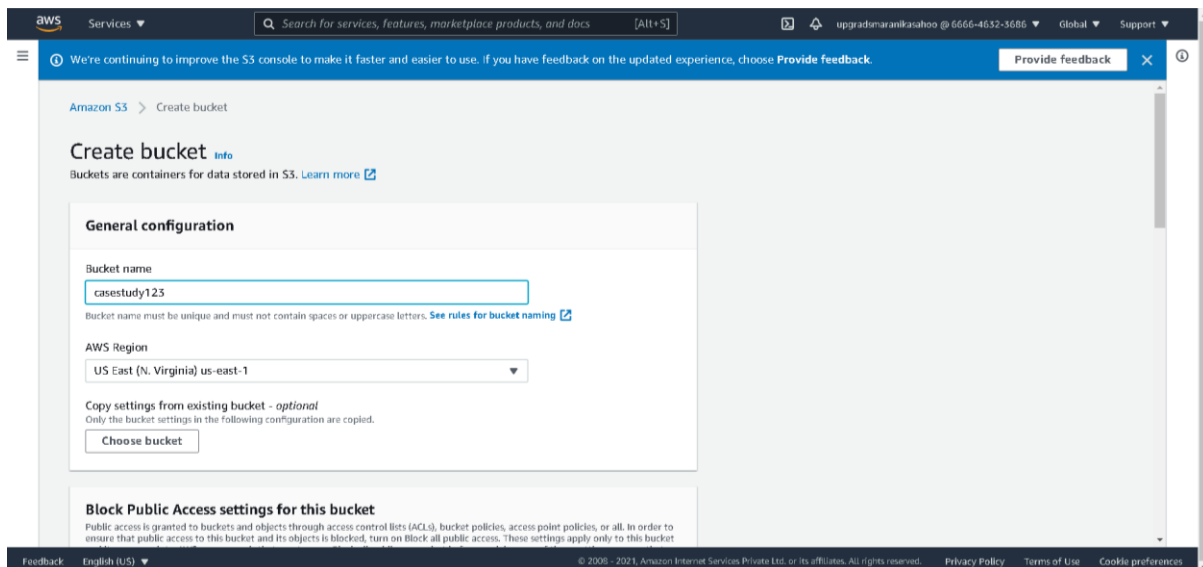


Hive case study- DA track

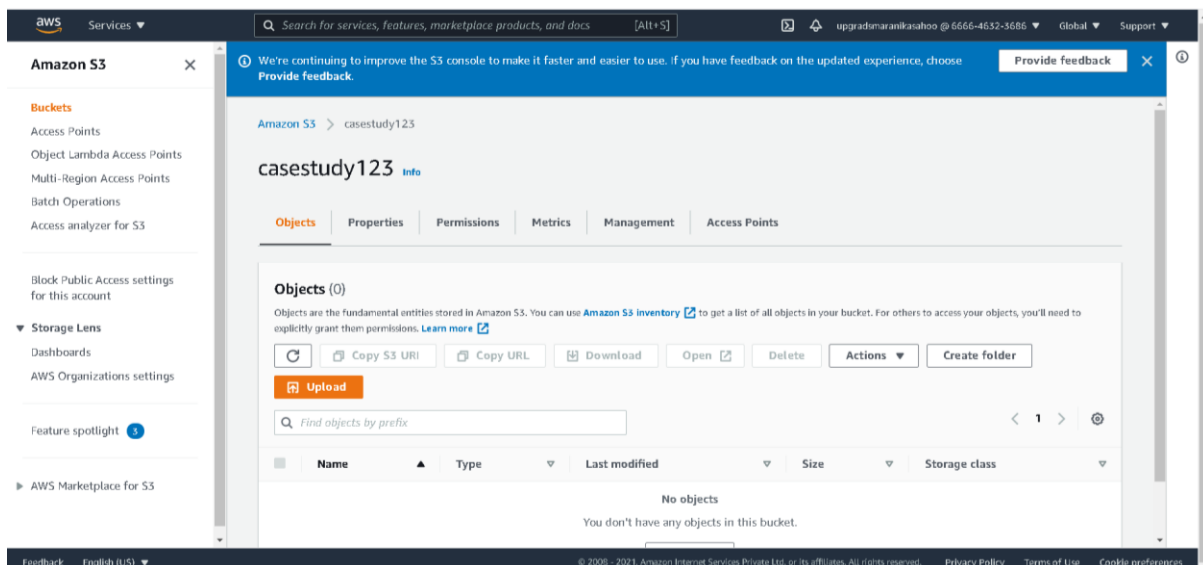
Submission By : Nitish Rathore & Smaranika Sahoo

S3 BUCKET:

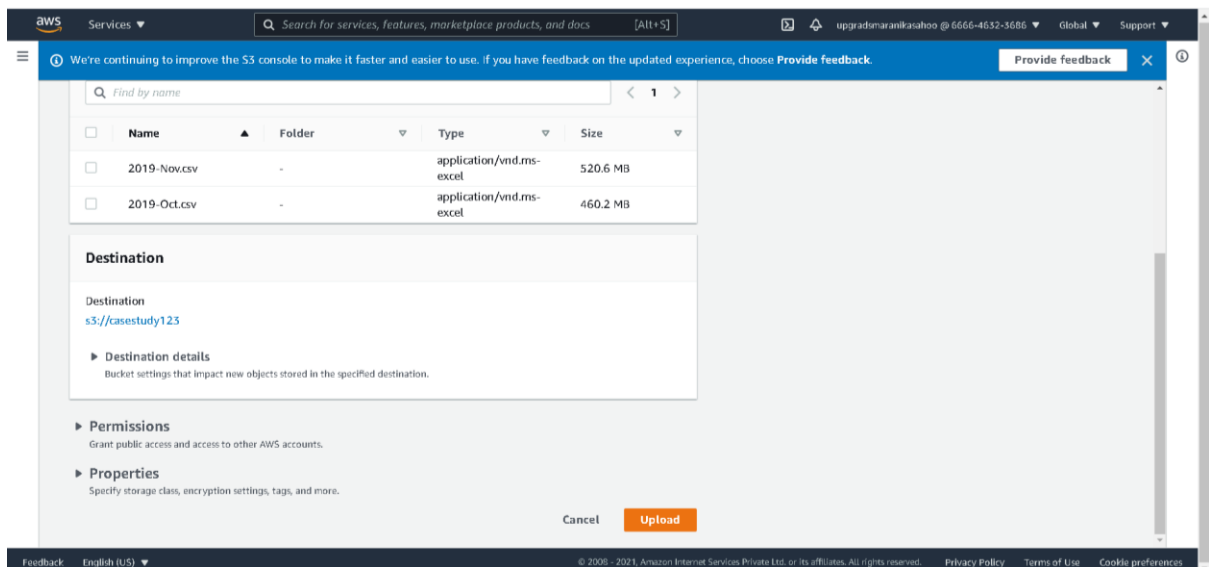
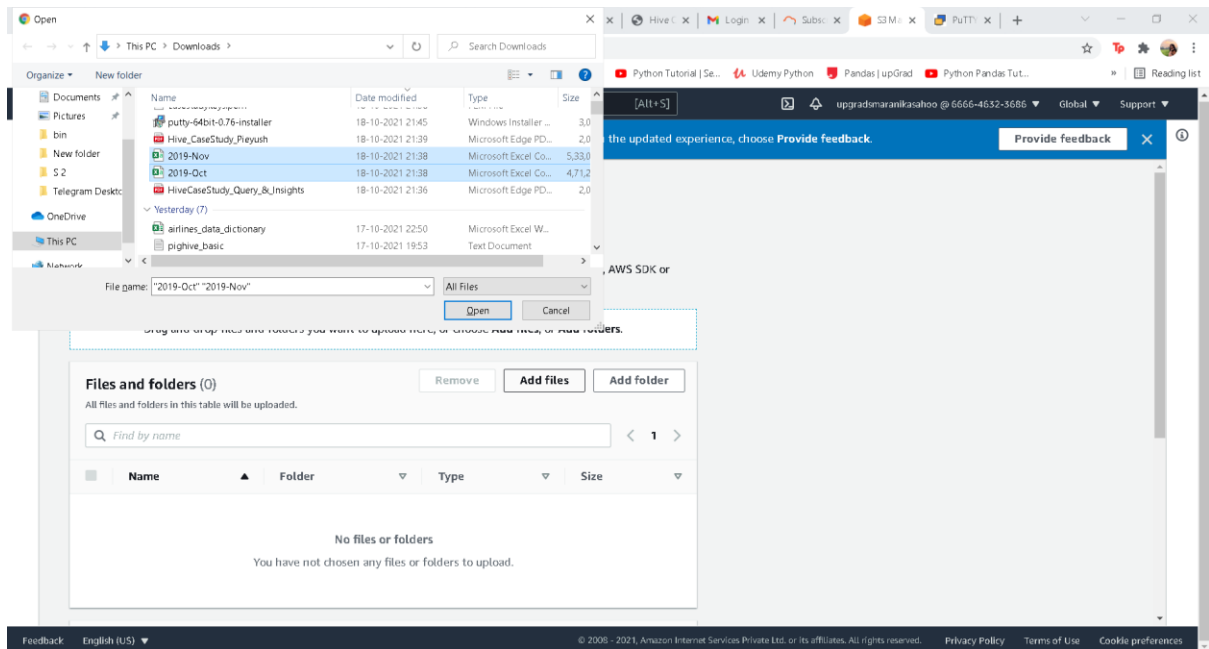
- To Store the data – Click on “Create Bucket”.

