



# Hive Case Study

Using HiveQL

# Case Study Steps

## Key-pair creation and generation

The screenshot displays the AWS Management Console interface for creating a new key pair. The first part shows the 'Create key pair' dialog box, and the second part shows the 'Key pairs (1/3)' list.

**Create key pair dialog:**

- Name:** hive\_case\_study\_key
- File format:** pem (selected), ppk (unselected)
- Tags (Optional):** No tags associated with the resource.
- Buttons:** Cancel, Create key pair

**Key pairs (1/3) list:**

	Name	Fingerprint	ID
<input type="checkbox"/>	demo_key_pair	47:8b:69:24:d3:9d:03:a2:6c:69:a7:e2:3...	key-072e1bf076205ba71
<input checked="" type="checkbox"/>	hive_case_study_key	41:31:34:45:32:b5:f7:1c:5e:fc:ff:5a:20:...	key-0afdb0271f3164dd9
<input type="checkbox"/>	Jay_060198243535_2020-12-14 14:1...	75:7b:8d:c5:3f:42:73:e2:70:ac:01:38:4...	key-01f41ef121e850931

**PUTTY Key Generator**

File Key Conversions Help

**Key**

Public key for pasting into OpenSSH authorized\_keys file:

```
ssh-rsa
AAAAB3NzaC1yc2EAAAADAQABAAQAUzBZX41TfMGczkx7Hin9JkBdtsGVyTsh8
RNwftBHQi8Y9M4RhZ8KgsKJ8FKTnO2MBanusJqyicL2lce5waFjio272ulz9aaHKd41VI
hqG/Fvkpi590bMiGyv7Q8WY4SP/fqlRbpIIDIZNswxajXxsiTgg/waDYrdu7wLQ42q/GE
PCNslwWH7lpS7kEMKcb/zz1k8zq3HoZbQnZaRQuTz/VPx7LFhRibUDub2++uybe
```

Key fingerprint: ssh-rsa 2048 8c:13:ab:a6:8f:7e:3f:59:b1:10:1e:02:ee:8b:70:b0

Key comment: imported-openssh-key

Key passphrase:

Confirm passphrase:

**Actions**

Generate a public/private key pair Generate

Load an existing private key file Load

Save the generated key Save public key Save private key

**Parameters**

Type of key to generate:

☒ RSA ☐ DSA ☐ ECDSA ☐ Ed25519 ☐ SSH-1 (RSA)

Number of bits in a generated key: 2048

hive_case_study_key.pem	1/18/2021 11:09 AM	PEM File	2 KB
hive_case_study_key.ppk	1/18/2021 11:10 AM	PPK File	2 KB

## Created EMR cluster 'Hive\_Case\_Study'

Subscription Details | Nuvepro x EMR - AWS Console

console.aws.amazon.com/elasticmapreduce/home?region=us-east-1#cluster-details:j-3B5QOU4EG0966

aws Services Search for services, features, marketplace products, and docs [Alt+S] upgradjaythakur @ 0601-9824-3535 N. Virginia Support

Amazon EMR

EMR on EC2

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

EMR on EKS

Virtual clusters

Help

What's new

Clone Terminate AWS CLI export

Cluster: Hive\_Case\_Study Waiting Cluster ready after last step completed.

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

**Summary**

ID: j-3B5QOU4EG0966

Creation date: 2021-01-17 22:58 (UTC+5:30)

Elapsed time: 1 hour, 22 minutes

After last step completes: Cluster waits

Termination protection: Off [Change](#)

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

Master public DNS: ec2-52-23-176-82.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-060198243535-us-east-1/elasticmapreduce/

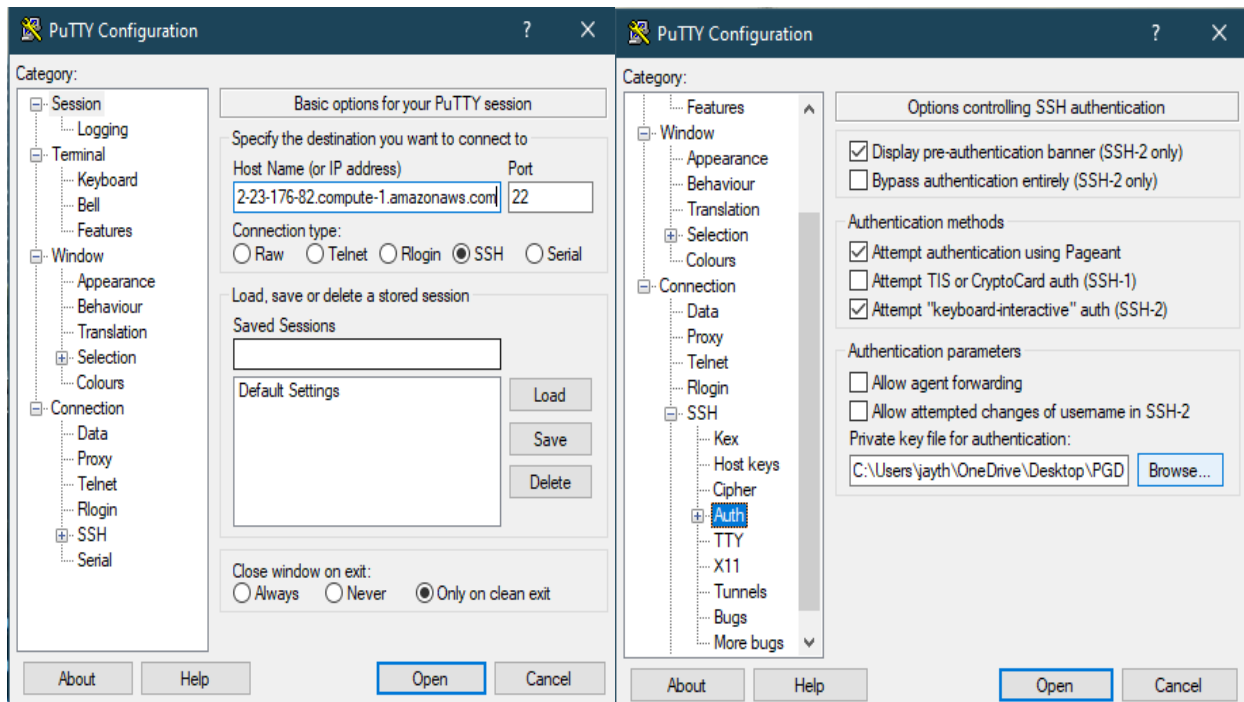
EMRFS consistent view: Disabled

Custom AMI ID: --

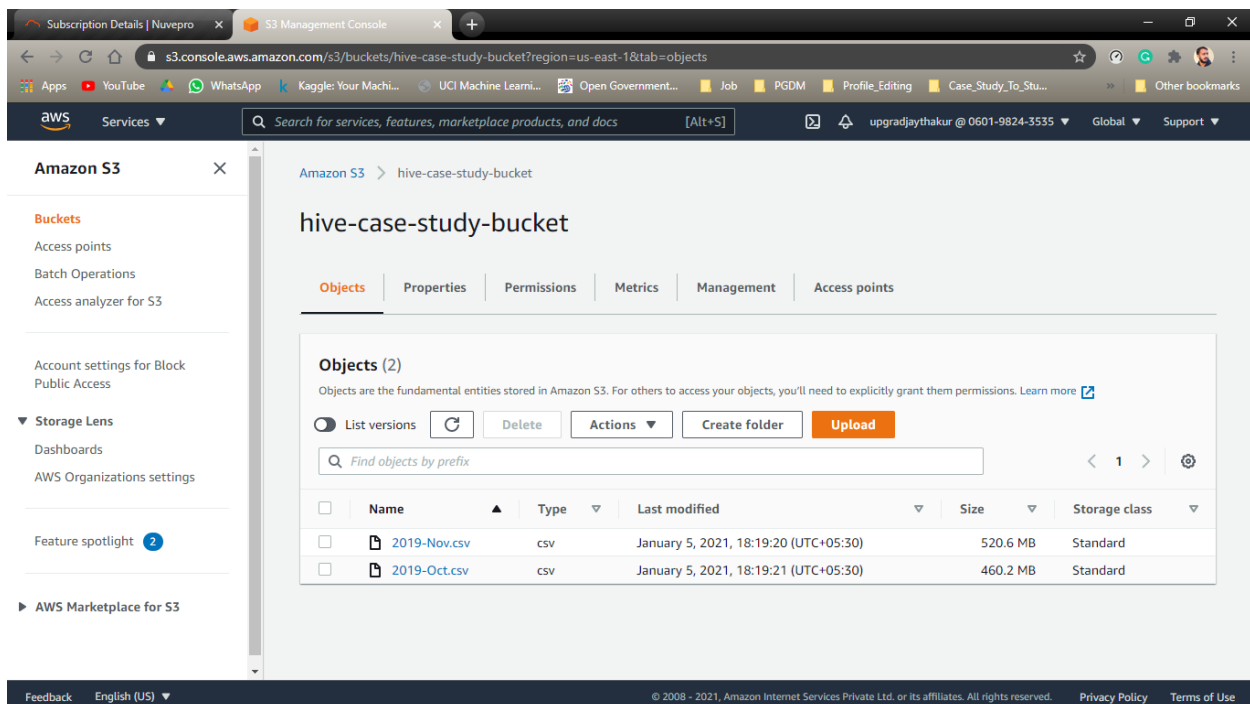
**Application user interfaces**

Feedback English (US) © 2008 - 2021, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved. Privacy Policy Terms of Use

