

Step1: Copying the data set into the HDFS

A EMR cluster has been created with the same configuration as done for previous case study.
PFB screenshot of the cluster.

The screenshot shows the AWS EMR Cluster Details page for 'C3 Assignment'. The cluster is in a 'Waiting' state. Key details include:

- Cluster ID:** j-E0W46K07IM65
- Creation date:** 2020-10-03 15:12 (UTC+5:30)
- Elapsed time:** 16 minutes
- After last step completes:** Cluster waits
- Termination protection:** Off
- Tags:** -- View All / Edit
- Master public DNS:** ec2-54-90-223-201.compute-1.amazonaws.com
- Log URL:** s3://aws-logs-779358638188-us-east-1/elasticmapreduce/
- EMRFS consistent view:** Disabled
- Custom AMI ID:** --

Application user interfaces: Persistent user interfaces: Spark history server. On-cluster user interface: Not Enabled. Enable an SSH Connection.

Network and hardware: Availability zone: us-east-1a. Subnet ID: subnet-4fe7d705. Master: Running 1 m4.large. Core: Running 1 m4.large. Task: --. Cluster scaling: Not enabled.

Security and access: Key name: SDK_keys. EC2 instance profile: EMR_EC2_DefaultRole. EMR role: EMR_DefaultRole. Auto Scaling role: EMR_AutoScaling_DefaultRole. Visible to all users: All. Security groups for Master: sg-0b4a9c78b8c05114f (ElasticMapReduce-master). Security groups for Core & Task: sg-0d98015fd5da89cf (ElasticMapReduce-slave).

Feedback English (US) ▾ © 2003–2020, Amazon Internet Services Private Ltd. or its affiliates. All rights reserved. Privacy Policy Terms of Use

```
Last login: Sat Oct 3 14:39:09 on console
kalaiselvanswamy@kalaiselvanswamy-MBP ~ % ssh -i ~/.ssh/SDK_keys.pem hadoop@ec2-54-90-223-201.compute-1.amazonaws.com
ssh: connect to host ec2-54-90-223-201.compute-1.amazonaws.com port 22: Operation timed out
kalaiselvanswamy@kalaiselvanswamy-MBP ~ % ssh -i ~/.ssh/SDK_keys.pem hadoop@ec2-54-90-223-201.compute-1.amazonaws.com
The authenticity of host 'ec2-54-90-223-201.compute-1.amazonaws.com (54.90.223.201)' can't be established.
ECDSA key fingerprint is SHA256:J9Wbv3QMINA6tWrLMPotjeNqgiOsavunS1sg/z1PQQ.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-54-90-223-201.compute-1.amazonaws.com,54.90.223.201' (ECDSA) to the list of known hosts.
Last login: Sat Oct 3 09:54:20 2020

--| ( --|-
 _| ( --|_ /  Amazon Linux AMI
---\----|---|--

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
44 package(s) needed for security, out of 73 available
Run "sudo yum update" to apply all updates.
-bash: warning: setlocale: LC_CTYPE: cannot change locale (UTF-8): No such file or directory

EEEEEEEEEEEEEEEEEEEEE MMMMMMM      MMMMMMM RRRRRRRRRRRRRRR
E:::;:::;:::;:::;:::;E M::::::::::M      M::::::::::M R:::::::::::R
EE:::::;EEEEE:::::;EE E M::::::::::M      M::::::::::M R::::::::::R
E:::::;E EEEEEE M::::::::::M      M::::::::::M R::::::::::R R:::::;R
E:::::;E M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;EEEEE::::EE M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;EEEEE::::EE M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;EEEEE::::EE M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;E EEEEEE M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;E EEEEEE E M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;EEEEE::::EE E M::::::::::M M::::::::::M R::::::::::R R:::::;R
E:::::;EEEEE::::EE E M::::::::::M M::::::::::M R::::::::::R R:::::;R
EEEEEEEEEEEEEEEEEEEEE MMMMMMM      MMMMMMM RRRRRRR RRRRRRRR
[hadoop@ip-172-31-18-181 ~]$
```

Now let us move the data from S3 to the EMR using hadoop distributed file copy. The command for the same is given below followed by the screenshots of the copy for both files.

```

hadoop distcp s3n://e-commerce-events-ml/2019-Oct.csv
/usr/hadoop/e-commerce-events/2019-Oct.csv
hadoop distcp s3n://e-commerce-events-ml/2019-Nov.csv
/usr/hadoop/e-commerce-events/2019-Nov.csv