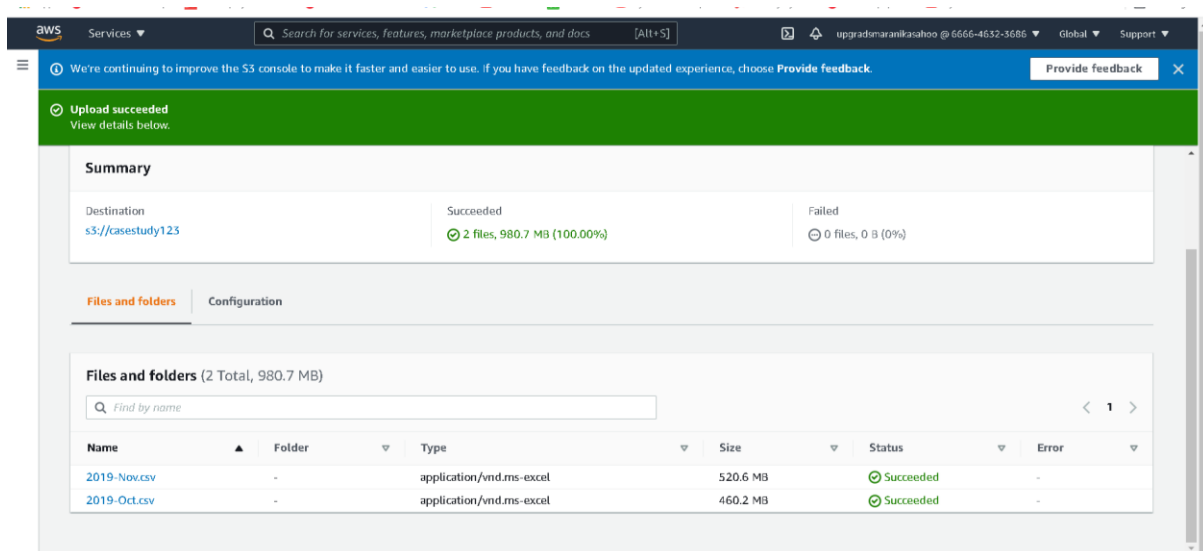


- Creating “casestudy123” with all default options.



- Bucket Successfully got created.

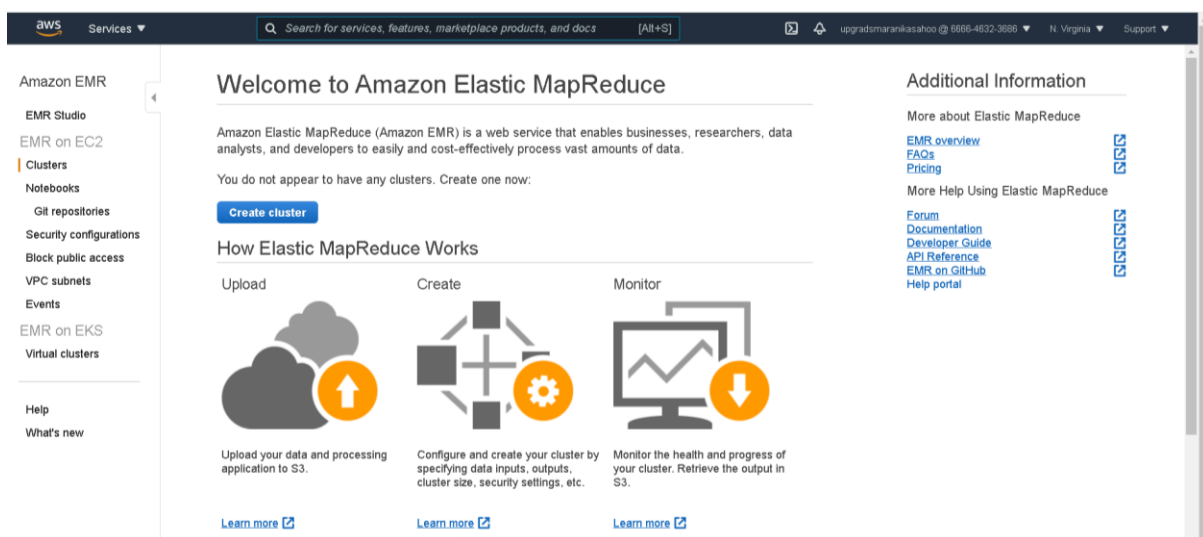




- Successfully uploaded the 2019 October and 2019 November csv file to S3 bucket.

EMR CLUSTER CREATION:

Click on “Create cluster” button to create the EMR cluster.



- Creating cluster with advanced options.

aws Services Search for services, features, marketplace products, and docs [Alt+S] upgradamaranikashoo@6866-4632-3886 N. Virginia Support

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps

Step 2: Hardware

Step 3: General Cluster Settings

Step 4: Security

Software Configuration

Release **emr-5.33.1**

<input checked="" type="checkbox"/> Hadoop 2.10.1	<input type="checkbox"/> Zeppelin 0.9.0	<input type="checkbox"/> Livy 0.7.0
<input type="checkbox"/> JupyterHub 1.1.0	<input type="checkbox"/> Tez 0.9.2	<input type="checkbox"/> Flink 1.12.1
<input type="checkbox"/> Ganglia 3.7.2	<input type="checkbox"/> HBase 1.4.13	<input checked="" type="checkbox"/> Pig 0.17.0
<input checked="" type="checkbox"/> Hive 2.3.7	<input type="checkbox"/> Presto 0.245.1	<input type="checkbox"/> ZooKeeper 3.4.14
<input type="checkbox"/> JupyterEnterpriseGateway 2.1.0	<input type="checkbox"/> MXNet 1.7.0	<input type="checkbox"/> Sqoop 1.4.7
<input type="checkbox"/> Mahout 0.13.0	<input checked="" type="checkbox"/> Hue 4.9.0	<input type="checkbox"/> Phoenix 4.14.3
<input type="checkbox"/> Oozie 5.2.0	<input type="checkbox"/> Spark 2.4.7	<input type="checkbox"/> HCatalog 2.3.7
<input type="checkbox"/> TensorFlow 2.4.1		

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

`classification=conf:file-name,properties=[myKey1=myValue1,myKey2=myValue2]`

- Hardware page: Changing the Master and Core nodes from m5.xlarge to “m4.large” with 1 instance each to make 2-node EMR cluster.

aws Services Search for services, features, marketplace products, and docs [Alt+S] upgradamaranikashoo@6866-4632-3886 N. Virginia Support

EC2 Subnet **subnet-06fc05d385c6e605** | Default in us-east-1a

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
Task Task - 3	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	<input type="text" value="0"/> Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

[+ Add task instance group](#)

Total core and task units: 1 Total units

Feedback English (US) © 2009 - 2021 Amazon Internet Services Private Ltd. or its affiliates. All rights reserved. Privacy Policy Terms of Use Cookie preferences

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps
Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

General Options

Cluster name:

☒ Logging ⓘ
S3 folder:

☐ Log encryption ⓘ
☒ 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: ⓘ

▶ Bootstrap Actions

- Security page: changing the EC2 key pair option to our created key pair – “CaseStudykeys”.

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps
Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

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

▶ EC2 security groups

[Cancel](#) [Previous](#) [Create cluster](#)

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: CaseStudycluster

Starting

Summary

Application user interfaces

Monitoring

Hardware

Configurations

Events

Steps

Bootstrap actions

Summary

Configuration details

Application user interfaces

Network and hardware

Security and access

ID: j-10Q2XLR2OXZ4Q

Creation date: 2021-10-18 22:10 (UTC+5:30)

Elapsed time: 1 second

After last step completes: Cluster waits

Termination protection: On [Change](#)

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

Master public DNS: --

Release label: emr-5.33.1

Hadoop distribution: Amazon 2.10.1

Applications: Hive 2.3.7, Pig 0.17.0, Hue 4.9.0

Log URI: s3://aws-logs-666646323686-us-east-1/elasticmapreduce/

EMRFS consistent view: Disabled

Custom AMI ID: --

Persistent user -- interfaces [\[+\]](#)

On-cluster user -- interfaces [\[+\]](#)

Availability zone: --

Subnet ID: subnet-06fc05d385c6e605 [\[+\]](#)

Master: Provisioning 1 m4.large

Core: Provisioning 1 m4.large

Task: --

Cluster scaling: Not enabled

Auto-termination: Not enabled

Key name: casestudykeys

EC2 instance profile: EMR_EC2_DefaultRole

EMR role: EMR_DefaultRole

Feedback

English (US)

© 2008 - 2021, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved.

Privacy Policy

Terms of Use

Cookie preferences

New EC2 Experience

EC2 Dashboard

EC2 Global View

Events

Tags

Limits

Instances

Instances **Now**

Instance Types

Launch Templates

Spot Requests

Savings Plans

Reserved Instances **Now**

Dedicated Hosts

Scheduled Instances

Capacity Reservations

Images

AMIs

Security Groups (1/2)

Info

Actions

Export security groups to CSV

Create security group

Filter security groups

search: sg-0fe49ab0968bddace

Clear filters

	Name	Security group ID	Security group name	VPC ID	Description	Owner
<input type="checkbox"/>	--	sg-0cbb2325b3b9c356c	ElasticMapReduce-slave	vpc-043e38e51f0b8ce95...	Slave group for Elastic...	66664632
<input checked="" type="checkbox"/>	--	sg-0fe49ab0968bddace	ElasticMapReduce-master	vpc-043e38e51f0b8ce95...	Master group for Elast...	66664632

sg-0fe49ab0968bddace - ElasticMapReduce-master

Details

Inbound rules

Outbound rules

Tags

You can now check network connectivity with Reachability Analyzer

Run Reachability Analyzer

Inbound rules (18)

Manage tags

Edit inbound rules

Feedback

English (US)

© 2008 - 2021, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved.