## Starting terminal using Putty



## S3 Bucket to store data files



## 1. Command to check for already present directories in HDFS

```
- hadoop fs -ls /
```

**Output:**

Found 4 items

```
drwxr-xr-x - hdfs hadoop      0 2021-01-17 17:34 /apps
```

```
drwxrwxrwt - hdfs hadoop    0 2021-01-17 17:36 /tmp
```

```
drwxr-xr-x - hdfs hadoop    0 2021-01-17 17:34 /user
```

```
drwxr-xr-x - hdfs hadoop    0 2021-01-17 17:34 /var
```

**Insights:**

- All the above directories are in-built in HDFS.
- Either these directories can be used to create our temporary directory to store data files or create a separate temporary directory.

[illegible]

## 2. Creating new temporary directory i.e., 'HiveCaseStudy' to store data file in the already present directory (Permanent) i.e., 'user'

```
- hadoop fs -mkdir /user/HiveCaseStudy/
```

```
hadoop@ip-172-31-93-164:~$ hadoop fs -mkdir /user/HiveCaseStudy/
```

### 3. Command to check creation of new temporary Directory in 'user' directory

- `hadoop fs -ls /user/`

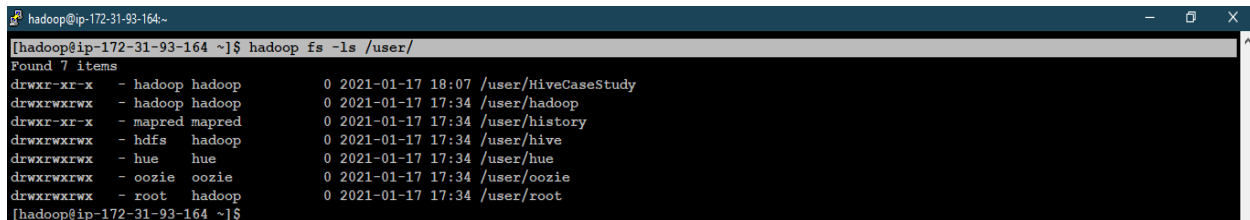
#### Output:

Found 7 items

```
drwxr-xr-x - hadoop hadoop      0 2021-01-17 18:07 /user/HiveCaseStudy
drwxrwxrwx - hadoop hadoop      0 2021-01-17 17:34 /user/hadoop
drwxr-xr-x - mapred mapred      0 2021-01-17 17:34 /user/history
drwxrwxrwx - hdfs hadoop        0 2021-01-17 17:34 /user/hive
drwxrwxrwx - hue hue            0 2021-01-17 17:34 /user/hue
drwxrwxrwx - oozie oozie        0 2021-01-17 17:34 /user/oozie
drwxrwxrwx - root hadoop        0 2021-01-17 17:34 /user/root
```

#### Insights:

- There will always be some files within the permanent directories of the HDFS.



```
hadoop@ip-172-31-93-164:~$ hadoop fs -ls /user/
Found 7 items
drwxr-xr-x - hadoop hadoop      0 2021-01-17 18:07 /user/HiveCaseStudy
drwxrwxrwx - hadoop hadoop      0 2021-01-17 17:34 /user/hadoop
drwxr-xr-x - mapred mapred      0 2021-01-17 17:34 /user/history
drwxrwxrwx - hdfs hadoop        0 2021-01-17 17:34 /user/hive
drwxrwxrwx - hue hue            0 2021-01-17 17:34 /user/hue
drwxrwxrwx - oozie oozie        0 2021-01-17 17:34 /user/oozie
drwxrwxrwx - root hadoop        0 2021-01-17 17:34 /user/root
hadoop@ip-172-31-93-164:~$
```

### 4. Command to load 1st data file '2019-Oct.csv' from S3 storage into HDFS storage as 'October.csv'

- hadoop distcp s3://hive-case-study-bucket/2019-Oct.csv /user/HiveCaseStudy/October.csv

```
hadoop@ip-172-31-80-31:~$ hadoop distcp s3://hive-case-study-bucket/2019-Oct.csv /user/HiveCaseStudy/October.csv
21/01/17 08:50:44 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, overwrite=false, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=100, sslConfigurationFile='null', copyStrategy='uniform size', preserveStatus=[], preserveRawAttrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://hive-case-study-bucket/2019-Oct.csv], targetPath=/user/HiveCaseStudy/October.csv, targetPathExists=false, filtersFile='null'}
21/01/17 08:50:44 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-80-31.ec2.internal/172.31.80.31:8032
21/01/17 08:50:49 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
21/01/17 08:50:49 INFO tools.SimpleCopyListing: Build file listing completed.
21/01/17 08:50:49 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
21/01/17 08:50:49 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
21/01/17 08:50:49 INFO tools.DistCp: Number of paths in the copy list: 1
21/01/17 08:50:49 INFO tools.DistCp: Number of paths in the copy list: 1
21/01/17 08:50:49 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-80-31.ec2.internal/172.31.80.31:8032
21/01/17 08:50:50 INFO mapreduce.JobSubmitter: number of splits:1
21/01/17 08:50:50 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1610869923603_0001
21/01/17 08:50:51 INFO Impl.YarnClientImpl: Submitted application application_1610869923603_0001
21/01/17 08:50:51 INFO mapreduce.Job: The url to track the job: http://ip-172-31-80-31.ec2.internal:20888/proxy/application_1610869923603_0001/
21/01/17 08:50:51 INFO tools.DistCp: DistCp job-id: job_1610869923603_0001
21/01/17 08:50:51 INFO mapreduce.Job: Running job: job_1610869923603_0001
21/01/17 08:51:01 INFO mapreduce.Job: Job job_1610869923603_0001 running in uber mode : false
21/01/17 08:51:01 INFO mapreduce.Job: map 0% reduce 0%
21/01/17 08:51:18 INFO mapreduce.Job: map 100% reduce 0%
21/01/17 08:51:21 INFO mapreduce.Job: Job job_1610869923603_0001 completed successfully
21/01/17 08:51:21 INFO mapreduce.Job: Counters: 38

File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=172475
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=361
  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)=579520
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=18110
Total vcore-milliseconds taken by all map tasks=18110
Total megabyte-milliseconds taken by all map tasks=18544640

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)=275
  CPU time spent (ms)=19580
  Physical memory (bytes) snapshot=589647872
  Virtual memory (bytes) snapshot=3304558592
  Total committed heap usage (bytes)=481820672

File Input Format Counters
  Bytes Read=225
File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=482542278
  Bytes Expected=482542278
  Files Copied=1

hadoop@ip-172-31-80-31:~$
```