```

```

E:::E M::::M M::::M R:::R R:::R
E:::E EEEEE M::::M M::::M R:::R R:::R
E:::EEEEE M::::M M::::M R:::R R:::R
E:::EEEEE E:::M M::::M R:::R R:::R
E:::EEEEE E:::M M::::M RR:::R R:::R
E:::::::E:::::::M M::::::M RRRRRR RRRRRR

[hadoop@ip-172-31-18-181 ~]$ hadoop distcp s3n://e-commerce-events-ml/2019-Oct.csv /user/hadoop/e-commerce-events/2019-Oct.csv
20/10/03 10:02:23 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolders=false, deleteMissing=false, ignoreFailures=false, overwrite=false, skipCRC=false, blocking=true, numListStatusThreads=0, maxMaps=20, mapBandwidth=100, sslComfigure=false, maxFilesize="100MB", preserveRawXattrs=false, atomicWorkPath=null, logPath=null, sourcePaths=[s3n://e-commerce-events-ml/2019-Oct.csv], targetPath=/user/hadoop/e-commerce-events/2019-Oct.csv, targetPathFilter="*"}, filtersFile="null"
20/10/03 10:02:23 INFO client.HMProxy: Connecting to ResourceManager at ip-172-31-18-181.ec2.internal/172.31.18.181:8882
20/10/03 10:02:27 INFO tools.SimpleCopyListing: Paths [files|dirs] cm = 1; dirFn = 0
20/10/03 10:02:27 INFO tools.SimpleCopyListing: file 1 has been completed.
20/10/03 10:02:27 INFO Configuration.deprecation: io.sort.mb is deprecated. Instead, use mapreduce.task.io.sort.mb
20/10/03 10:02:27 INFO Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.factor
20/10/03 10:02:27 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:02:27 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:02:27 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:02:27 INFO mapreduce.JobSubmitter: number of splits:1
20/10/03 10:02:28 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1601718673301_0001
20/10/03 10:02:28 INFO impl.YarnClientImpl: Submitted application application_1601718673301_0001
20/10/03 10:02:28 INFO tools.DistCp: The url to track the job: http://ip-172-31-18-181.ec2.internal:20888/proxy/application_1601718673301_0001/
20/10/03 10:02:28 INFO tools.DistCp: DistCp job-id: job_1601718673301_0001
20/10/03 10:02:28 INFO mapreduce.Job: Running job: job_1601718673301_0001
20/10/03 10:02:28 INFO mapreduce.Job: Job ID : job_1601718673301_0001
20/10/03 10:02:28 INFO mapreduce.Job: map 0% reduce 0%
20/10/03 10:02:56 INFO mapreduce.Job: map 100% reduce 0%
20/10/03 10:02:57 INFO mapreduce.Job: Job job_1601718673301_0001 completed successfully
20/10/03 10:02:57 INFO mapreduce.Job: Counters
File System Counters
FILE: Number of bytes read=0
FILE: Number of bytes written=17886
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=366
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
SN: Number of bytes read=42278
SN: Number of bytes written=0
SN: Number of read operations=0
SN: Number of large read operations=0
SN: Number of write operations=0
Job Counters
Launched map tasks=1
Other local map tasks=1
Total time spent by all map tasks in occupied slots (ms)=10112
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=15941
Total vcore-milliseconds taken by all map tasks=10941
Total vcore-milliseconds taken by all map tasks=16323584
Map-Reduce Framework
Map input records=1
Map output records=0
Input split bytes=16384
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time spent (ms)=278
CPU time spent (ms)=17300
Physical memory (bytes) snapshot=456536272
Virtual memory (bytes) snapshot=3286282816
Total committed heap usage (bytes)=364477312
File Input Format Counters
Bytes Read=224
File Output Format Counters
Bytes Written=0
DistCp Counters
Bytes Copied=482542278
Bytes Expended=482542278
Files Copied=1

[hadoop@ip-172-31-18-181 ~]$