Privacy Policy

Terms of Use

Cookie preferences

Custom TCP

Custom TCP

Custom TCP

Custom TCP

SSH

207.171.167.25/32

54.240.217.16/29

sg-0cbb2325b3b9c356c

sg-0fe49ab0968bddace

54.240.217.8/29

0.0.0.0/0

8443

0 - 65535

0 - 65535

8443

22

Custom

Custom

Custom

Custom

Anywh...

Q

Q

Q

Q

Q

Delete

Delete

Delete

Delete

Delete

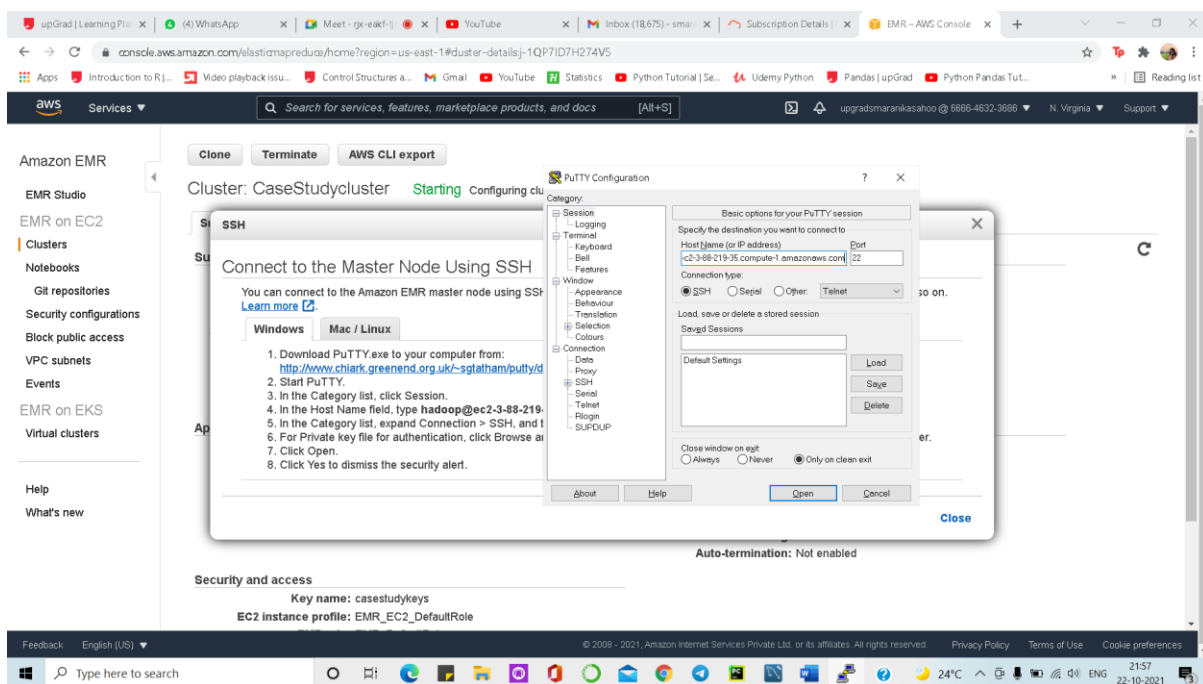
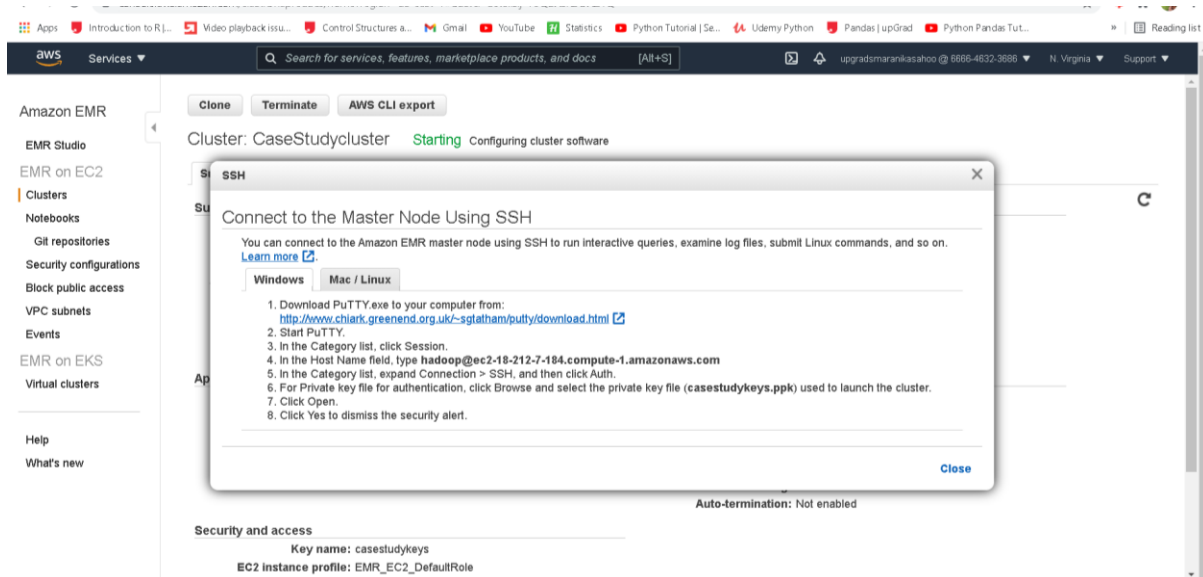
Add rule

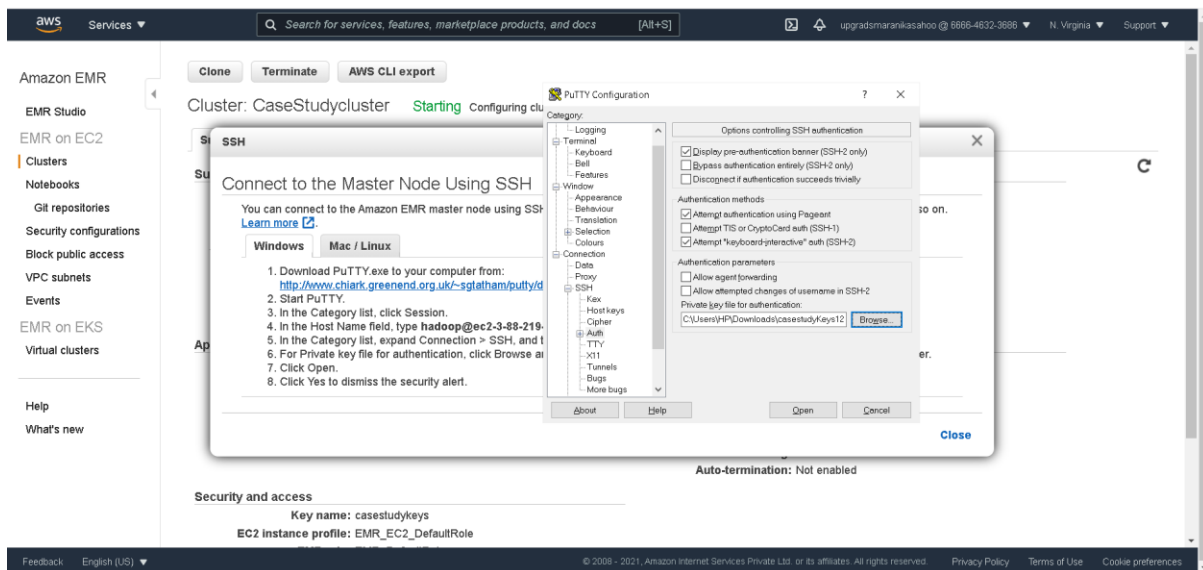
Cancel

Preview changes

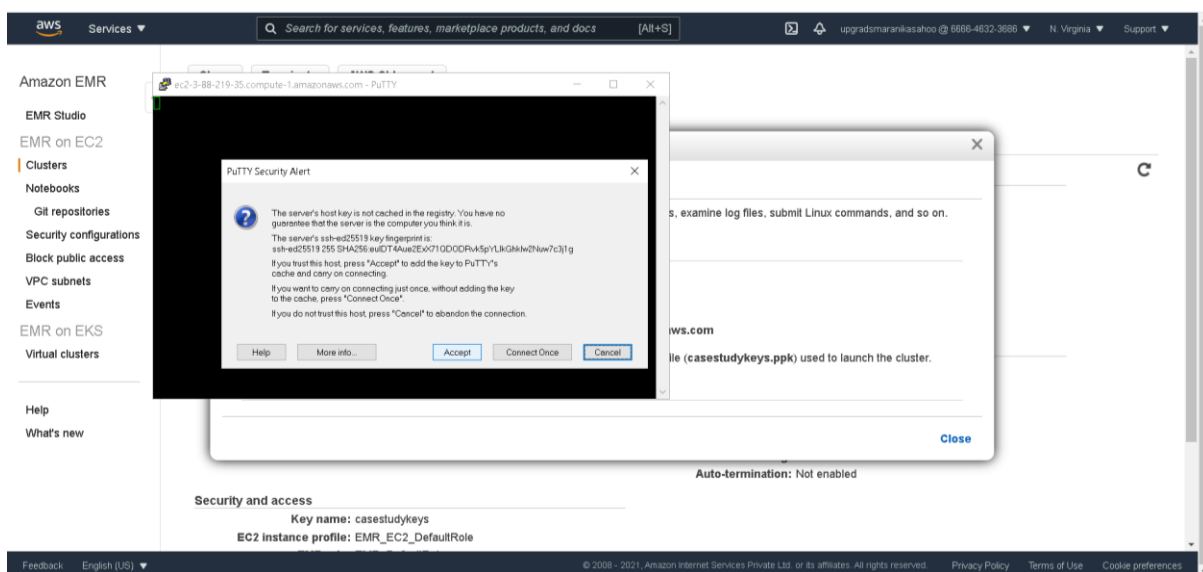
Save rules

- Then save the SSH rule to the inbound rules **CONNECT TO MASTER NODE**: Open the putty and enter the Host Name as `hadoop@ec2-18-212-7-184.compute-1.amazonaws.com` and navigate to Connection > SSH > Auth then browse and select the private key, which we creating initially.