5. Command to load 2nd data file '2019-Nov.csv' from S3 storage into HDFS storage as 'November.csv'

- hadoop distcp s3://hive-case-study-bucket/2019-Nov.csv /user/HiveCaseStudy/November.csv

```
hadoop@ip-172-31-80-31:~$ hadoop distcp s3://hive-case-study-bucket/2019-Nov.csv /user/HiveCaseStudy/November.csv
21/01/17 08:55:02 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='uniform size', preserveStatus=[], preserveRawAttrs=false, atomicWorkPath=null, logPath=null, sourceFileListing=null, sourcePaths=[s3://hive-case-study-bucket/2019-Nov.csv], targetPath=/user/HiveCaseStudy/November.csv, targetPathExists=false, filtersFile='null'}
21/01/17 08:55:02 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-80-31.ec2.internal/172.31.80.31:8032
21/01/17 08:55:07 INFO tools.SimpleCopyListing: Paths (files+dirs) cnt = 1; dirCnt = 0
21/01/17 08:55:07 INFO tools.SimpleCopyListing: Build file listing completed.
21/01/17 08:55:07 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
21/01/17 08:55:07 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
21/01/17 08:55:07 INFO tools.DistCp: Number of paths in the copy list: 1
21/01/17 08:55:07 INFO tools.DistCp: Number of paths in the copy list: 1
21/01/17 08:55:07 INFO client.RMPProxy: Connecting to ResourceManager at ip-172-31-80-31.ec2.internal/172.31.80.31:8032
21/01/17 08:55:07 INFO mapreduce.JobSubmitter: number of splits:1
21/01/17 08:55:08 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1610869923603_0002
21/01/17 08:55:08 INFO Impl.YarnClientImpl: Submitted application application_1610869923603_0002
21/01/17 08:55:08 INFO mapreduce.Job: The url to track the job: http://ip-172-31-80-31.ec2.internal:20888/proxy/application_1610869923603_0002/
21/01/17 08:55:08 INFO tools.DistCp: DistCp job-id: job_1610869923603_0002
21/01/17 08:55:08 INFO mapreduce.Job: Running job: job_1610869923603_0002
21/01/17 08:55:16 INFO mapreduce.Job: Job job_1610869923603_0002 running in uber mode : false
21/01/17 08:55:16 INFO mapreduce.Job: map 0% reduce 0%
21/01/17 08:55:34 INFO mapreduce.Job: map 100% reduce 0%
21/01/17 08:55:35 INFO mapreduce.Job: Job job_1610869923603_0002 completed successfully
21/01/17 08:55:36 INFO mapreduce.Job: Counters: 38

File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=172480
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=362
  HDFS: Number of bytes written=545839412
  HDFS: Number of read operations=12
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=4
  S3: Number of bytes read=545839412
  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)=539200
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=16850
Total vcore-milliseconds taken by all map tasks=16850
Total megabyte-milliseconds taken by all map tasks=17254400

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)=300
  CPU time spent (ms)=20600
  Physical memory (bytes) snapshot=596889600
  Virtual memory (bytes) snapshot=3302703104
  Total committed heap usage (bytes)=503316480

File Input Format Counters
  Bytes Read=225
File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=545839412
  Bytes Expected=545839412
  Files Copied=1

hadoop@ip-172-31-80-31:~$
```

## 6. Command to check successful loading of data files into the already created new temporary directory of HDFS i.e., 'HiveCaseStudy'

- hadoop fs -ls /user/HiveCaseStudy/

**Output:**

Found 2 items

-rw-r--r-- 1 hadoop hadoop 545839412 2021-01-17 14:54 /user/HiveCaseStudy/November.csv

-rw-r--r-- 1 hadoop hadoop 482542278 2021-01-17 14:51 /user/HiveCaseStudy/October.csv



```
hadoop@ip-172-31-94-188:~$ hadoop fs -ls /user/HiveCaseStudy/
Found 2 items
-rw-r--r-- 1 hadoop hadoop 545839412 2021-01-17 14:54 /user/HiveCaseStudy/November.csv
-rw-r--r-- 1 hadoop hadoop 482542278 2021-01-17 14:51 /user/HiveCaseStudy/October.csv
hadoop@ip-172-31-94-188:~$
```

## 7. Command to start Hive system

- hive

## 8. Creating an External table i.e., 'Shopping' which will hold the data for both the data files stored in temporary directory of HDFS.

- CREATE EXTERNAL TABLE IF NOT EXISTS Shopping (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");

### Output:

OK

Time taken: 0.968 seconds

```
hadoop@ip-172-31-94-188:~$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> CREATE EXTERNAL TABLE IF NOT EXISTS Shopping(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: 0.968 seconds
hive>
```

## 9. Command to enable heading in the output

- set hive.cli.print.header=True;

## 10. Simple HiveQL command to check successful creation of table and storage of data from both data files into table

### Query:

SELECT \* FROM Shopping LIMIT 5;

### Output:

OK

shopping.event_time	shopping.event_type	shopping.product_id	shopping.category_id	shopping.category_code	shopping.brand	shopping.price	shopping.user_id	shopping.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	
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-a1cc-af0575a34ffb		

Time taken: 0.181 seconds, Fetched: 5 row(s)

```

hadoop@ip-172-31-94-188:~
hive> SELECT * FROM Shopping LIMIT 5;
OK
shopping.event_time  shopping.event_type  shopping.product_id  shopping.category_id  shopping.category_code  shopping.brand  shopping.pric
e      shopping.user_id      shopping.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
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-a1cc-af057
5a34ffb
Time taken: 0.181 seconds, Fetched: 5 row(s)
hive>

```

## Questions

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

**Query:**

```
SELECT SUM(price) AS Total_Revenue_October
FROM Shopping
WHERE date_format(event_time, 'MM')=10
AND
event_type='purchase';
```

**Output:**

Query ID = hadoop\_20210117151333\_a5f18170-c287-4638-95b7-e98a5e28fe0d

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_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: 121.90 s

-----  
OK

total\_revenue\_october

1211538.4299997438

Time taken: 133.224 seconds, Fetched: 1 row(s)

```
hadoop@ip-172-31-94-188:~$
hive> SELECT SUM(price) AS Total_Revenue_October
> FROM Shopping
> WHERE date_format(event_time, 'MM')=10
> AND
> event_type='purchase';
Query ID = hadoop_20210117152447_1d6470de-e091-4264-97f4-11a30899f096
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_1610894517504_0005)

-----
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: 124.93 s
-----
OK
total revenue october
1211538.4299997438
Time taken: 133.224 seconds, Fetched: 1 row(s)
hive>
```

### Insights:

- The total revenue generated based on Purchase in the month of October of 2019 was **1,211,538.43** /-.

**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 Shopping
WHERE event_type='purchase'
GROUP BY date_format(event_time, 'MM');
```

### Output:

```
Query ID = hadoop_20210117153532_bdf9c31f-a5b2-4e95-8489-160046b2db17
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_1610894517504_0006)
```

```

-----
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: 62.73 s
-----

```

OK

months sum\_of\_purchases

10    245624

11    322417

Time taken: 71.767 seconds, Fetched: 2 row(s)

```

hadoop@ip-172-31-94-188:~$
hive> SELECT date format(event_time, 'MM') AS Months, COUNT(event_type) AS Sum_of_Purchases
> FROM Shopping
> WHERE event_type='purchase'
> GROUP BY date format(event_time, 'MM');
Query ID = hadoop_20210117153532_bdf9c31f-a5b2-4e95-8489-160046b2db17
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_1610894517504_0006)

-----
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: 62.73 s
-----

OK
months sum_of_purchases
10      245624
11      322417
Time taken: 71.767 seconds, Fetched: 2 row(s)
hive>

```

## Insights:

- It seems to be that there was **more purchase made in the month of November (11) i.e., 322,417 than in the month of October (10) i.e., 245,624.**
- Looking at these figures **we could assume that the month of November must be more profitable than the month of October.** But we can verify our assumption by conducting further investigations.

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

**Query:**

```
WITH Monthly_Revenue AS (  
SELECT  
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 shopping  
WHERE event_type= 'purchase'  
AND date_format(event_time, 'MM') in ('10', '11')  
)  
SELECT Nov_Revenue, Oct_Revenue, (Nov_Revenue - Oct_Revenue) AS Revenue_Difference FROM  
Monthly_Revenue;
```

**Output:**

Query ID = hadoop\_20210117154514\_6595f895-55b0-4187-b279-7e62ab7eb67b

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_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: 71.32 s  
-----
```

OK

nov\_revenue    oct\_revenue    revenue\_difference

1531016.900000122    1211538.4299997438    319478.4700003781

Time taken: 71.887 seconds, Fetched: 1 row(s)

```
hive> WITH Monthly_Revenue AS (  
  > SELECT  
  > 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 shopping  
  > WHERE event_type= 'purchase'  
  > AND date_format(event_time, 'MM') in ('10', '11')  
  > )  
  > SELECT Nov_Revenue, Oct_Revenue, (Nov_Revenue - Oct_Revenue) AS Revenue_Difference FROM Monthly_Revenue;  
Query ID = hadoop_20210117154514_6595f895-55b0-4187-b279-7e62ab7eb67b  
Total jobs = 1  
Launching Job 1 out of 1  
Status: Running (Executing on YARN cluster with App id application_1610894517504_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: 71.32 s  
-----  
OK  
nov_revenue    oct_revenue    revenue_difference  
1531016.900000122    1211538.4299997438    319478.4700003781  
Time taken: 71.887 seconds, Fetched: 1 row(s)  
hive>
```

Insights:

- On the basis of the results **considering purchase as event**, we could conclude that the **revenue generated in November of 2019 was more than** the revenue generated in the month of **October**. In other words, **November was more profitable for the company than October**.
- Company had a better sale in November, 2019.