```

```

File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=48254278
  Bytes Expected=48254278
  Files Copied=1
[hadoop@ip-172-31-18-181 ~]$ hadoop distcp s3n://e-commerce-events-ml/2019-Nov.csv /user/hadoop/e-commerce-events/2019-Nov.csv
20/10/03 10:00:18 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, sslConnfigureOnOpen=null, copyStrategy="copy", filterSize=1, preserveStatuses[], preserveRawXAttrs=false, atomicWorkPath=null, logPath=null, sourcePaths=[s3n://e-commerce-events-ml/2019-Nov.csv], targetPath=/user/hadoop/e-commerce-events/2019-Nov.csv}
20/10/03 10:00:18 INFO client.RMProxy: Connecting to ResourceManager at ip-172-31-18-181.ec2.internal/172.31.18.181:8082
20/10/03 10:00:24 INFO tools.SimpleCopyListing: Paths [files+dirs] cmt = 1; dirCnt = 0
20/10/03 10:00:24 INFO tools.SimpleCopyListing: Build file listing completed
20/10/03 10:00:24 INFO configuration.Configuration: Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapreduce.task.io.sort.mb
20/10/03 10:00:24 INFO configuration.Configuration: Configuration.deprecation: io.sort.factor is deprecated. Instead, use mapred.task.io.sort.factor
20/10/03 10:00:24 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:00:24 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:00:24 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:00:24 INFO tools.DistCp: Number of paths in the copy list: 1
20/10/03 10:00:24 INFO madruse.JobSubmitter: number of splits=1
20/10/03 10:00:24 INFO madruse.JobSubmitter: Submitting tokens for job: job_1601718673301_0002
20/10/03 10:00:24 INFO impl.YarnClientImpl: Submitting application application_1601718673301_0002
20/10/03 10:00:24 INFO madruse.JobSubmitter: The url to track the job: http://ip-172-31-18-181.ec2.internal:20888/proxy/application_1601718673301_0002/
20/10/03 10:00:25 INFO madruse.DistCp: DistCp job-id: job_1601718673301_0002
20/10/03 10:00:25 INFO madruse.Job: Running job: job_1601718673301_0002
20/10/03 10:00:25 INFO madruse.Job: Job job_1601718673301_0002 running in uber mode : false
20/10/03 10:00:25 INFO madruse.Job: map 0% reduce 0%
20/10/03 10:00:51 INFO madruse.Job: map 100% reduce 0%
20/10/03 10:00:53 INFO madruse.Job: Job job_1601718673301_0002 completed successfully
20/10/03 10:00:53 INFO madruse.Job: Counters: 38
File System Counters
  FILE: Number of bytes read=0
  FILE: Number of bytes written=17866
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=366
  HDFS: Number of bytes written=546839412
  HDFS: Number of read operations=12
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
  S3N: Number of bytes read=366
  S3N: Number of bytes written=546839412
  S3N: Number of read operations=0
  S3N: Number of large read operations=0
  S3N: Number of write operations=0
Job Counters
  Launched map tasks=1
  Other local map tasks=1
  Total map memory allocated for all maps in occupied slots (mb)=553440
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=17295
  Total wcore-milliseconds taken by all map tasks=17295
  Total wio-milliseconds taken by all map tasks=17710800
Map-Reduce Framework
  Map input records=1
  Map output records=0
  Input split bytes=17866
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time spent (ms)=277
  CPU time spent (ms)=9770
  Physical memory (bytes) snapshot=640827648
  Virtual memory (bytes) snapshot=3294138176
  Total heap capacity used in heap usage (bytes)=486014976
File Input Format Counters
  Bytes Read=224
File Output Format Counters
  Bytes Written=0
DistCp Counters
  Bytes Copied=548539412
  Bytes Expected=548539412
  Files Copied=1
[hadoop@ip-172-31-18-181 ~]#