- Click on “open” and then Accept the connection.




```
[hadoop@ip-172-31-81-254 ~]$ hadoop fs -mkdir /user/HiveCaseStudy/
[hadoop@ip-172-31-81-254 ~]$ hadoop fs -ls /user/
Found 7 items
drwxr-xr-x - hadoop hdfsadmingroup 0 2021-10-18 17:13 /user/HiveCaseStudy
drwxr-xr-x - hadoop hdfsadmingroup 0 2021-10-18 16:48 /user/hadoop
drwxr-xr-x - mapred mapred 0 2021-10-18 16:48 /user/history
drwxr-xr-x - hdfs hdfsadmingroup 0 2021-10-18 16:48 /user/hive
drwxr-xr-x - hue hue 0 2021-10-18 16:48 /user/hue
drwxr-xr-x - cozie cozie 0 2021-10-18 16:51 /user/cozie
drwxr-xr-x - root hdfsadmingroup 0 2021-10-18 16:48 /user/root
[hadoop@ip-172-31-81-254 ~]$
```

- New directory is successfully created.

LOADING THE DATA FROM S3 BUCKET to HDFS:

- Distributed copy command is used to copy the data from S3 to HDFS.
- For 2019 October:
“hadoop distcp s3://casestudy123/2019-Oct.csv /user/HiveCaseStudy/October.csv”
- For 2019 November:
“hadoop distcp s3://casestudy123/2019-Nov.csv /user/HiveCaseStudy/November.csv”
- Below are the screenshots for copying October 2019 and November 2019 data individually.

```
[hadoop@ip-172-31-81-254 ~]$ hadoop distcp s3://casestudy123/2019-Oct.csv /user/HiveCaseStudy/october.csv
21/10/18 17:18:27 INFO tools.CyclicProgress: parseChunkSize: blockPerChunk false
21/10/18 17:18:28 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, useDiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=100, s3ConfigurationFile='null', copyStrategy='unifor
size', preserveStatus=[], preserveRawKattis=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://casestudy123/2019-Oct.csv], targetPath=/user/HiveCaseStudy/o
ctober.csv, targetPathExists=false, filterFiles='null', blocksPerChunk=0, copyBufferSize=8192, verboseLog=false}
21/10/18 17:18:28 INFO client.RMProxy: Connecting to Resource Manager at ip-172-31-81-254.ec2.internal/172.31.81.254:8032
21/10/18 17:18:28 INFO client.AMRProxy: Connecting to Application History server at ip-172-31-81-254.ec2.internal/172.31.81.254:10200
21/10/18 17:18:32 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
21/10/18 17:18:32 INFO tools.SimpleCopyListing: Build file listing completed.
21/10/18 17:18:32 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
21/10/18 17:18:32 INFO tools.DistCp: Number of paths in the copy list: 1
21/10/18 17:18:33 INFO client.RMProxy: Connecting to Resource Manager at ip-172-31-81-254.ec2.internal/172.31.81.254:8032
21/10/18 17:18:33 INFO client.AMRProxy: Connecting to Application History server at ip-172-31-81-254.ec2.internal/172.31.81.254:10200
21/10/18 17:18:33 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1634575743345_0001
21/10/18 17:18:34 INFO conf.Configuration: resource-types.xml not found
21/10/18 17:18:34 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/10/18 17:18:34 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/10/18 17:18:34 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/10/18 17:18:34 INFO impl.YarnClientImpl: Submitted application application_1634575743345_0001
21/10/18 17:18:34 INFO mapreduce.Job: The url to track the job: http://ip-172-31-81-254.ec2.internal:20888/proxy/application_1634575743345_0001/
21/10/18 17:18:34 INFO tools.DistCp: Distcp job-id: job_1634575743345_0001
21/10/18 17:18:34 INFO mapreduce.Job: Running job: job_1634575743345_0001
```

```
hadoop@ip-172-31-81-254-
21/10/18 17:18:44 INFO mapreduce.Job: Job job_1634575743345_0001 running in uber mode : false
21/10/18 17:18:44 INFO mapreduce.Job: map 0% reduce 0%
21/10/18 17:19:02 INFO mapreduce.Job: map 100% reduce 0%
21/10/18 17:19:06 INFO mapreduce.Job: Job job_1634575743345_0001 completed successfully
21/10/18 17:19:06 INFO mapreduce.Job: Counters: 38
File System Counters
FILE: Number of bytes read=0
FILE: Number of bytes written=223523
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=364
HDFS: Number of bytes written=482542278
HDFS: Number of read operations=11
HDFS: Number of large read operations=0
HDFS: Number of write operations=5
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)=606688
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=18959
Total wordcount-milliseconds taken by all map tasks=18959
Total megabyte-milliseconds taken by all map tasks=19414016
Map-Reduce Framework
Map input records=1
Map output records=0
Input split bytes=134
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=344
CPU time spent (ms)=20210
Physical memory (bytes) snapshot=629489664
Virtual memory (bytes) snapshot=3326230528
Total committed heap usage (bytes)=503840768
File Input Format Counters
Bytes Read=230
File Output Format Counters
Bytes Written=0
DistCp Counters
Bytes Copied=482542278
Bytes Expected=482542278
Files Copied=1
[hadoop@ip-172-31-81-254 ~]$
```

```

[hadoop@ip-172-31-81-254 ~]$
[hadoop@ip-172-31-81-254 ~]$ hadoop distcp s3:///casestudy123/2019-Nov.csv /user/HiveCaseStudy/November.csv
21/10/18 17:21:08 INFO tools.OptionsParser: parseChunkSize: blockPerChunk false
21/10/18 17:21:09 INFO tools.DistCp: Input Options: DistCpOptions(atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, append=false, useDiff=false, useRdiff=false, fromSnapshot=null, toSnapshot=null, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=100, sslConfigurationFile='null', copyStrategy='uniformize', preserveStatus=[], preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3:///casestudy123/2019-Nov.csv], targetPath=/user/HiveCaseStudy/November.csv, targetPathExists=false, filtersFile='null', blocksPerChunk=0, copyBufferSize=0192, verboseCopy=false)
21/10/18 17:21:10 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-81-254.ec2.internal/172.31.81.254:8032
21/10/18 17:21:10 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-81-254.ec2.internal/172.31.81.254:10200
21/10/18 17:21:14 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
21/10/18 17:21:14 INFO tools.SimpleCopyListing: Build file listing completed.
21/10/18 17:21:14 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
21/10/18 17:21:14 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
21/10/18 17:21:14 INFO tools.DistCp: Number of paths in the copy list: 1
21/10/18 17:21:14 INFO tools.DistCp: Number of paths in the copy list: 1
21/10/18 17:21:14 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-81-254.ec2.internal/172.31.81.254:8032
21/10/18 17:21:14 INFO client.AHSProxy: Connecting to Application History server at ip-172-31-81-254.ec2.internal/172.31.81.254:10200
21/10/18 17:21:15 INFO mapreduce.JobSubmitter: number of splits:1
21/10/18 17:21:15 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1634575743345_0002
21/10/18 17:21:15 INFO conf.Configuration: resource-types.xml not found
21/10/18 17:21:15 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
21/10/18 17:21:15 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
21/10/18 17:21:15 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
21/10/18 17:21:15 INFO impl.YarnClientImpl: Submitted application application_1634575743345_0002
21/10/18 17:21:16 INFO mapreduce.Job: The url to track the job: http://ip-172-31-81-254.ec2.internal:20808/proxy/application_1634575743345_0002/
21/10/18 17:21:16 INFO tools.DistCp: DistCp job-id: job_1634575743345_0002
21/10/18 17:21:16 INFO mapreduce.Job: Running job: job_1634575743345_0002

```

```

hadoop@ip-172-31-81-254~$
21/10/18 17:21:24 INFO mapreduce.Job: Job job_1634575743345_0002 running in uber mode : false
21/10/18 17:21:24 INFO mapreduce.Job: map 0% reduce 0%
21/10/18 17:21:41 INFO mapreduce.Job: map 100% reduce 0%
21/10/18 17:21:44 INFO mapreduce.Job: Job job_1634575743345_0002 completed successfully
21/10/18 17:21:44 INFO mapreduce.Job: Counters: 38
File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=222534
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=367
  HDFS: Number of bytes written=545839412
  HDFS: Number of read operations=11
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=5
  S3: Number of bytes read=545839412
  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)=584000
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=58250
  Total vcore-milliseconds taken by all map tasks=18250
  Total megabyte-milliseconds taken by all map tasks=18688000
Map-Reduce Framework
  Map input records=1
  Map output records=0
  Input split bytes=137
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=254
  CPU time spent (ms)=22260
  Physical memory (bytes) snapshot=617050112
  Virtual memory (bytes) snapshot=3330248704
  Total committed heap usage (bytes)=479199232
File Input Format Counters
  Bytes Read=230
File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=545839412
  Bytes Expected=545839412
  Files Copied=1
hadoop@ip-172-31-81-254 ~$