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

**Query:**

SELECT DISTINCT SPLIT(category\_code, '\\.')[0] AS Category

FROM Shopping

WHERE SPLIT(category\_code, '\\.')[0] <> '';

**Output:**

Query ID = hadoop\_20210117154910\_654d0efc-2bcc-44a2-b180-92f5ef08f141

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_0006)

```
-----
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: 59.71 s
-----
```

OK

category

furniture

appliances

accessories

apparel

sport

stationery

Time taken: 60.323 seconds, Fetched: 6 row(s)

```
hadoop@ip-172-31-94-188:~$
hive> SELECT DISTINCT SPLIT(category_code,'\\\.')[0] AS Category
> FROM Shopping
> WHERE SPLIT(category_code,'\\\.')[0] <> '';
Query ID = hadoop_20210117154910_654d0efc-2bcc-44a2-b180-92f5ef08f141
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_0006)

-----
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: 59.71 s
-----

OK
category
furniture
appliances
accessories
apparel
sport
stationery
Time taken: 60.323 seconds, Fetched: 6 row(s)
hive>
```

Insights:



- There is total **6 different categories** under which company sells their different products.

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

**Query:**

```
SELECT SPLIT(category_code,'\\\.')[0] AS Category, COUNT(product_id) AS No_of_products
FROM Shopping
WHERE SPLIT(category_code,'\\\.')[0] <> ''
GROUP BY SPLIT(category_code,'\\\.')[0]
ORDER BY No_of_products DESC;
```

**Output:**

Query ID = hadoop\_20210117155132\_d38f4ba6-947d-41ca-b972-2c5530428355

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_0006)

```
-----
      VERTICES   MODE    STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
-----
Map 1 ..... container  SUCCEEDED   2      2      0      0      0      0
Reducer 2 ..... container  SUCCEEDED   5      5      0      0      0      0
Reducer 3 ..... container  SUCCEEDED   1      1      0      0      0      0
-----
```

VERTICES: 03/03 [=====>>] 100% ELAPSED TIME: 59.60 s

OK

```
category      no_of_products
```

appliances	61736
stationery	26722
furniture	23604
apparel	18232
accessories	12929
sport	2

Time taken: 60.311 seconds, Fetched: 6 row(s)

```

hive> SELECT SPLIT(category_code,'\\\.') [0] AS Category, COUNT(product_id) AS No_of_products
> FROM Shopping
> WHERE SPLIT(category_code,'\\\.') [0] <> ''
> GROUP BY SPLIT(category_code,'\\\.') [0]
> ORDER BY No_of_products DESC;

Query ID = hadoop_20210117155132_d38f4ba6-947d-41ca-b972-2c5530428355
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_0006)

-----
VERTICES      MODE        STATUS      TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED      2         2           0         0         0         0
Reducer 2 ..... container  SUCCEEDED      5         5           0         0         0         0
Reducer 3 ..... container  SUCCEEDED      1         1           0         0         0         0
-----
VERTICES: 03/03  [=====>>>] 100%  ELAPSED TIME: 59.60 s
-----
OK
category      no_of_products
appliances    61736
stationery    26722
furniture     23604
apparel 18232
accessories   12929
sport         2
Time taken: 60.311 seconds, Fetched: 6 row(s)
hive>

```

### Insights:

- 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 has only 2 products registered**. This must be due to low cosmetic products in the sports market.

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

### Query:

WITH Max\_Sales\_Brand AS (  
SELECT brand,

```

SUM(CASE WHEN date_format(event_time, 'MM')=10 THEN price ELSE 0 END) AS Oct_Sales,
SUM(CASE WHEN date_format(event_time, 'MM')=11 THEN price ELSE 0 END) AS Nov_Sales
FROM Shopping
WHERE (
event_type='purchase'
AND
date_format(event_time, 'MM') in ('10','11')
AND
brand <> '')
GROUP BY brand
)
SELECT brand, Nov_Sales + Oct_Sales AS Total_Sales
FROM Max_Sales_Brand
ORDER BY Total_Sales DESC
LIMIT 1;

```

**Output:**

Query ID = hadoop\_20210117155441\_e5643e59-8162-4068-a271-a8e536398dbc

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_0006)

---

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
<hr/>								
Map 1 .....	container	SUCCEEDED	2	2	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: 63.74 s

OK

brand total\_sales

runail 148297.9400000003

Time taken: 64.31 seconds, Fetched: 1 row(s)

```
hive> WITH Max_Sales_Brand AS (
  > SELECT brand,
  > SUM(CASE WHEN date_format(event_time, 'MM')=10 THEN price ELSE 0 END) AS Oct_Sales,
  > SUM(CASE WHEN date_format(event_time, 'MM')=11 THEN price ELSE 0 END) AS Nov_Sales
  > FROM Shopping
  > WHERE (
  > event_type='purchase'
  > AND
  > date_format(event_time, 'MM') in ('10','11')
  > AND
  > brand <> ''
  > GROUP BY brand
  > )
  > SELECT brand, Nov_Sales + Oct_Sales AS Total_Sales
  > FROM Max_Sales_Brand
  > ORDER BY Total_Sales DESC
  > LIMIT 1;

Query ID = hadoop_20210117155441_e5643e59-8162-4068-a271-a8e536398dbc
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_0006)

-----
VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         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: 63.74 s
-----
OK
brand  total_sales
runail 148297.9400000003
Time taken: 64.31 seconds, Fetched: 1 row(s)
hive>
```

Insights:

- Runail is the brand that has highest / maximum sales in the month of October and November of 2019 combined.
- It seems that Runail brand has high popularity among cosmetic lovers and bringing in more products related to Runail brand could help in increasing their profit.

## 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 Shopping
WHERE event_type='purchase'

AND

date_format(event_time, 'MM') IN ('10', '11')

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;

```

**Output:**

Query ID = hadoop\_20210117155852\_282b0369-324c-4c04-91c0-102abc59add0

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_0006)

```

-----
      VERTICES   MODE    STATUS TOTAL COMPLETED RUNNING PENDING FAILED KILLED
-----
Map 1 ..... container  SUCCEEDED   2     2     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: 69.69 s

OK

brand	oct_revenue	nov_revenue	sales_difference
ovale	2.54	3.1	0.56

cosima	20.23	20.92999999999993	0.699999999999922
grace	100.92000000000002	102.61000000000001	1.689999999999977
helloganic	0.0	3.1	3.1
skinity	8.88	12.440000000000001	3.560000000000005
bodyton	1376.3399999999974	1380.6399999999992	4.3000000000017735
moyou	5.71	10.280000000000001	4.570000000000001
neoleor	43.41	51.7	8.290000000000006
soleo	204.20000000000003	212.52999999999998	8.329999999999501
jaguar	1102.11	1110.6500000000003	8.540000000000418
tertio	236.16000000000008	245.79999999999978	9.639999999999702
fly	17.14	27.17	10.030000000000001
rasyan	18.79999999999997	28.93999999999994	10.13999999999997
deoproce	316.84	329.17000000000001	12.330000000000098
barbie	0.0	12.39	12.39
supertan	50.37000000000001	66.51000000000002	16.140000000000008
treaclemoon	163.36999999999995	181.48999999999995	18.120000000000005
kamill	63.00999999999999	81.49000000000002	18.480000000000032
juno	0.0	21.08	21.08
veraclara	50.109999999999985	71.21000000000001	21.100000000000023
glysolid	69.72999999999998	91.58999999999997	21.86
godefroy	401.22000000000002	425.12000000000006	23.89999999999864
binacil	0.0	24.259999999999998	24.259999999999998
blixz	38.949999999999996	63.39999999999998	24.44999999999998
profepil	93.36000000000003	118.02000000000005	24.660000000000025
estelare	444.80999999999943	471.87000000000009	27.06000000000148
orly	902.38000000000005	931.09000000000003	28.70999999999981
biore	60.650000000000006	90.31	29.65999999999997
beautyblender	78.74000000000001	109.41	30.66999999999987
vilenta	197.60000000000002	231.21000000000002	33.61000000000014

mavala	409.03999999999985	446.32	37.28000000000014
likato	296.05999999999999	340.96999999999999	44.910000000000025
ladykin	125.64999999999999	170.57	44.92
foamie	35.04	80.49	45.449999999999996
elskin	251.090000000000057	307.650000000000055	56.559999999999974
balbcare	155.32999999999996	212.380000000000025	57.050000000000296
koelcia	55.5	112.750000000000003	57.250000000000003
profhenna	679.22999999999999	736.850000000000005	57.620000000000057
kares	0.0	59.45	59.45
marutaka-foot	49.219999999999999	109.33	60.110000000000001
dewal	0.0	61.29	61.29
inm	288.02	351.210000000000001	63.190000000000011
laboratorium	246.49999999999991	312.52	66.020000000000007
cutrin	299.36999999999995	367.62	68.250000000000006
egomania	77.47	146.040000000000002	68.570000000000002
konad	739.82999999999991	810.670000000000003	70.840000000000117
nirvel	163.03999999999996	234.32999999999984	71.28999999999988
koelf	422.72999999999985	507.290000000000002	84.560000000000034
plazan	101.37	194.010000000000002	92.640000000000001
aura	83.95	177.51	93.55999999999999
kerasys	430.90999999999985	525.200000000000002	94.290000000000003
enjoy	41.349999999999994	136.570000000000002	95.220000000000003
depilflax	2707.0699999999994	2803.7799999999975	96.710000000000367
eos	54.339999999999996	152.61	98.270000000000001
carmex	145.08	243.36	98.28
batiste	772.39999999999999	874.16999999999994	101.76999999999953
osmo	645.58	762.310000000000002	116.730000000000013
dizao	819.130000000000012	945.50999999999998	126.37999999999852
igrobeauty	513.660000000000009	645.06999999999999	131.40999999999906

