK
distinct_category_list

furniture
appliances
accessories
apparel
sport
stationery
Time taken: 75.222 seconds, Fetched: 7 row(s)
```

- We can observe that there are 6 distinct categories namely, furniture, appliances, accessories, apparel, sport and stationery.

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

```
Command: SELECT SPLIT (category_code, '\\\.')[0] AS category, COUNT(product_id) AS
total_prod
>FROM cosmetics_par
>GROUP BY SPLIT(category_code, '\\\.')[0]
>ORDER BY total_prod;
```

Output:

```
hive> SELECT SPLIT(category_code,'\\\.')[0] AS category, COUNT(product_id) AS total_prod
> FROM cosmetics_par
> GROUP BY SPLIT(category_code,'\\\.')[0]
> ORDER BY total_prod;
Query ID = hadoop_20221201172308_d41edd06-4e6b-4b40-9b0c-5252c4285efd
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 74.83 s
-----
OK
category      total_prod
sport      2
accessories  12929
apparel 18232
furniture    23604
stationery   26722
appliances   61736
8594871
Time taken: 75.845 seconds, Fetched: 7 row(s)
```

The total number of products under each category are as follows:

- Sport : 2
- Accessories : 12929
- Apparel : 18232
- Furniture : 23604
- Stationery : 26722
- Appliances : 61736

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

Command: *SELECT brand, ROUND(SUM(price),2) AS sales_total*
 >FROM cosmetics_par
 >WHERE event_type= 'purchase'
 >GROUP BY brand
 > ORDER BY sales_total DESC
 >LIMIT 5;

Output:

```
hive> SELECT brand,ROUND(SUM(price),2) as sales_total
> FROM cosmetics_par
> WHERE event_type='purchase'
> GROUP BY brand
> ORDER BY sales_total DESC
> LIMIT 5;
Query ID = hadoop_20221201173150_1986cdf5-d678-408e-b7b6-66e6c8ca9de5
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 23.76 s
-----
OK
brand    sales_total
runail   148297.94
grattol  106918.25
irisk    92538.0
uno      86341.78
Time taken: 24.547 seconds, Fetched: 5 row(s)
```

- From the output, we can observe that ‘Runail’ is a popular brand with high sales in both months. The total sales is Rs.148297.94

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

Command: *WITH sales_total 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 cosmetics_par
>WHERE event_type = 'purchase' AND month(event_time) in (10,11)
>GROUP BY brand )
>SELECT brand, ROUND(oct_sales,2), ROUND(nov_sales,2), ROUND((nov_sales-
oct_sales),2) AS increased_sales
>FROM sales_total
>WHERE (nov_sales-oct_sales)>0
>ORDER BY increased_sales DESC;
```

Output:

```
hive> WITH sales_total 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 cosmetics_par  
  > WHERE event_type='purchase' AND  
  > month(event_time) in (10,11)  
  > GROUP BY brand  
  > )  
  > SELECT brand, ROUND(oct_sales,2), ROUND(nov_sales,2), ROUND((nov_sales-oct_sales),2) AS increased_  
sales  
  > FROM sales_total  
  > WHERE (nov_sales-oct_sales)>0  
  > ORDER BY increased_sales DESC;
```