```

- Verifying whether the data is successfully copied into HDFS from S3 buckets Command:
hadoop fs -ls /user/ HiveCaseStudy

```

Files Copied=1
[hadoop@ip-172-31-80-138 ~]$ hadoop fs -ls /user/HiveCaseStudy/
Found 2 items
-rw-r--r-- 1 hadoop hdfsadmin group 545839412 2021-10-19 16:35 /user/HiveCaseStudy/November.csv
-rw-r--r-- 1 hadoop hdfsadmin group 482542278 2021-10-19 16:34 /user/HiveCaseStudy/October.csv

```

Moving to hive:

```

[hadoop@ip-172-31-80-138 ~]$ hadoop fs -ls /user/HiveCaseStudy/October.csv
[hadoop@ip-172-31-80-138 ~]$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> CREATE EXTERNAL TABLE IF NOT EXISTS retailstore (event_time timestamp, event_type string,
  > product_id string, category_id string, category_code string, brand string, price float, user_id bigint,
  > user_session string) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS
  > TEXTFILE LOCATION '/user/HiveCaseStudy/' tblproperties ("skip.header.line.count"="1");
OK
Time taken: 1.039 seconds
hive> show tables ;
OK
retailstore
Time taken: 0.242 seconds, Fetched: 1 row(s)
hive>

```

Creating RetailDB Database:

```
hive> create database if not exists RetailDB;
OK
Time taken: 0.566 seconds
hive> describe database RetailDB ;
OK
retaildb                hdfs://ip-172-31-91-69.ec2.internal:8020/user/hive/warehouse/retaildb.db      hadoop  USER
Time taken: 0.203 seconds, Fetched: 1 row(s)
```

```
hive> use RetailDB;
OK
Time taken: 0.035 seconds
```

CREATING AN EXTERNAL TABLE IN HIVE:

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

```
[hadoop@ip-172-31-80-130 ~]$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> CREATE EXTERNAL TABLE IF NOT EXISTS retailstore (event_time timestamp, event_type string,
  > product_id string, category_id string, category_code string, brand string, price float, user_id bigint,
  > user_session string) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS
  > TEXTFILE LOCATION '/user/HiveCaseStudy/' tblproperties("skip.header.line.count"="1");
OK
Time taken: 1.039 seconds
hive> show tables ;
OK
retailstore
Time taken: 0.242 seconds, Fetched: 1 row(s)
hive> set hive.cli.print.header=True;
hive> select * from retailstore limit 3 ;
OK
retailstore.event_time retailstore.event_type retailstore.product_id retailstore.category_id      retailstore.category_code      retailstore.brand      retailstore.price      retail
istore.user_id retailstore.user_session
2019-11-01 00:00:02 UTC view      5802432 1487580009286598681      0.32      562076640      09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:09 UTC cart      5844397 1487580006317032337      2.38      553329724      2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:10 UTC view      5837166 1783999064103190764      pnb      22.22      556138645      57ed222e-a54a-4907-9944-5a875c2d7f4f
Time taken: 2.502 seconds, Fetched: 3 row(s)
hive>
```

- set hive.cli.print.header=True;
- Checked the data in the table by querying the below.

```
hive> set hive.cli.print.header=True;
hive> select * from retailstore limit 3 ;
OK
retailstore.event_time retailstore.event_type retailstore.product_id retailstore.category_id      retailstore.category_code      retailstore.brand      retailstore.price      retail
istore.user_id retailstore.user_session
2019-11-01 00:00:02 UTC view      5802432 1487580009286598681      0.32      562076640      09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:09 UTC cart      5844397 1487580006317032337      2.38      553329724      2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:10 UTC view      5837166 1783999064103190764      pnb      22.22      556138645      57ed222e-a54a-4907-9944-5a875c2d7f4f
Time taken: 0.207 seconds, Fetched: 3 row(s)
```

Questions & Answers

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

Query: SELECT SUM(price) AS total_revenue_October FROM retailstore WHERE date_format(event_time,'MM')=10 AND event_type='purchase';

```

Time taken: 0.132 seconds, Fetched: 1 row(s)
hive> SELECT SUM(price) AS total_revenue_October FROM retailstore WHERE date_format(event_time,'MM')=10 AND event_type='purchase';
Query ID = hadoop_20211019165603_10765737-92bc-45b7-9335-79e30cd447e8
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634660723820_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: 126.79 s
-----
OK
1211538.4299997438
Time taken: 132.79 seconds, Fetched: 1 row(s)

```

Insight:

The total revenue generated based on Purchase in the month of October of 2019 was 1,211,538.42999/

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

Query: SELECT date_format(event_time, 'MM') AS Months, COUNT(event_type) AS Sum_of_Purchases FROM retailstore WHERE event_type='purchase' GROUP BY date_format(event_time, 'MM');

```

hive> SELECT date_format(event_time, 'MM') AS Months, COUNT(event_type) AS Sum_of_Purchases FROM retailstore WHERE event_type='purchase' GROUP BY date_format(event_time, 'MM');
Query ID = hadoop_20211019170526_58e26359-7b4b-4438-ab0e-2f178c76fc82
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634660723820_0005)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    6         6         0         0         0         0
-----
VERTICES: 02/02 [=====>>>] 100% ELAPSED TIME: 62.31 s
-----
OK
10      245624
11     322417
Time taken: 66.994 seconds, Fetched: 2 row(s)

```

Insight:

The total sum of purchases per month, October (10) i.e., 245,624, November (11) i.e., 322,417

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

Query: with diff_revenue as(select sum(case when MONTH(event_time) ='10' then price else 0 end) as oct_revenue,sum(case when MONTH(event_time) ='11' then price else 0 end) as nov_revenue from retailstore where event_type='purchase') select (nov_revenue - oct_revenue) as revenue_difference from diff_revenue;

```

hive> with diff_revenue as(select sum(case when MONTH(event_time) ='10' then price else 0 end) as oct_revenue,
> sum(case when MONTH(event_time) ='11' then price else 0 end) as nov_revenue
> from retailstore where event_type='purchase')
> select (nov_revenue - oct_revenue) as revenue_difference from diff_revenue;
Query ID = hadoop_20211024130838_fa0c4c6c-4714-4576-8b09-1f5675626936
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1635075264890_0006)

-----
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: 56.73 s
-----
OK
319478.4700003781
Time taken: 61.75 seconds, Fetched: 1 row(s)

```

Insight:

From the query we got the difference in revenue from October to November is 319478.47

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

Query: select distinct(category_code) from retailstore where category_code != " ";

```

hive> select distinct(category_code)from retailstore where category_code != ' ' ;
Query ID = hadoop_20211019172124_d9695756-ab8a-46eb-9fe5-180045826ff9
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634660723820_0007)

-----
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: 55.86 s
-----
OK
accessories.bag
accessories.cosmetic_bag
apparel.glove
appliances.environment.air_conditioner
appliances.environment.vacuum
appliances.personal.hair_cutter
furniture.bathroom.bath
furniture.living_room.cabinet
furniture.living_room.chair
sport.diving
stationery.cartridge
Time taken: 61.594 seconds, Fetched: 11 row(s)

```

Insight:

- Total we got 6 distinct categories are – furniture, appliances, accessories, apparel, sport, stationary.

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

Query: select count(product_id) , category_code from retailstore where category_code IS NOT NULL group by category_code;

```
hive> select count(product id) , category_code from retailstore where category_code IS NOT NULL group by category_code;
Query ID = hadoop_20211019172400_837299f4-bald-453e-ba84-95d85564a7e6
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634660723820_0007)

-----
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: 51.61 s
-----
OK
8594895
11681  accessories.bag
1248   accessories.cosmetic_bag
18232  apparel.glove
332    appliances.environment.air_conditioner
59761  appliances.environment.vacuum
1643   appliances.personal.hair_cutter
9857   furniture.bathroom.bath
13439  furniture.living_room.cabinet
308    furniture.living_room.chair
2      sport.diving
26722  stationery.cartridge
Time taken: 52.376 seconds, Fetched: 12 row(s)
```

Insight:

- Company has more products registered under Appliances category i.e., 61,736 products than any other categories.
- Then it is followed by stationery as second with 26,722 products, furniture as third with 23,604 products, apparel as fourth with 18232 products registered, accessories as fifth with 12929 products
- Sports category is least available with 2 products.

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

Query: select brand ,sum(price) as sales from retailstore where brand != "" and event_type= "purchase" group by brand order by sales desc limit 1;

```
hive> select brand ,sum(price) as sales from retailstore where brand != "" and event_type= "purchase" group by brand order by sales desc limit 1;
Query ID = hadoop_20211024131531_c4a5a511-39da-429d-8112-ddc5f782ed53
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_1635075264890_0007)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    6         6         0         0         0         0
Reducer 3 ..... container  SUCCEEDED    1         1         0         0         0         0
-----
VERTICES: 03/03  [=====>>>] 100%  ELAPSED TIME: 59.49 s
-----
OK
runail 148297.94000000003
Time taken: 68.936 seconds, Fetched: 1 row(s)
```

Insight:

From the query we got that runail had the maximum sales in October and November combined.

Question 7: Which brands increased their sales from October to November?