finish	98.38	230.38000000000008	132.00000000000009
nefertiti	233.52000000000007	366.64	133.11999999999992
elizavecca	70.53	204.3	133.77
miskin	158.04	293.07000000000005	135.03000000000006
latinoil	249.52	384.59	135.06999999999996
farmona	1692.4599999999996	1843.4300000000007	150.97000000000116
cristalinas	427.6299999999999	584.9499999999999	157.31999999999914
chi	358.94000000000002	538.61000000000002	179.67000000000002
matreshka	0.0	182.67000000000002	182.67000000000002
freshbubble	318.70000000000001	502.34000000000015	183.64000000000004
mane	66.78999999999999	260.26	193.47
keen	236.35000000000005	435.62	199.26999999999995
ecocraft	41.160000000000004	241.95	200.79
fedua	52.38	263.81000000000006	211.43000000000006
provoc	827.99000000000009	1063.8200000000006	235.82999999999997
skinlite	651.94000000000002	890.44999999999979	238.50999999999772
entity	479.71000000000015	719.2599999999993	239.5499999999978
trind	298.07000000000005	542.96000000000002	244.89000000000001
protokeratin	201.25000000000003	456.79000000000013	255.54000000000001
beauugreen	511.5099999999999	768.35	256.84000000000015
bluesky	10307.239999999858	10565.529999999713	258.28999999985535
candy	534.9599999999999	799.3799999999993	264.4199999999994
insight	1443.70000000000012	1721.9600000000003	278.2599999999991
kocostar	310.85000000000001	594.93000000000003	284.08000000000002
happyfons	801.92000000000006	1091.5900000000001	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	3318.959999999995	3657.4300000000026	338.47000000000753
ellips	245.8499999999999	606.0399999999996	360.1899999999997
lador	2083.6100000000004	2471.5300000000007	387.92000000000028
naomi	0.0	389.0	389.0
kiss	421.54999999999944	817.3299999999994	395.7799999999999
yu-r	271.41	673.7099999999998	402.2999999999998
sophin	1067.8600000000001	1515.5200000000011	447.6600000000001
farmavita	837.3699999999984	1291.9700000000003	454.60000000000184
bioaqua	942.8899999999996	1398.1199999999997	455.23
greymy	29.21	489.49	460.28000000000003
gehwol	1089.07	1557.6799999999982	468.6099999999983
matrix	3243.249999999999	3726.7400000000007	483.4900000000016
limoni	1308.9000000000003	1796.5999999999997	487.69999999999936
s.care	412.68	913.0699999999999	500.38999999999993
coifin	903.0000000000001	1428.4899999999998	525.4899999999997
uskusi	5142.2700000000017	5690.3100000000005	548.0399999999981
airnails	5118.8999999999939	5691.5199999999996	572.62000000000572
browxenna	14331.369999999995	14916.729999999976	585.3600000000026
kinetics	6334.2499999999945	6945.2600000000017	611.0100000000022
kosmekka	1181.4400000000003	1813.37	631.9299999999996
kaaral	4412.4299999999985	5086.069999999992	673.6399999999994
refectocil	2716.1800000000005	3475.5800000000007	759.4000000000024
rosi	3077.0399999999927	3841.5600000000013	764.52000000000204
solomeya	1899.6999999999992	2685.7999999999991	786.0999999999999
missha	1293.8299999999995	2150.2799999999984	856.4499999999989
levissime	2227.5000000000064	3085.3099999999977	857.8099999999913
art-visage	2092.7100000000001	2997.8000000000011	905.0900000000001
ecolab	262.85000000000001	1214.2999999999988	951.4499999999987
nagaraku	4369.7400000000054	5327.6800000000063	957.9400000000087

sanoto	157.14	1209.6799999999998	1052.54
markell	1768.7499999999998	2834.4300000000007	1065.6800000000019
metzger	5373.4500000000006	6457.1599999999988	1083.70999999999818
de.lux	1659.6999999999967	2775.5099999999968	1115.8100000000009
swarovski	1887.9299999999983	3043.1600000000003	1155.2300000000157
beauty-free	554.1700000000006	1782.8600000000163	1228.6900000000155
zeitun	708.6600000000004	2009.63	1300.9699999999998
joico	705.52	2015.1000000000015	1309.5800000000015
severina	4775.88	6120.4800000000023	1344.6000000000023
irisk	45591.960000000588	46946.0400000002184	1354.07999999963056
oniq	8425.410000000003	9841.6500000000018	1416.2399999999987
levrana	2243.5600000000002	3664.0999999999998	1420.5399999999959
roubloff	3491.3600000000003	4913.7699999999991	1422.4099999999985
smart	4457.2600000000004	5902.1400000000017	1444.8800000000128
shik	3341.2	4839.7200000000007	1498.5200000000068
domix	10472.049999999994	12009.1700000000022	1537.12000000000827
artex	2730.6399999999998	4327.2500000000017	1596.6100000000192
beautix	10493.9499999999966	12222.9499999999913	1728.9999999999472
milv	3904.9399999999964	5642.0100000000008	1737.07000000000838
masura	31266.079999999821	33058.469999999708	1792.38999999988753
f.o.x	6624.2299999999982	8577.2800000000004	1953.0500000000022
kapous	11927.1599999999898	14093.0800000000158	2165.9200000000026
concept	11032.1399999999925	13380.399999999993	2348.26000000000057
estel	21756.7500000000342	24142.6700000000022	2385.9199999999878
kaypro	881.3399999999998	3268.6999999999995	2387.3599999999995
benovy	409.62000000000002	3259.9700000000001	2850.3500000000001
italwax	21940.2399999999732	24799.3699999999893	2859.1300000000161
yoko	8756.9099999999949	11707.879999999996	2950.97000000000466
haruyama	9390.6899999999991	12352.910000000013	2962.22000000001394

marathon	7280.749999999997	10273.1	2992.350000000003
lovely	8704.379999999952	11939.060000000045	3234.680000000093
bpw.style	11572.150000001699	14837.440000000812	3265.289999999113
staleks	8519.730000000003	11875.610000000008	3355.8800000000774
freedecor	3421.779999999971	7671.800000000175	4250.020000000204
runail	71539.27999999933	76758.660000000098	5219.380000001649
polarus	6013.720000000003	11371.930000000018	5358.2100000000155
cosmoprofi	8322.810000000007	14536.990000000016	6214.1800000000089
jessnail	26287.839999999916	33345.22999999992	7057.3900000000007
strong	29196.629999999994	38671.269999999924	9474.639999999985
ingarden	23161.390000000138	33566.210000000009	10404.819999999949
lianail	5892.839999999975	16394.240000000245	10501.400000000027
uno	35302.02999999977	51039.749999998035	15737.719999998262
grattol	35445.5400000011	71472.710000000068	36027.169999999576
	474679.0599999623	619509.2399999934	144830.18000003108

Time taken: 70.259 seconds, Fetched: 161 row(s)

```

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 Shopping
> WHERE event_type='purchase'
> AND
> date_format(event_time, 'MM') IN ('10', '11')
> 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;

Query ID = hadoop_20210117155852_282b0369-324c-4c04-91c0-102abc59add0
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_0006)

-----
VERTICES      MODE        STATUS      TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED      2         2           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: 69.69 s
-----
OK
brand  oct_revenue  nov_revenue  sales_difference
ovale  2.54         3.1          0.56
cosima 20.23        20.92999999999993  0.699999999999922
grace  100.92000000000002  102.61000000000001  1.689999999999977
helloganic  0.0         3.1          3.1
skinity 8.88         12.440000000000001  3.5600000000000005
bodyton 1376.339999999974  1380.639999999992  4.3000000000017735
moyou  5.71         10.280000000000001  4.570000000000001
neoleor 43.41        51.7         8.290000000000006
soleo  204.20000000000003  212.5299999999998  8.329999999999501
jaguar 1102.11      1110.6500000000003  8.5400000000000418
tertio 236.16000000000008  245.7999999999978  9.639999999999702
fly    17.14        27.17        10.030000000000001
rasyan 18.7999999999997  28.93999999999994  10.13999999999997