Query ID = hadoop_20221201175202_0fb2b8e6-a47c-4960-8afa-0183562c6b88

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_1669909606668_0011)

```
-----  
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: 34.66 s

OK

brand	_c1	_c2	increased_sales
	474677.31	619502.46	144825.15
grattol	35445.54	71472.71	36027.17
uno	35302.03	51039.75	15737.72
lianail	5892.84	16394.24	10501.4
ingarden	23161.39	33566.21	10404.82
strong	29196.63	38671.27	9474.64
jessnail	26287.84	33345.23	7057.39
cosmoprofi	8322.81	14536.99	6214.18
polarus	6013.72	11371.93	5358.21
runail	71539.28	76758.66	5219.38
freedecor	3421.78	7671.8	4250.02
staleks	8519.73	11875.61	3355.88
bpw.style	11572.15	14836.36	3264.21

staleks	8519.73	11875.61	3355.88	
bpw.style		11572.15	14836.36	3264.21
lovely	8704.38	11939.06	3234.68	
marathon		7280.75	10273.1	2992.35
haruyama		9390.69	12352.91	2962.22
yoko	8756.91	11707.88	2950.97	
italwax	21940.24		24799.37	2859.13
benovy	409.62	3259.97	2850.35	
kaypro	881.34	3268.7	2387.36	
estel	21756.75		24142.67	2385.92
concept	11032.14		13380.4	2348.26
kapous	11927.16		14093.08	2165.92
f.o.x	6624.23	8577.28	1953.05	
masura	31263.71		33056.74	1793.03
milv	3904.94	5642.01	1737.07	
beautix	10493.95		12222.95	1729.0
artex	2730.64	4327.25	1596.61	
domix	10472.05		12009.17	1537.12
shik	3341.2	4839.72	1498.52	
smart	4457.26	5902.14	1444.88	
roubloff		3491.36	4913.77	1422.41
levrana	2243.56	3664.1	1420.54	
oniq	8425.41	9841.65	1416.24	
irisk	45591.96		46946.04	1354.08
severina		4775.88	6120.48	1344.6
joico	705.52	2015.1	1309.58	
zeitun	708.66	2009.63	1300.97	
beauty-free	554.17	1782.86	1228.69	
swarovski		1887.93	3043.16	1155.23
de.lux	1659.7	2775.51	1115.81	
metzger	5373.45	6457.16	1083.71	
markell	1768.75	2834.43	1065.68	
sanoto	157.14	1209.68	1052.54	
nagaraku		4369.74	5327.68	957.94
ecolab	262.85	1214.3	951.45	
art-visage		2092.71	2997.8	905.09
levissime		2227.5	3085.31	857.81
missha	1293.83	2150.28	856.45	
solomeya		1899.7	2685.8	786.1
rosi	3077.04	3841.56	764.52	
refectocil		2716.18	3475.58	759.4
kaaral	4412.43	5086.07	673.64	
kosmekka		1181.44	1813.37	631.93
kinetics		6334.25	6945.26	611.01
browxenna		14331.37	14916.73	585.36
airnails		5118.9	5691.52	572.62

kosmekka	1181.44	1813.37	631.93	
kinetics	6334.25	6945.26	611.01	
browxenna	14331.37		14916.73	585.36
airnails	5118.9	5691.52	572.62	
uskusi	5142.27	5690.31	548.04	
coifin	903.0	1428.49	525.49	
s.care	412.68	913.07	500.39	
limoni	1308.9	1796.6	487.7	
matrix	3243.25	3726.74	483.49	
gehwol	1089.07	1557.68	468.61	
greymy	29.21	489.49	460.28	
bioaqua	942.89	1398.12	455.23	
farmavita		837.37	1291.97	454.6
sophin	1067.86	1515.52	447.66	
yu-r	271.41	673.71	402.3	
kiss	421.55	817.33	395.78	
naomi	0.0	389.0	389.0	
lador	2083.61	2471.53	387.92	
ellips	245.85	606.04	360.19	
jas	3318.96	3657.43	338.47	
lowence	242.84	567.75	324.91	
nitrile	847.28	1162.68	315.4	
shary	871.96	1176.49	304.53	
kims	330.04	632.04	302.0	
happyfons		801.92	1091.59	289.67
kocostar		310.85	594.93	284.08
insight	1443.7	1721.96	278.26	
candy	534.96	799.38	264.42	
bluesky	10307.24		10565.53	258.29
beauugreen		511.51	768.35	256.84
protokeratin		201.25	456.79	255.54
trind	298.07	542.96	244.89	
entity	479.71	719.26	239.55	
skinlite		651.94	890.45	238.51
provoc	827.99	1063.82	235.83	
fedua	52.38	263.81	211.43	
ecocraft		41.16	241.95	200.79
keen	236.35	435.62	199.27	
mane	66.79	260.26	193.47	
freshbubble		318.7	502.34	183.64
matreshka		0.0	182.67	182.67
chi	358.94	538.61	179.67	
cristalinas		427.63	584.95	157.32
farmona	1692.46	1843.43	150.97	

cristalinas	427.63	584.95	157.32
farmona	1692.46	1843.43	150.97
latinoil	249.52	384.59	135.07
maskin	158.04	293.07	135.03
elizavecca	70.53	204.3	133.77
nefertiti	233.52	366.64	133.12
finish	98.38	230.38	132.0
igrobeauty	513.66	645.07	131.41
dizao	819.13	945.51	126.38
osmo	645.58	762.31	116.73
batiste	772.4	874.17	101.77
carmex	145.08	243.36	98.28
eos	54.34	152.61	98.27
depilflax	2707.07	2803.78	96.71
enjoy	41.35	136.57	95.22
kerasys	430.91	525.2	94.29
aura	83.95	177.51	93.56
plazan	101.37	194.01	92.64
koelf	422.73	507.29	84.56
nirvel	163.04	234.33	71.29
konad	739.83	810.67	70.84
egomania	77.47	146.04	68.57
cutrin	299.37	367.62	68.25
laboratorium	246.5	312.52	66.02
inm	288.02	351.21	63.19
dewal	0.0	61.29	61.29
marutaka-foot	49.22	109.33	60.11
kares	0.0	59.45	59.45
profhenna	679.23	736.85	57.62
koelcia	55.5	112.75	57.25
balbcare	155.33	212.38	57.05
elskin	251.09	307.65	56.56
foamie	35.04	80.49	45.45
ladykin	125.65	170.57	44.92
likato	296.06	340.97	44.91
mavala	409.04	446.32	37.28
vilenta	197.6	231.21	33.61
beautyblender	78.74	109.41	30.67
biore	60.65	90.31	29.66
orly	902.38	931.09	28.71