Query: WITH Monthly_Revenue AS (SELECT brand, SUM(CASE WHEN date_format(event_time, 'MM')=10 THEN price ELSE 0 END) AS Oct_Revenue, SUM(CASE WHEN date_format(event_time, 'MM')=11 THEN price ELSE 0 END) AS Nov_Revenue FROM retailstore WHERE event_type='purchase' GROUP BY brand) SELECT brand, Oct_Revenue, Nov_Revenue, Nov_Revenue-Oct_Revenue AS Sales_Difference FROM Monthly_Revenue WHERE (Nov_Revenue - Oct_Revenue)>0 ORDER BY Sales_Difference ;

```
hive>
+ WITH Monthly_Revenue AS ( SELECT brand, SUM(CASE WHEN date_format(event_time, 'MM')=10 THEN price ELSE 0 END) AS Oct_Revenue, SUM(CASE WHEN date_format(event_time, 'MM')=11 THEN price
ELSE 0 END) AS Nov_Revenue FROM retailstore WHERE event_type='purchase' GROUP BY brand ) SELECT brand, Oct_Revenue, Nov_Revenue, Nov_Revenue-Oct_Revenue AS Sales_Difference FROM Monthly_R
venue WHERE (Nov_Revenue - Oct_Revenue)>0 ORDER BY Sales_Difference ;
Query ID = hadoop_20211024131813_4cd99144-ebaa-4111-9cd5-373681e798a1
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1635075264890_0007)

-----
VERTICES      MODE          STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    6         6         0         0         0         0
Reducer 3 ..... container  SUCCEEDED    1         1         0         0         0         0
-----
VERTICES: 03/03 [=====]>>] 100% ELAPSED TIME: 70.52 s
-----
OK
ovale 2.54 3.1 0.56
cosima 20.23 20.929999999999993 0.69999999999999922
piace 100.920000000000002 102.61000000000001 1.6899999999999977
halloganic 0.0 3.1 3.1
skinity 8.08 12.440000000000001 3.5600000000000005
bodyton 1376.3399999999974 1380.6399999999992 4.3000000000017735
mayou 5.71 10.280000000000001 4.570000000000001
necleor 42.41 51.7 9.290000000000006
salo 204.20000000000003 212.52999999999998 8.329999999999501
jaguar 1102.11 1110.6500000000003 8.5400000000000418
tertio 236.16000000000008 245.79999999999978 9.639999999999702
rhy 17.14 27.17 10.030000000000001
rasyan 10.799999999999997 28.839999999999994 10.139999999999997
deoproce 316.84 329.17000000000001 12.330000000000098
barbie 0.0 12.39 12.39
supertan 50.370000000000001 66.510000000000002 16.140000000000008
treaclemoon 163.36999999999995 191.48999999999995 18.120000000000005
kamill 63.009999999999999 81.490000000000002 18.480000000000032
juno 0.0 21.08 21.08
veraclarla 50.109999999999985 71.210000000000001 21.1000000000000023
glysolid 69.72999999999998 91.58999999999997 21.86
podeftoy 401.22000000000002 425.12000000000006 23.899999999999864
binacil 0.0 24.259999999999998 24.259999999999998
blixr 38.949999999999996 63.399999999999998 24.449999999999998
profepil 93.360000000000003 118.02000000000005 24.660000000000025
estelare 444.80999999999943 471.8700000000009 27.060000000000148
cely 902.38000000000005 931.09000000000003 28.709999999999981
biore 60.650000000000006 90.31 29.659999999999997
beautyblender 78.740000000000001 109.41 30.669999999999987
vilenta 197.60000000000002 231.21000000000002 33.610000000000014
navaia 409.03999999999985 446.32 37.280000000000014
```


mavala	409.0399999999985	446.32	37.28000000000014	
likato	296.0599999999999	340.9699999999999	44.910000000000025	
ladykin	125.6499999999999	170.57	44.92	
foamie	35.04	80.49	45.44999999999996	
elskin	251.090000000000057	307.650000000000055	56.55999999999974	
balbcare	155.32999999999996	212.380000000000025	57.050000000000296	
koelcia	55.5	112.75000000000003	57.25000000000003	
profhenna	679.2299999999999	736.85000000000005	57.620000000000057	
kares	0.0	59.45	59.45	
marutaka-foot	49.21999999999999	109.33	60.110000000000001	
dewal	0.0	61.29	61.29	
inm	288.02	351.21000000000001	63.190000000000011	
laboratorium	246.4999999999991	312.52	66.020000000000007	
cutrin	299.3699999999995	367.62	68.25000000000006	
egomania	77.47	146.040000000000002	68.570000000000002	
konad	739.8299999999991	810.67000000000003	70.840000000000117	
nirvel	163.03999999999996	234.32999999999984	71.28999999999988	
koelf	422.72999999999985	507.29000000000002	84.560000000000034	
plazan	101.37	194.010000000000002	92.640000000000001	
aura	83.95	177.51	93.55999999999999	
kerasys	430.90999999999985	525.20000000000002	94.29000000000003	
enjoy	41.34999999999994	136.570000000000002	95.220000000000003	
depilflax	2707.069999999994	2803.7799999999975	96.710000000000367	
eos	54.33999999999996	152.61	98.270000000000001	
carmex	145.08	243.36	98.28	
batiste	772.3999999999999	874.1699999999994	101.76999999999953	
osmo	645.58	762.31000000000002	116.730000000000013	
dizao	819.13000000000012	945.5099999999998	126.37999999999982	
igrobeauty	513.66000000000009	645.0699999999999	131.40999999999906	
finish	98.38	230.38000000000008	132.00000000000009	
nefertiti	233.52000000000007	366.64	133.11999999999992	
elizavecca	70.53	204.3	133.77	
maskin	158.04	293.07000000000005	135.03000000000006	
latinoil	249.52	384.59	135.06999999999996	
farmona	1692.4599999999996	1843.43000000000007	150.970000000000116	
crystalinas	427.6299999999999	584.9499999999999	157.31999999999914	
chi	358.94000000000002	538.61000000000002	179.67000000000002	
matreshka	0.0	182.67000000000002	182.67000000000002	
freshbubble	318.70000000000001	502.340000000000015	183.64000000000004	
mane	66.78999999999999	260.26	193.47	
keen	236.350000000000005	435.62	199.26999999999995	
ecocraft	41.160000000000004	241.95	200.79	
fedua	52.38	263.81000000000006	211.43000000000006	
provoc	827.99000000000009	1063.82000000000006	235.8299999999997	
skinlite	651.94000000000002	890.4499999999979	238.50999999999772	
entity	479.71000000000015	719.2599999999993	239.5499999999978	
trind	298.070000000000005	542.96000000000002	244.89000000000001	
protokeratin	201.25000000000003	456.790000000000013	255.54000000000001	
beauugreen	511.5099999999999	768.35	256.840000000000015	
bluesky	10307.2399999999858	10565.5299999999713	258.28999999998535	