```

deoproce	316.84	329.17000000000001	12.330000000000098
barbie	0.0	12.39	12.39
supertan	50.370000000000001	66.510000000000002	16.140000000000008
treaclemoon	163.36999999999995	181.48999999999995	18.120000000000005
kamill	63.009999999999999	81.490000000000002	18.480000000000032
juno	0.0	21.08	21.08
veraclara	50.109999999999985	71.210000000000001	21.100000000000023
glysolid	69.72999999999998	91.58999999999997	21.86
godefroy	401.22000000000002	425.12000000000006	23.899999999999864
binacil	0.0	24.259999999999998	24.259999999999998
bliizz	38.949999999999996	63.399999999999998	24.449999999999998
profepil	93.360000000000003	118.02000000000005	24.660000000000025
estelare	444.80999999999943	471.87000000000009	27.060000000000148
orly	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
mavala	409.03999999999985	446.32	37.280000000000014
likato	296.0599999999999	340.9699999999999	44.910000000000025
ladykin	125.64999999999999	170.57	44.92
foamie	35.04	80.49	45.449999999999996
elskin	251.090000000000057	307.650000000000055	56.559999999999974
balbcare	155.32999999999996	212.380000000000025	57.0500000000000296
koelcia	55.5	112.75000000000003	57.250000000000003
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	246.49999999999991	312.52	66.020000000000007
cutrin	299.36999999999995	367.62	68.250000000000006
egomania	77.47	146.040000000000002	68.570000000000002
konad	739.82999999999991	810.67000000000003	70.8400000000000117
nirvel	163.03999999999996	234.32999999999984	71.289999999999988
koelf	422.72999999999985	507.29000000000002	84.560000000000034
plazan	101.37	194.010000000000002	92.640000000000001
aura	83.95	177.51	93.559999999999999
kerasys	430.90999999999985	525.20000000000002	94.290000000000003
enjoy	41.349999999999994	136.570000000000002	95.220000000000003
depilflax	2707.0699999999994	2803.7799999999975	96.710000000000367
eos	54.339999999999996	152.61	98.270000000000001
carmex	145.08	243.36	98.28
batiste	772.3999999999999	874.1699999999994	101.769999999999953

osmo	645.58	762.310000000000002	116.730000000000013
dizao	819.130000000000012	945.5099999999998	126.379999999999852
igrobeauty	513.660000000000009	645.0699999999999	131.409999999999906
finilah	98.38	230.380000000000008	132.000000000000009
nefertiti	233.520000000000007	366.64	133.11999999999992
elizavecca	70.53	204.3	133.77
maskin	158.04	293.070000000000005	135.030000000000006
latinoil	249.52	384.59	135.06999999999996
farmona	1692.4599999999996	1843.43000000000007	150.970000000000116
crystalinas	427.6299999999999	584.9499999999999	157.319999999999914
chi	358.940000000000002	538.61000000000002	179.670000000000002
matreshka	0.0	182.670000000000002	182.670000000000002
freshbubble	318.700000000000001	502.340000000000015	183.640000000000004
mane	66.78999999999999	260.26	193.47
keen	236.350000000000005	435.62	199.26999999999995
ecocraft	41.160000000000004	241.95	200.79
fedua	52.38	263.810000000000006	211.430000000000006
provoc	827.99000000000009	1063.82000000000006	235.82999999999997
skinlite	651.940000000000002	890.4499999999979	238.509999999999772
entity	479.710000000000015	719.2599999999993	239.54999999999978
trind	298.070000000000005	542.96000000000002	244.890000000000001
protokeratin	201.250000000000003	456.790000000000013	255.540000000000001
beauugreen	511.5099999999999	768.35	256.840000000000015
bluesky	10307.2399999999858	10565.5299999999713	258.289999999985535
candy	534.9599999999999	799.3799999999993	264.4199999999994
insight	1443.70000000000012	1721.96000000000003	278.25999999999991
kocostar	310.850000000000001	594.93000000000003	284.080000000000002
happyfons	801.92000000000006	1091.59000000000001	289.66999999999995
kims	330.03999999999996	632.04000000000001	302.000000000000001
shary	871.9599999999994	1176.4899999999989	304.52999999999995
nitrile	847.2799999999999	1162.679999999999	315.4
lowence	242.84	567.7499999999997	324.9099999999996
jas	3318.9599999999995	3657.43000000000026	338.470000000000753
ellips	245.8499999999999	606.0399999999996	360.1899999999997
lador	2083.6100000000004	2471.5300000000007	387.920000000000028
naomi	0.0	389.0	389.0
kiss	421.54999999999944	817.3299999999984	395.7799999999999
yu-r	271.41	673.7099999999998	402.2999999999998
sophin	1067.86000000000001	1515.52000000000011	447.660000000000001
farmavita	837.36999999999984	1291.97000000000003	454.600000000000184
bloaqua	942.8899999999996	1398.1199999999997	455.23
greymy	29.21	489.49	460.280000000000003
gehwoi	1089.07	1557.6799999999982	468.60999999999983

```
hadoop@ip-172-31-94-188~
```

matrix	3243.249999999999	3726.7400000000007	483.49000000000016
limoni	1308.9000000000003	1796.5999999999997	487.69999999999936
s.care	412.68	913.0699999999999	500.38999999999993
coifin	903.0000000000001	1428.4899999999998	525.4899999999997
uskuai	5142.2700000000017	5690.3100000000005	548.03999999999881
airmails	5118.8999999999939	5691.519999999996	572.620000000000572
browxenna	14331.369999999995	14916.729999999976	585.3600000000026
kinetics	6334.2499999999945	6945.2600000000017	611.0100000000022
kosmekka	1181.4400000000003	1813.37	631.9299999999996
kaaral	4412.4299999999985	5086.069999999992	673.6399999999994
refectocil	2716.1800000000005	3475.5800000000007	759.40000000000024
rosi	3077.0399999999927	3841.5600000000013	764.52000000000204
solomeya	1899.6999999999992	2685.799999999991	786.0999999999999
missha	1293.8299999999995	2150.2799999999984	856.4499999999989
levissime	2227.5000000000064	3085.3099999999977	857.8099999999913
art-visage	2092.710000000001	2997.8000000000011	905.0900000000001
ecolab	262.8500000000001	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	1083.7099999999818
de.lux	1659.6999999999967	2775.5099999999968	1115.8100000000009
swarovski	1887.9299999999983	3043.1600000000003	1155.23000000000157
beauty-free	554.1700000000006	1782.86000000000163	1228.69000000000155
zeitun	708.6600000000004	2009.63	1300.9699999999998
joico	705.52	2015.1000000000015	1309.58000000000015
severina	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	1420.5399999999959
roubloff	3491.3600000000003	4913.7699999999991	1422.40999999999885
smart	4457.2600000000004	5902.1400000000017	1444.88000000000128
shik	3341.2	4839.7200000000007	1498.52000000000068
domix	10472.049999999994	12009.1700000000022	1537.120000000000827
artex	2730.6399999999998	4327.2500000000017	1596.61000000000192
beautix	10493.949999999966	12222.949999999913	1728.9999999999472
milv	3904.9399999999964	5642.010000000008	1737.07000000000838
masura	31266.079999999821	33058.469999999708	1792.38999999988753
f.o.x	6624.2299999999982	8577.2800000000004	1953.0500000000022
kapous	11927.1599999999898	14093.0800000000158	2165.9200000000026
concept	11032.1399999999925	13380.399999999993	2348.26000000000057
estel	21756.7500000000342	24142.6700000000022	2385.9199999999878
kaypro	881.3399999999998	3268.699999999995	2387.3599999999995
benovy	409.62000000000002	3259.9700000000001	2850.35000000000001
italwax	21940.2399999999732	24799.3699999999893	2859.13000000000161
yoko	8756.9099999999949	11707.8799999999996	2950.97000000000466
haryuama	9390.6899999999991	12352.910000000013	2962.220000000001394
marathon	7280.7499999999997	10273.1	2992.3500000000003
lovely	8704.3799999999952	11939.0600000000045	3234.68000000000093
bpw.style	11572.1500000001699	14837.4400000000812	3265.2899999999113
staleks	8519.7300000000003	11875.610000000008	3355.88000000000774
freedecor	3421.7799999999971	7671.8000000000175	4250.02000000000204
runail	71539.279999999933	76758.660000000098	5219.3800000001649
polarus	6013.7200000000003	11371.9300000000018	5358.21000000000155
cosmoprofi	8322.8100000000007	14536.9900000000016	6214.18000000000089
jessnail	26287.8399999999916	33345.229999999992	7057.39000000000007
strong	29196.629999999994	38671.2699999999924	9474.6399999999985
ingarden	23161.3900000000138	33566.210000000009	10404.8199999999949
lianail	5892.8399999999975	16394.2400000000245	10501.4000000000027
uno	35302.029999999977	51039.7499999998035	15737.7199999998262
grattol	35445.54000000011	71472.710000000068	36027.1699999999576
	474679.05999999623	619509.23999999934	144830.180000003108