```

ladykin 125.65 170.57 44.92
likato 296.06 340.97 44.91
mavala 409.04 446.32 37.28
vilenta 197.6 231.21 33.61
beautyblender 78.74 109.41 30.67
biore 60.65 90.31 29.66
orly 902.38 931.09 28.71
estelare 444.81 471.87 27.06
profepil 93.36 118.02 24.66
blixz 38.95 63.4 24.45
binacil 0.0 24.26 24.26
godefroy 401.22 425.12 23.9
glysolid 69.73 91.59 21.86
veraclara 50.11 71.21 21.1
juno 0.0 21.08 21.08
kamill 63.01 81.49 18.48
treaclemoon 163.37 181.49 18.12
supertan 50.37 66.51 16.14
barbie 0.0 12.39 12.39
deoproce 316.84 329.17 12.33
rasyan 18.8 28.94 10.14
fly 17.14 27.17 10.03
tertio 236.16 245.8 9.64
jaguar 1102.11 1110.65 8.54
soleo 204.2 212.53 8.33
neoleor 43.41 51.7 8.29
moyou 5.71 10.28 4.57
bodyton 1376.34 1380.64 4.3
skinity 8.88 12.44 3.56
helloganic 0.0 3.1 3.1
grace 100.92 102.61 1.69
cosima 20.23 20.93 0.7
ovale 2.54 3.1 0.56
Time taken: 45.29 seconds, Fetched: 161 row(s)

```

- From the output, we can observe that 161 brands increased their sales from the month of October to November.
- ‘Grattol’ brand has the highest per month increase of Rs. 36027.17
- ‘Ovale’ brand has the lowest per month increase of Rs. 0.56

viii. 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.

Command: *SELECT user_id, ROUND(SUM(price),2) AS spent_money*
>FROM cosmetics_par
>WHERE event_type= 'purchase'
>GROUP BY user_id
>ORDER BY spent_money DESC
>LIMIT 10;

Output:

```
hive> SELECT user_id, ROUND(SUM(price),2) AS spent_money
> FROM cosmetics_par
> WHERE event_type='purchase'
> GROUP BY user_id
> ORDER BY spent_money DESC
> LIMIT 10;
Query ID = hadoop_20221201173522_9d4770b6-1149-458a-a367-16c45b53c2f3
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 27.07 s
-----
OK
user_id spent_money
557790271      2715.87
150318419      1645.97
562167663      1352.85
531900924      1329.45
557850743      1295.48
522130011      1185.39
561592095      1109.7
431950134      1097.59
566576008      1056.36
521347209      1040.91
Time taken: 27.757 seconds, Fetched: 10 row(s)
```

The top 10 users of the website are as follows:

- 557790271 : 2715.87
- 150318419 : 1645.97
- 562167663 : 1352.85
- 531900924 : 1329.45

- 557850743 : 1295.48
- 522130011 : 1185.39
- 561592095 : 1109.7
- 431950134 : 1097.59
- 566576008 : 1056.36
- 521347209 : 1040.91

6. Improvement of performance after using optimization on tables

The execution time of query are reduced by using Hive optimization techniques. We had created an optimized table “cosmetics_par” through partitioning and bucketing in the start of the case study.

We will now observe the importance of optimization techniques, by querying two questions using the table “sales_details” and “cosmetics_par”. We have selected Question 6 & Question 8 for the demonstration.

Q1. Write a query to yield the total sum of purchases per month in a single output.

Command: *SELECT Month(event_time) as Month, ROUND(SUM(price),2) as total_purchases*
>FROM sales_details
>WHERE event_type='purchase'
>GROUP BY month(event_time);

Output:

```
hive> SELECT month(event_time) AS month, ROUND(SUM(price),2) AS total_purchases
> FROM sales_details
> WHERE event_type='purchase'
> GROUP BY month(event_time);
Query ID = hadoop_20221201180819_e42f05ee-3fc6-48d6-b082-525e955a6a3a
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0012)

-----
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: 59.43 s
-----
OK
month  total_purchases
10     1211538.43
11     1531016.9
Time taken: 60.192 seconds, Fetched: 2 row(s)
hive>
```

Command: *SELECT Month(event_time) as Month, ROUND(SUM(price),2) as total_purchases*
>FROM cosmetics_par
>WHERE event_type='purchase'
>GROUP BY month(event_time); (Optimised querying)

Output:

```
hive> SELECT month(event_time) as month, ROUND(SUM(price),2) AS total_purchases
> FROM cosmetics_par
> WHERE event_type='purchase'
> GROUP BY month(event_time);
Query ID = hadoop_20221201171315_af38aabd-820a-4376-90fb-eb057953f474
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)

-----
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: 25.97 s
-----
OK
month  total_purchases
10     1211534.31
11     1531007.31
Time taken: 26.668 seconds, Fetched: 2 row(s)
```

- We can observe that the results are similar for both queries but the query execution time differs.
- We see the query execution time for non optimized query is 60.192 seconds. We observed that, the query execution time for optimized query is 26.668 seconds. This is a huge difference, as the data size increase the increase in time will be substantially higher.

Q2. 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.

Command: *SELECT user_id, ROUND(SUM(price),2) AS spent_money*
>FROM sales_details
>WHERE event_type= 'purchase'
>GROUP BY user_id
>ORDER BY spent_money DESC
>LIMIT 10;

Output:

```
hive> SELECT user_id, ROUND(SUM(price),2) AS spent_money
> FROM sales_details
> WHERE event_type='purchase'
> GROUP BY user_id
> ORDER BY spent_money DESC
> LIMIT 10;
Query ID = hadoop_20221201175852_2da3a69c-bd1a-4902-a346-4d77986518dc
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_1669909606668_0012)
```

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

```
VERTICES: 03/03  [=====>>] 100% ELAPSED TIME: 62.68 s
OK
user_id spent_money
557790271      2715.87
150318419      1645.97
562167663      1352.85
531900924      1329.45
557850743      1295.48
522130011      1185.39
561592095      1109.7
431950134      1097.59
566576008      1056.36
521347209      1040.91
Time taken: 71.187 seconds, Fetched: 10 row(s)
```


Command: *SELECT user_id, ROUND(SUM(price),2) AS spent_money*
>FROM cosmetics_par
>WHERE event_type= 'purchase'
>GROUP BY user_id
>ORDER BY spent_money DESC
>LIMIT 10; (Optimized query)

Output:

```
hive> SELECT user_id, ROUND(SUM(price),2) AS spent_money
> FROM cosmetics_par
> WHERE event_type='purchase'
> GROUP BY user_id
> ORDER BY spent_money DESC
> LIMIT 10;
Query ID = hadoop_20221201173522_9d4770b6-1149-458a-a367-16c45b53c2f3
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1669909606668_0010)
```

	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: 27.07 s
OK
user_id spent_money
557790271      2715.87
150318419      1645.97
562167663      1352.85
531900924      1329.45
557850743      1295.48
522130011      1185.39
561592095      1109.7
431950134      1097.59
566576008      1056.36
521347209      1040.91
Time taken: 27.757 seconds, Fetched: 10 row(s)
```

- We can observe that the results are similar for both queries but the query execution time differs.
- We see the query execution time for non optimized query is 71.187 seconds. We observed that, the query execution time for optimized query is 27.757 seconds.

This is a huge difference, as the data size increase the increase in time will be substantially higher.

- Thus, partitioning and bucketing plays a crucial role to reduce the query execution time.

7. Cleaning-up :

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

```
hive> show databases;
OK
default
hive_casestudy
Time taken: 0.03 seconds, Fetched: 2 row(s)
hive> drop database hive_casestudy CASCADE;
OK
Time taken: 0.356 seconds
hive> show databases;
OK
default
Time taken: 0.017 seconds, Fetched: 1 row(s)
```