bluesky	10307.239999999858	10565.529999999713	258.28999999985535
candy	534.9599999999999	799.3799999999993	264.4199999999994
insight	1443.70000000000012	1721.96000000000003	278.25999999999991
kocostar	310.85000000000001	594.93000000000003	284.08000000000002
happyfons	801.92000000000006	1091.59000000000001	289.66999999999995
kims	330.03999999999996	632.04000000000001	302.00000000000001
shary	871.9599999999994	1176.4899999999989	304.52999999999995
nitrile	847.2799999999999	1162.6799999999999	315.4
lowence	242.84 567.7499999999997	324.90999999999996	
jas	3318.9599999999995	3657.43000000000026	338.470000000000753
ellips	245.84999999999999	606.03999999999996	360.18999999999997
lador	2083.6100000000004	2471.5300000000007	387.92000000000028
naomi	0.0 389.0 389.0		
kiss	421.54999999999944	817.3299999999994	395.77999999999999
yu-r	271.41 673.7099999999998	402.2999999999998	
sophin	1067.86000000000001	1515.52000000000011	447.66000000000001
farmavita	837.3699999999984	1291.97000000000003	454.600000000000184
bioaqua	942.8899999999996	1398.1199999999997	455.23
greymy	29.21 489.49 460.28000000000003		
gehwol	1089.07 1557.6799999999982	468.6099999999983	
matrix	3243.2499999999999	3726.7400000000007	483.49000000000016
limoni	1308.90000000000003	1796.5999999999997	487.69999999999936
s.care	412.68 913.0699999999999	500.38999999999993	
coifin	903.00000000000001	1428.4899999999998	525.4899999999997
uskusi	5142.2700000000017	5690.3100000000005	548.03999999999881
airnails	5118.8999999999939	5691.5199999999996	572.620000000000572
browxenna	14331.369999999995	14916.729999999976	585.36000000000026
kinetics	6334.2499999999945	6945.2600000000017	611.01000000000022
kosmekka	1181.44000000000003	1813.37 631.9299999999996	
kaaral	4412.4299999999985	5086.0699999999992	673.6399999999994
refectocil	2716.1800000000005	3475.5800000000007	759.40000000000024
rosi	3077.0399999999927	3841.5600000000013	764.520000000000204
solomeya	1899.6999999999992	2685.7999999999991	786.0999999999999
missha	1293.8299999999995	2150.2799999999984	856.4499999999989
levissime	2227.50000000000064	3085.3099999999977	857.80999999999913
art-visage	2092.7100000000001	2997.8000000000011	905.09000000000001
ecolab	262.85000000000001	1214.2999999999988	951.4499999999987
nagaraku	4369.7400000000054	5327.6800000000063	957.94000000000087
sanoto	157.14 1209.6799999999998	1052.54	
markell	1768.7499999999989	2834.4300000000007	1065.68000000000019
metzger	5373.4500000000006	6457.1599999999988	1083.70999999999818
de.lux	1659.6999999999967	2775.5099999999968	1115.81000000000009
swarovski	1887.92999999999873	3043.1600000000003	1155.230000000000157
beauty-free	554.17000000000006	1782.86000000000163	1228.690000000000155
zeitun	708.66000000000004	2009.63 1300.9699999999998	
joico	705.52 2015.10000000000015	1309.58000000000015	
severina	4775.88 6120.4800000000023	1344.6000000000023	

```

hive>
levislime      2221.50000000000064      3085.30999999999977      857.80999999999913
art-vilage     2092.710000000001      2997.8000000000011      905.0900000000001
esolab 262.850000000001      1214.2999999999998      951.4499999999997
nagataku      4369.7400000000054      5327.6800000000063      957.94000000000087
sanoto 157.14      1209.6799999999998      1052.54
mackell 1768.7499999999989      2834.4300000000007      1065.68000000000019
marzger 5373.4500000000006      6457.159999999998      1083.70999999999818
de.lux 1659.6999999999967      2775.5099999999968      1115.0100000000009
smarovski 1887.92999999999873      3043.1600000000003      1155.23000000000157
beauty-free 554.17000000000006      1782.86000000000163      1228.69000000000155
zeitun 708.6600000000004      2009.63      1300.9699999999998
joico 705.52      2015.1000000000015      1209.58000000000015
esaverina 4775.88 6120.4800000000023      1344.60000000000023
irisk 45591.960000000588      46946.0400000002184      1354.07999999963056
oniq 8425.410000000003      9841.6500000000018      1416.2399999999987
levrana 2243.5600000000002      3664.0999999999998      1420.53999999999959
combluff 5491.26000000000005      4943.7699999999991      1422.40999999999985
smart 4457.2600000000004      5902.1400000000017      1444.88000000000128
shik 3341.2 4839.720000000007      1498.52000000000068
domix 10472.049999999994      12009.1700000000022      1537.120000000000827
artex 2730.699999999998      4327.2500000000017      1596.61000000000192
beautik 10493.949999999966      12222.949999999913      1728.9999999999472
milv 3904.9399999999964      5642.0100000000008      1737.070000000000838
mamura 31266.079999999821      33058.469999999708      1792.3899999998753
f.o.x 6624.229999999982      8577.2800000000004      1953.0500000000022
kapous 11927.159999999988      14093.0800000000158      2165.9200000000026
concept 11032.139999999925      13380.399999999993      2348.26000000000057
estel 21756.7500000000342      24142.6700000000022      2385.9199999999878
kaypro 881.3399999999998      3268.6999999999995      2387.3599999999995
benovy 499.6200000000002      3259.9760000000001      2850.3500000000001
itelwax 21940.239999999722      24709.369999999993      2859.1500000000161
yoko 8756.909999999949      11707.879999999996      2950.97000000000466
haruyama 9390.689999999991      12352.910000000013      2962.22000000001394
marathon 7280.749999999997      10273.1 2992.3500000000003
lovely 8704.379999999952      1139.0600000000045      2234.68000000000093
hpw.style 11572.150000001699      14837.4400000000812      3265.2899999999113
staleks 8519.730000000003      11875.610000000008      3355.88000000000774
freedecor 3421.7799999999971      7671.8000000000175      4250.0200000000204
runail 71539.279999999923      76758.660000000098      5219.3800000001649
polarus 6913.720000000003      11371.9300000000018      5358.21000000000155
cosmoprofi 8322.810000000007      14536.990000000016      6214.1800000000089
jessnail 26287.839999999916      33345.229999999992      7057.3900000000007
strong 29196.629999999994      38671.269999999924      9474.6399999999985
magarden 22161.3900000000138      33366.210000000009      10464.8199999999949
lianail 5892.839999999975      16394.2400000000245      10501.400000000027
uno 35302.829999999977      51039.7499999998035      15737.7199999998262
grattol 35445.54000000011      71472.710000000068      36027.1699999999576
474679.05999999623      619509.2399999934      144830.180000003108
Time taken: 71.737 seconds, Fetched: 161 row(s)
hive>

```

Insight:

- There are so many brands which have the increment from October to November. Among which 'Grattol' brand has the highest total increment and 'Ovale' seems to have least increment.

Question 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.

Query: SELECT user_id, SUM(price) as Total_Expenditure FROM retailstore WHERE event_type='purchase' GROUP BY user_id ORDER BY Total_Expenditure DESC LIMIT 10;

```

hive> SELECT user_id, SUM(price) as Total_Expenditure FROM retailstore WHERE event_type='purchase' GROUP BY user_id ORDER BY Total_Expenditure DESC LIMIT 10;
Query ID = hadoop_20211019173930_fe2a4d12-d6bb-4ca5-8edd-470bad470227
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634660723820_0008)

-----
      VERTICES      MODE      STATUS      TOTAL      COMPLETED      RUNNING      PENDING      FAILED      KILLED
-----
Map 1 ..... container      SUCCEEDED      2      2      0      0      0      0
Reducer 2 ..... container      SUCCEEDED      6      6      0      0      0      0
Reducer 3 ..... container      SUCCEEDED      1      1      0      0      0      0
-----
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 59.12 s
-----
OK
557790271      2715.8699999999991
150318419      1645.97
562167663      1352.8500000000004
531900924      1329.4500000000003
557850743      1285.4800000000002
522130011      1105.3899999999994
561592095      1109.6999999999996
421950134      1097.5899999999995
566576008      1056.3600000000017
521347209      1040.9099999999999
Time taken: 59.961 seconds, Fetched: 10 row(s)

```

Insight:

- Here is the list of the top 10 users who have spent the most.

To create table with Partitioning and Bucketing below commands need to be executed one by one separately.

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

```

[had00p@ip-172-31-84-44 ~]$ set hive.exec.dynamic.partition.mode=nonstrict;
[had00p@ip-172-31-84-44 ~]$ set hive.exec.dynamic.partition=true;
[had00p@ip-172-31-84-44 ~]$ set hive.enforce.bucketing=true;

```

Table optimization steps:-

- Command to create table '**Dyn_retailstore_part**' with partition on '**event_type**' attribute and bucket(cluster) on '**price**' attribute.

Query: CREATE TABLE IF NOT EXISTS Dyn_retailstore_part (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 7 BUCKETS ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS TEXTFILE ;

```

hive> CREATE TABLE IF NOT EXISTS Dyn_retailstore_part ( 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 7 BUCKETS ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' STORED AS TEXTFILE ;
OK
Time taken: 0.202 seconds
hive> INSERT INTO TABLE Dyn_retailstore_part PARTITION (event_type) SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type FROM retailstore ;
FAILED: SemanticException [Error 10096]: Dynamic partition strict mode requires at least one static partition column. To turn this off set hive.exec.dynamic.partition.mode=nonstrict
hive> set hive.exec.dynamic.partition.mode=nonstrict;
hive> set hive.exec.dynamic.partition=true;
hive> set hive.enforce.bucketing=true;
hive> INSERT INTO TABLE Dyn_retailstore_part PARTITION (event_type) SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type FROM retailstore ;
Query ID = hadoop_20211020164051_03c02e7b-2f13-4468-ae42-73eb310e1904
Total jobs = 1
Launching Job 1 out of 1
TeX session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_163474653704_0004)

-----
VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    5         5         0         0         0         0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 165.81 s
-----
Loading data to table default.dyn_retailstore_part partition (event_type=null)
Loaded : 4/4 partitions.
Time taken to load dynamic partitions: 0.817 seconds
Time taken for adding to write entity : 0.006 seconds
OK
event_time    product_id    category_id    category_code    brand    price    user_id    user_session    event_type
Time taken: 178.87 seconds

```

- To add data into partitioned and bucketed table we need to get it from already created table '**Dyn_retailstore_part**'.

Query: INSERT INTO TABLE Dyn_retailstore_part PARTITION (event_type) SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type FROM retailstore ;

- Command to check the successful creation of partitioned and bucketed table first we need to exit from Hive environment by executing 'EXIT;' command and then run below mentioned command.

```
hive> exit;
[hadoop@ip-172-31-92-176 ~]$
```

```
[hadoop@ip-172-31-84-44 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_retailssstore_part
Found 4 items
drwxrwxrwt - hadoop hdfsadmingroup 0 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart
drwxrwxrwt - hadoop hdfsadmingroup 0 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase
drwxrwxrwt - hadoop hdfsadmingroup 0 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart
drwxrwxrwt - hadoop hdfsadmingroup 0 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view
```

```
[hadoop@ip-172-31-84-44 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase
Found 7 items
-rwxrwxrwt 1 hadoop hdfsadmingroup 13052654 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000000_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 9399111 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000001_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 12636711 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000002_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 10650131 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000003_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 7226455 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000004_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 10737803 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000005_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 7825305 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=purchase/000006_0
```

```
[hadoop@ip-172-31-84-44 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_retailssstore_part/event_type=cart
Found 7 items
-rwxrwxrwt 1 hadoop hdfsadmingroup 57724286 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000000_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 43094161 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000001_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 56823661 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000002_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 49030059 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000003_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 31050141 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000004_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 48253679 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000005_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 34272441 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=cart/000006_0
[hadoop@ip-172-31-84-44 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart
Found 7 items
-rwxrwxrwt 1 hadoop hdfsadmingroup 39017824 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000000_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 29421828 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000001_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 38713899 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000002_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 31959876 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000003_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 19751571 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000004_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 31335021 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000005_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 22175799 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=remove_from_cart/000006_0
[hadoop@ip-172-31-84-44 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_retailssstore_part/event_type=view
Found 7 items
-rwxrwxrwt 1 hadoop hdfsadmingroup 88831872 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000000_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 73953212 2021-10-20 16:42 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000001_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 85620113 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000002_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 71874121 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000003_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 48335545 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000004_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 72515614 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000005_0
-rwxrwxrwt 1 hadoop hdfsadmingroup 56694677 2021-10-20 16:43 /user/hive/warehouse/dyn_retailssstore_part/event_type=view/000006_0
```