Time taken: 70.259 seconds, Fetched: 161 row(s)

```
hive>
```

## Insights:

- Here are some **161 brands with increment** in the selling from October to November.
- **‘Grattol’ brand has the highest total increment i.e., 36,027 /-** and **‘Ovale’ seems to have least increment of 0.56 /-** from October to November.
- Among all these brands list, **‘Runail’** which was the best brand in terms of selling in October and November combined **is also in the top 10 brands with high increment for October (71539.28 /-) to November (76758.61 /-) i.e., increment of total 5219.38 /-.**
- This implies that **‘Runail’ is the best and popular brand among all other brands within people.**

**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 Shopping
WHERE event_type='purchase'
GROUP BY user_id
ORDER BY Total_Expenditure DESC
LIMIT 10;
```

**Output:**

Query ID = hadoop\_20210117161116\_a5fd0524-a0de-4ac7-9013-121790c67e18

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\_1610894517504\_0007)

```
-----
      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
Reducer 3 ..... container SUCCEEDED   1     1     0     0     0     0
-----
VERTICES: 03/03 [=====>>] 100% ELAPSED TIME: 60.76 s
-----
```

OK

user\_id total\_expenditure

557790271	2715.8699999999991
150318419	1645.97
562167663	1352.8500000000004
531900924	1329.4500000000003
557850743	1295.4800000000002
522130011	1185.3899999999994
561592095	1109.6999999999996
431950134	1097.5899999999995
566576008	1056.3600000000017
521347209	1040.9099999999999

Time taken: 69.753 seconds, Fetched: 10 row(s)

```
hadoop@ip-172-31-94-188:~$
hive> SELECT user_id, SUM(price) as Total_Expenditure
> FROM Shopping
> WHERE event_type='purchase'
> GROUP BY user_id
> ORDER BY Total_Expenditure DESC
> LIMIT 10;

Query ID = hadoop_20210117161116_a5fd0524-a0de-4ac7-9013-121790c67e18
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_1610894517504_0007)

-----
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
Reducer 3 ..... container  SUCCEEDED   1       1         0         0         0         0
-----
VERTICES: 03/03  [=====>>>] 100% ELAPSED TIME: 60.76 s
-----

OK
user_id total_expenditure
557790271      2715.8699999999991
150318419      1645.97
562167663      1352.8500000000004
531900924      1329.4500000000003
557850743      1295.4800000000002
522130011      1185.3899999999994
561592095      1109.6999999999996
431950134      1097.5899999999995
566576008      1056.3600000000017
521347209      1040.9099999999999
Time taken: 69.753 seconds, Fetched: 10 row(s)
hive>
```

### Insights:

- Here is the list of the top 10 users or buyers who have spend the most and could be rewarded with a Golden Customer plan to attract more people in the coming future.
- We are **selecting this query to be executed using Optimized table** to check that does optimized table reduces execution time with proper partitioning and bucketing.
- **Time taken to execute this query on Base table (non-optimized table) is 69.753 seconds.**

## Optimized Table

**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;

A screenshot of a terminal window with a dark background. The window title is 'hadoop@ip-172-31-94-188:~'. The terminal shows three lines of Hive commands being entered at the 'hive>' prompt: 'set hive.exec.dynamic.partition.mode=nonstrict;', 'set hive.exec.dynamic.partition=true;', and 'set hive.enforce.bucketing=true;'. The prompt 'hive>' is visible on the line following the last command.

```
hadoop@ip-172-31-94-188:~  
hive> set hive.exec.dynamic.partition.mode=nonstrict;  
hive> set hive.exec.dynamic.partition=true;  
hive> set hive.enforce.bucketing=true;  
hive>
```

### **Table optimization steps:-**

**1. Command to create table 'Dyn\_Part\_Buck\_Shopping' with partition on 'event\_type' attribute and bucket(cluster) on 'price' attribute.**

#### **Query:**

```
CREATE TABLE IF NOT EXISTS Dyn_Part_Buck_Shopping(  
    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;
```

#### **Output:**

OK



Time taken: 0.159 seconds

```
hadoop@ip-172-31-94-188:~$
hive> CREATE TABLE IF NOT EXISTS Dyn_Part_Buck_Shopping(
> 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: 0.159 seconds
hive>
```

## 2. To add data into partitioned and bucketed table we need to get it from already created table i.e., 'Shopping'

### Query:

INSERT INTO TABLE Dyn\_Part\_Buck\_Shopping

PARTITION (event\_type)

SELECT event\_time, product\_id, category\_id, category\_code, brand, price, user\_id, user\_session,  
event\_type

FROM Shopping;

### Output:

Query ID = hadoop\_20210117162425\_57023bb0-e16e-4665-8c81-ab7f87859fd7

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_0011)

```
-----
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: 163.41 s
-----
```

Loading data to table default.dyn\_part\_buck\_shopping partition (event\_type=null)

Loaded : 4/4 partitions.

Time taken to load dynamic partitions: 0.697 seconds

Time taken for adding to write entity : 0.003 seconds

OK

Time taken: 170.452 seconds

```
hive> INSERT INTO TABLE Dyn_Part_Buck_Shopping
> PARTITION (event_type)
> SELECT event_time, product_id, category_id, category_code, brand, price, user_id, user_session, event_type
> FROM Shopping;
Query ID = hadoop_20210117162425_57023bb0-e16e-4665-8c81-ab7f87859fd7
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_0011)

-----
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: 163.41 s
-----
Loading data to table default.dyn_part_buck_shopping partition (event_type=null)

Loaded : 4/4 partitions.
Time taken to load dynamic partitions: 0.697 seconds
Time taken for adding to write entity : 0.003 seconds
OK
Time taken: 170.452 seconds
hive>
```

**3. 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 commands**

### 3.1. Command to exit Hive environment

- EXIT;

```
hive> EXIT;
```

### 3.2. Command to check successful existence of Partitioned and Bucketed table 'Dyn\_Part\_Buck\_Shopping' in hive warehouse.

- hadoop fs -ls /user/hive/warehouse/Dyn\_Part\_Buck\_Shopping

**Output:**

Fount 4 items

```
drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart

drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase

drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart

drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view
```

```
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping
Found 4 items
drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart
drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase
drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart
drwxrwxrwt - hadoop hadoop      0 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase
Found 7 items
```

### 3.3. Command to check existence of partitions (event\_type = purchase) in the table

```
hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase
```

#### Output:

```
Found 7 items
```

```
-rwxrwxrwt 1 hadoop hadoop 13052654 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000000_0

-rwxrwxrwt 1 hadoop hadoop 9399111 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000001_0

-rwxrwxrwt 1 hadoop hadoop 12636711 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000002_0

-rwxrwxrwt 1 hadoop hadoop 10650131 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000003_0

-rwxrwxrwt 1 hadoop hadoop 7226455 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000004_0

-rwxrwxrwt 1 hadoop hadoop 10737803 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000005_0
```

```
-rwxrwxrwt 1 hadoop hadoop 7825305 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000006_0
```

```
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase
Found 7 items
-rwxrwxrwt 1 hadoop hadoop 13052654 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000000_0
-rwxrwxrwt 1 hadoop hadoop 9399111 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000001_0
-rwxrwxrwt 1 hadoop hadoop 12636711 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000002_0
-rwxrwxrwt 1 hadoop hadoop 10650131 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000003_0
-rwxrwxrwt 1 hadoop hadoop 7226455 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000004_0
-rwxrwxrwt 1 hadoop hadoop 10737803 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000005_0
-rwxrwxrwt 1 hadoop hadoop 7825305 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=purchase/000006_0
```

### 3.4. Command to check existence of partitions (event\_type = cart) in the table

```
hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart
```

#### Output:

```
Found 7 items
```

```
-rwxrwxrwt 1 hadoop hadoop 57724286 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000000_0
```

```
-rwxrwxrwt 1 hadoop hadoop 43094161 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000001_0
```

```
-rwxrwxrwt 1 hadoop hadoop 56823661 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000002_0
```

```
-rwxrwxrwt 1 hadoop hadoop 49030059 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000003_0
```

```
-rwxrwxrwt 1 hadoop hadoop 31050141 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000004_0
```

```
-rwxrwxrwt 1 hadoop hadoop 48253679 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000005_0
```

```
-rwxrwxrwt 1 hadoop hadoop 34272441 2021-01-17 16:26
```

```
/user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000006_0
```

```
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart
Found 7 items
-rwxrwxrwt 1 hadoop hadoop 57724286 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000000_0
-rwxrwxrwt 1 hadoop hadoop 43094161 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000001_0
-rwxrwxrwt 1 hadoop hadoop 56823661 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000002_0
-rwxrwxrwt 1 hadoop hadoop 49030059 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000003_0
-rwxrwxrwt 1 hadoop hadoop 31050141 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000004_0
-rwxrwxrwt 1 hadoop hadoop 48253679 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000005_0
-rwxrwxrwt 1 hadoop hadoop 34272441 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=cart/000006_0
```

### 3.5. Command to check existence of partitions (event\_type = remove\_from\_cart) in the table