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

```
hive> select sum(price) AS Total_Revenue_Oct from dyn_retailssstore_part where month(event_time)=10 and event_type='purchase' MINUS select sum(price) AS Total_Revenue_Nov from dyn_retailssstore_part where month(event_time)=11 and event_type='purchase' ;
Query ID = hadoop_20211020170255_0a9c44b8-4514-4a6b-a304-c058bc862063
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634746553704_0005)

-----
VERTICES      MODE        STATUS      TOTAL    COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    7         7         0         0         0         0
Map 6 ..... container  SUCCEEDED    7         7         0         0         0         0
Reducer 2 .... container  SUCCEEDED    1         1         0         0         0         0
Reducer 3 .... container  SUCCEEDED    2         2         0         0         0         0
Reducer 5 .... container  SUCCEEDED    2         2         0         0         0         0
Reducer 7 .... container  SUCCEEDED    1         1         0         0         0         0
Reducer 8 .... container  SUCCEEDED    2         2         0         0         0         0
-----
VERTICES: 07/07 [=====>] 100% ELAPSED TIME: 33.31 s
-----
OK
1211538.4300000786
Time taken: 34.723 seconds, Fetched: 1 row(s)
```

Insight:

Earlier the query was taking 134 seconds to complete but with the partitioning and bucketing the query took only 35 seconds to complete.

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

```
hive> SELECT user_id, SUM(price) as Total_Expenditure FROM dyn_retailstore_part WHERE event_type='purchase' GROUP BY user_id ORDER BY Total_Expenditure DESC LIMIT 10;
Query ID = hadoop_20211020170810_98d31001-3983-4a1d-ac09-273ee796e16d
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1634746553704_0005)

-----
VERTICES      MODE        STATUS      TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED      7         7         0         0         0         0
Reducer 2 ..... container  SUCCEEDED      2         2         0         0         0         0
Reducer 3 ..... container  SUCCEEDED      1         1         0         0         0         0
-----
VERTICES: 03/03 [=====>>>] 100% ELAPSED TIME: 25.00 s
-----
OK
557790271      2715.8699999999996
150318419      1645.97
562167663      1352.85
831900924      1229.45
557850743      1295.48
522130011      1185.3900000000003
561592095      1109.7000000000003
431950134      1097.59
566576008      1056.3599999999997
521347209      1040.9099999999999
Time taken: 26.545 seconds, Fetched: 10 row(s)
```

Insight:

Earlier the query was taking 60 seconds to complete but with the partitioning and bucketing the query took only 27 seconds to complete. Hence, optimized table gives better performance in execution time.

TERMINATION PROCESS:

Dropping Database:

```
hive> drop database RetailDB CASCADE;
OK
Time taken: 0.706 seconds
hive> show databases ;
NoViableAltException(240[846:1: ddlStatement : ( createDatabaseStatement | switchDatabaseStatement | dropDatabaseStatement | createTableStatement | dropTableStatement | truncateTableStatement | alterStatement | descStatement | showStatement | metastoreCheck | createViewStatement | createMaterializedViewStatement | dropViewStatement | dropMaterializedViewStatement | createFunctionStatement | createMacroStatement | createIndexStatement | dropIndexStatement | dropFunctionStatement | reloadFunctionStatement | dropMacroStatement | analyzeStatement | lockStatement | unlockStatement | lockDatabase | createRoleStatement | dropRoleStatement | ( grantPrivileges )=> grantPrivileges | ( revokePrivileges )=> revokePrivileges | showGrants | showRolePrincipals | showRoles | grantRole | revokeRole | setRole | showCurrentRole | abortTransactionStatement );])
    at org.antlr.runtime.DFA.noViableAlt(DFA.java:158)
    at org.antlr.runtime.DFA.predict(DFA.java:144)
    at org.apache.hadoop.hive.q1.parse.HiveParser.ddlStatement(HiveParser.java:3757)
    at org.apache.hadoop.hive.q1.parse.HiveParser.executeStatement(HiveParser.java:2382)
    at org.apache.hadoop.hive.q1.parse.HiveParser.statement(HiveParser.java:1332)
    at org.apache.hadoop.hive.q1.parse.ParseDriver.parse(ParseDriver.java:208)
    at org.apache.hadoop.hive.q1.parse.ParseUtils.parse(ParseUtils.java:77)
    at org.apache.hadoop.hive.q1.parse.ParseUtils.parse(ParseUtils.java:70)
    at org.apache.hadoop.hive.q1.Driver.compileInternal(Driver.java:469)
    at org.apache.hadoop.hive.q1.Driver.compileInternal(Driver.java:1317)
    at org.apache.hadoop.hive.q1.Driver.runInternal(Driver.java:1457)
    at org.apache.hadoop.hive.q1.Driver.run(Driver.java:1237)
    at org.apache.hadoop.hive.q1.Driver.run(Driver.java:1237)
    at org.apache.hadoop.hive.cli.CliDriver.processLocalCmd(CliDriver.java:233)
    at org.apache.hadoop.hive.cli.CliDriver.processCmd(CliDriver.java:184)
    at org.apache.hadoop.hive.cli.CliDriver.processLine(CliDriver.java:403)
    at org.apache.hadoop.hive.cli.CliDriver.executeDriver(CliDriver.java:821)
    at org.apache.hadoop.hive.cli.CliDriver.run(CliDriver.java:759)
    at org.apache.hadoop.hive.cli.CliDriver.main(CliDriver.java:686)
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.RunJar.run(RunJar.java:244)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:158)
FAILED: ParseException line 1:5 cannot recognize input near 'show' 'databases' '<EOF>' in ddl statement
hive> show databases;
OK
default
Time taken: 0.025 seconds, Fetched: 1 row(s)
hive>
```

Terminating the EMR Cluster:

- After completing our analysis, we should terminate the EMR cluster.

Amazon EMR

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

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-1GU5E92OPZ8II
Creation date: 2021-10-20 21:38 (UTC+5:30)
Elapsed time: 1 hour, 2 minutes
After last step completes: Cluster waits
Termination protection: Off [Change](#)
Tags: -- [View All / Edit](#)
Master public DNS: ec2-3-84-164-145.compute-1.amazonaws.com [Connect to the Master Node Using SSH](#)

Configuration details

Release label: emr-5.33.1
Hadoop distribution: Amazon 2.10.1
Applications: Hive 2.3.7, Pig 0.17.0, Hue 4.9.0
Log URI: s3://aws-logs-666646323686-us-east-1/elasticmapreduce/
EMRFS consistent view: Disabled
Custom AMI ID: --

Application user interfaces

Persistent user [YARN timeline server, Tez UI](#)
Interfaces [\[?\]](#)
On-cluster user Not Enabled [Enable an SSH Connection](#)
Interfaces [\[?\]](#)

Network and hardware

Availability zone: us-east-1a
Subnet ID: [subnet-96fc05d385c6e605](#)
Master: **Running** 1 m4.large
Core: **Running** 1 m4.large
Task: --
Cluster scaling: Not enabled
Auto-termination: Not enabled

Security and access

Key name: casestudykeys
EC2 instance profile: EMR_EC2_DefaultRole

Amazon EMR

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

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-1GU5E92OPZ8II
Creation date: 2021-10-20 21:38 (UTC+5:30)
Elapsed time: 1 hour, 2 minutes
After last step completes: Cluster waits
Termination protection: Off [Change](#)
Tags: -- [View All / Edit](#)
Master public DNS: ec2-3-84-164-145.compute-1.amazonaws.com [Connect to the Master Node Using SSH](#)

Configuration details

Release label: emr-5.33.1
Hadoop distribution: Amazon 2.10.1
Applications: Hive 2.3.7, Pig 0.17.0, Hue 4.9.0
Log URI: s3://aws-logs-666646323686-us-east-1/elasticmapreduce/
EMRFS consistent view: Disabled
Custom AMI ID: --

Application user interfaces

Persistent user [YARN timeline server, Tez UI](#)
Interfaces [\[?\]](#)
On-cluster user Not Enabled [Enable an SSH Connection](#)
Interfaces [\[?\]](#)

Network and hardware

Availability zone: us-east-1a
Subnet ID: [subnet-96fc05d385c6e605](#)
Master: **Running** 1 m4.large
Core: **Running** 1 m4.large
Task: --
Cluster scaling: Not enabled
Auto-termination: Not enabled

Security and access

Key name: casestudykeys
EC2 instance profile: EMR_EC2_DefaultRole

Terminate cluster

Are you sure you want to terminate this cluster?

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

[Cancel](#) [Terminate](#)

Services

Search for services, features, marketplace products, and docs

[Alt+S]

upgradamarankasahoo@6866-4632-3686

N. Virginia

Support

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

Create clusterView detailsCloneTerminate

Filter: All clustersFilter clusters ...3 clusters (all loaded)

	Name	ID	Status	Creation time (UTC+5:30)	Elapsed time	Normalized instance hours
<input type="checkbox"/>	CaseStudycluster	j-1GU5E92OPZ8ll	Terminating User request	2021-10-20 21:38 (UTC+5:30)	1 hour, 3 minutes	0
<input type="checkbox"/>	CaseStudycluster	j-197J1CU03ZUK	Terminated User request	2021-10-19 21:47 (UTC+5:30)	1 hour, 25 minutes	16
<input type="checkbox"/>	CaseStudycluster	j-10Q2XLR2OXZ4Q	Terminated User request	2021-10-18 22:10 (UTC+5:30)	47 minutes	8