```
hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart
```

#### Output:

Found 7 items

```
-rwxrwxrwt 1 hadoop hadoop 39017824 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000000_0

-rwxrwxrwt 1 hadoop hadoop 29421828 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000001_0

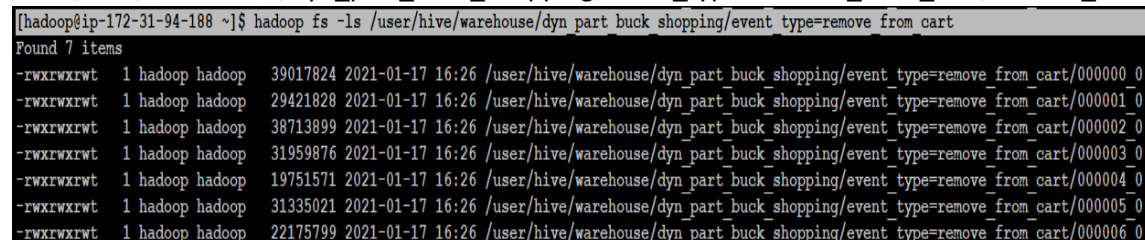
-rwxrwxrwt 1 hadoop hadoop 38713899 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000002_0

-rwxrwxrwt 1 hadoop hadoop 31959876 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000003_0

-rwxrwxrwt 1 hadoop hadoop 19751571 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000004_0

-rwxrwxrwt 1 hadoop hadoop 31335021 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000005_0

-rwxrwxrwt 1 hadoop hadoop 22175799 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000006_0
```



```
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart
Found 7 items
-rwxrwxrwt 1 hadoop hadoop 39017824 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000000_0
-rwxrwxrwt 1 hadoop hadoop 29421828 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000001_0
-rwxrwxrwt 1 hadoop hadoop 38713899 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000002_0
-rwxrwxrwt 1 hadoop hadoop 31959876 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000003_0
-rwxrwxrwt 1 hadoop hadoop 19751571 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000004_0
-rwxrwxrwt 1 hadoop hadoop 31335021 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000005_0
-rwxrwxrwt 1 hadoop hadoop 22175799 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=remove_from_cart/000006_0
```

### 3.6. Command to check existence of partitions (event\_type = view) in the table

```
hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=view
```

#### Output:

Found 2 items

```
-rwxrwxrwt 1 hadoop hadoop 88831872 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000000_0

-rwxrwxrwt 1 hadoop hadoop 73953212 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000001_0
```

```
-rwxrwxrwt 1 hadoop hadoop 85620113 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000002_0

-rwxrwxrwt 1 hadoop hadoop 71874121 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000003_0

-rwxrwxrwt 1 hadoop hadoop 48335545 2021-01-17 16:26
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000004_0

-rwxrwxrwt 1 hadoop hadoop 72515614 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000005_0

-rwxrwxrwt 1 hadoop hadoop 56694677 2021-01-17 16:27
/user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000006_0
```

```
[hadoop@ip-172-31-94-188 ~]$ hadoop fs -ls /user/hive/warehouse/dyn_part_buck_shopping/event_type=view
Found 7 items
-rwxrwxrwt 1 hadoop hadoop 88831872 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000000_0
-rwxrwxrwt 1 hadoop hadoop 73953212 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000001_0
-rwxrwxrwt 1 hadoop hadoop 85620113 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000002_0
-rwxrwxrwt 1 hadoop hadoop 71874121 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000003_0
-rwxrwxrwt 1 hadoop hadoop 48335545 2021-01-17 16:26 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000004_0
-rwxrwxrwt 1 hadoop hadoop 72515614 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000005_0
-rwxrwxrwt 1 hadoop hadoop 56694677 2021-01-17 16:27 /user/hive/warehouse/dyn_part_buck_shopping/event_type=view/000006_0
[hadoop@ip-172-31-94-188 ~]$
```

**4. Now we need to re-enter the Hive environment to execute Query No 8 which we have selected to run on Optimized table.**

- hive

**5. Running the same query for Question 8 on Optimized as executed on Base table to understand the execution time of same query on Base table and Optimized table.**

**(Optimized) 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 Dyn_Part_Buck_Shopping
```

```
WHERE event_type='purchase'

GROUP BY user_id

ORDER BY Total_Expenditure DESC

LIMIT 10;
```

**Output:**

Query ID = hadoop\_20210117164116\_05c7be3c-12d0-479f-8890-fd815730dff6

Total jobs = 1

Launching Job 1 out of 1

Status: Running (Executing on YARN cluster with App id application\_1610894517504\_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
Reducer 3 ..... container  SUCCEEDED   1     1     0     0     0     0
-----
```

VERTICES: 03/03 [=====>>] 100% ELAPSED TIME: 26.83 s

OK

user_id	total_expenditure
557790271	2715.8699999999996
150318419	1645.97
562167663	1352.8500000000001
531900924	1329.4500000000003
557850743	1295.4800000000005
522130011	1185.3899999999999
561592095	1109.7

431950134      1097.59000000000001

566576008      1056.36000000000006

521347209      1040.91000000000003

Time taken: 27.634 seconds, Fetched: 10 row(s)

```
hadoop@ip-172-31-94-188:~$
hive> SELECT user_id, SUM(price) AS Total_Expenditure
> FROM Dyn_Part_Buck_Shopping
> WHERE event_type='purchase'
> GROUP BY user_id
> ORDER BY Total_Expenditure DESC
> LIMIT 10;

Query ID = hadoop_20210117164116_05c7be3c-12d0-479f-8890-fd815730dff6
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1610894517504_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
Reducer 3 ..... container  SUCCEEDED      1        1          0         0         0         0
-----
VERTICES: 03/03  [=====>>>] 100% ELAPSED TIME: 26.83 s
-----
OK
user_id total_expenditure
557790271      2715.8699999999996
150318419      1645.97
562167663      1352.8500000000001
531900924      1329.4500000000003
557850743      1295.4800000000005
522130011      1185.3899999999999
561592095      1109.7
431950134      1097.5900000000001
566576008      1056.3600000000006
521347209      1040.9100000000003
Time taken: 27.634 seconds, Fetched: 10 row(s)
hive>
```

## Insights:

- After creating an optimized table by **Partitioning on 'event\_type'** attribute and **Bucketing (Clustering) on 'price'** we have executed same query of Question No. 8 on this table.
- We can the result is same as we have got when executed on Base table (Non-Optimized table).
- Secondly, most importantly we can see there is significant drop in the execution time of the same query i.e., **previously the execution was measured as 69.753 seconds and now it is 27.634 seconds with the difference of 42.119 seconds.**
- **Hence, with proper partitioning and bucketing on table we can reduce execution time of the query.**



# Terminating EMR Cluster (Hive\_Case\_Study)

The screenshot displays the AWS Management Console interface for an Amazon EMR cluster. The browser address bar shows the URL: `console.aws.amazon.com/elasticmapreduce/home?region=us-east-1#cluster-details:j-3B5QOU4EG0966`. The left-hand navigation pane lists various services under 'Amazon EMR', including 'EMR on EC2', 'Clusters', 'Notebooks', 'Git repositories', 'Security configurations', 'Block public access', 'VPC subnets', 'Events', 'EMR on EKS', and 'Virtual clusters'. The main content area is titled 'Cluster: Hive\_Case\_Study' and indicates the cluster is 'Terminated' with the reason 'Terminated by user request'. Above the cluster name are buttons for 'Clone', 'Terminate', and 'AWS CLI export'. Below the name is a tabbed interface with 'Summary' selected, followed by 'Application user interfaces', 'Monitoring', 'Hardware', 'Configurations', 'Events', 'Steps', and 'Bootstrap actions'. The 'Summary' tab displays the following information:

- ID:** j-3B5QOU4EG0966
- Creation date:** 2021-01-17 22:58 (UTC+5:30)
- End date:** 2021-01-18 00:44 (UTC+5:30)
- Elapsed time:** 1 hour, 45 minutes
- After last step completes:** Cluster waits
- Termination protection:** Off
- Tags:** --
- Master public DNS:** ec2-52-23-176-82.compute-1.amazonaws.com [Connect to the Master Node Using SSH](#)

The 'Configuration details' section shows:

- 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-060198243535-us-east-1/elasticmapreduce/
- EMRFS consistent view:** Disabled
- Custom AMI ID:** --

The footer of the console includes 'Feedback', 'English (US)', and copyright information: '© 2008 - 2021, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved.' along with links to 'Privacy Policy' and 'Terms of Use'.