```

```

[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls e-commerce-events
Found 2 items
-rw-r--r-- 1 hadoop hadoop 545839412 2020-10-03 10:06 e-commerce-events/2019-Nov.csv
-rw-r--r-- 1 hadoop hadoop 48254278 2020-10-03 10:02 e-commerce-events/2019-Oct.csv
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -cat e-commerce-events/2019-Nov.csv | head
event_time,event_type,product_id,category_id,category_code,brand,price,user_id,user_session
2019-11-01 00:00:02 UTC,view,5802432,1487580009286598681,,,0.32,562076640,09fafd6c-6c99-46b1-834f-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,5837166,178399964103198764,,pnb,22,22,556138645,57ed222e-a54a-4907-9944-5a875c2d7f4f
2019-11-01 00:00:11 UTC,cart,5876812,148758001010293687,,jessnail,3.16,564586666,186c1951-8052-4b37-adce-dd9644b1d5f7
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,3.33,553329724,2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:24 UTC,remove_from_cart,5826182,1487580007483048900,,3.33,553329724,2067216c-31b5-455d-a1cc-af0575a34ffb
2019-11-01 00:00:25 UTC,view,5856189,148758000902651821,,runail,15.71,562076640,09fafd6c-6c99-46b1-834f-33527f4de241
2019-11-01 00:00:32 UTC,view,5837835,1933472286753424063,,3.49,514649199,432a4e95-375c-4b40-bd36-0fc039e77580
2019-11-01 00:00:34 UTC,remove_from_cart,5870838,148758000767598893,,milv,0.79,429913900,2f0bbf3c-252f-4fe6-afcd-5d8a6a92839a
cat: Unable to write to output stream.
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -cat e-commerce-events/2019-Oct.csv | head
event_time,event_type,product_id,category_id,category_code,brand,price,user_id,user_session
2019-10-01 00:00:00 UTC,cart,5773203,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:03 UTC,cart,5773353,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:07 UTC,cart,5881589,2151191071051219817,,lovely,13.48,429681830,49e8d843-adf3-428b-a2c3-fe8bc6a307c9
2019-10-01 00:00:07 UTC,cart,5723490,1487580005134238553,,runail,2.62,463240011,26dd6e6e-4dac-4778-8d2c-92e149dab885
2019-10-01 00:00:15 UTC,cart,5881449,1487580013522845895,,lovely,0.56,429681830,49e8d843-adf3-428b-a2c3-fe8bc6a307c9
2019-10-01 00:00:16 UTC,cart,5857269,1487580005134238553,,runail,2.62,430174032,73dea1e7-664e-43f4-8b30-d32b9d5af04f
2019-10-01 00:00:19 UTC,cart,5739055,1487580008246412266,,kapous,4.75,377667011,81326ac6-daa4-4f0a-b488-fd0956a78733
2019-10-01 00:00:24 UTC,cart,5825598,1487580009445982239,,0.56,467916806,2f5b5546-b8cb-9ee7-7ecd-84276f8ef486
2019-10-01 00:00:25 UTC,cart,5698989,1487580006317032337,,1.27,385985999,d30965e8-1101-44ab-b45d-cc1bb9fae694
cat: llnable to write to output stream.

```

Step 2: Creating the database and launching Hive queries on your EMR cluster

Let us first create a database with some comments.

```
create database if not exists sks_c3_assignment comment "Kalaiselvan  
C3 assignment Database";  
use sks_c3_assignment;
```

```
[hadoop@ip-172-31-18-181 ~]$ hive  
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false  
hive> create database if not exists sks_c3_assignment comment "Kalaiselvan C3 assignment Database";  
OK  
Time taken: 0.752 seconds  
hive> use sks_c3_assignment;  
OK  
Time taken: 0.051 seconds  
hive> show databases;  
OK  
default  
sks_c3_assignment  
Time taken: 0.178 seconds, Fetched: 2 row(s)  
hive> █
```

Set Hive CLI headers:

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

Create table for CSV delimited by ',' and lines separated by newline. This is an ordinary table without any special storage format and without any partition.

NOTE: I am having event_time as string as there are string values like UTC in the file as shown above.

```
create table if not exists ecom_events_2019 (event_time string,  
event_type string, product_id string, category_id string,  
category_code string, brand string, price float, user_id bigint,  
user_session string) row format delimited fields terminated by ','  
lines terminated by '\n' tblproperties("skip.header.line.count"="1");
```

```
hive> set hive.cli.print.header=true;  
hive> create table if not exists ecom_events_2019 (event_time string, event_type string, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string) row format delimited fields terminated by ',' lines terminated by '\n' tblproperties("skip.header.line.count"="1");  
OK  
Time taken: 0.634 seconds  
hive> show tables;  
OK  
tab_name  
ecom_events_2019  
Time taken: 0.055 seconds, Fetched: 1 row(s)  
hive> █
```

Now let us load both the data into the same table as there are queries which require data from both files.

```
load data inpath 'e-commerce-events/2019-Oct.csv' into table
ecom_events_2019;
load data inpath 'e-commerce-events/2019-Nov.csv' into table
ecom_events_2019;
```

```
|hive> load data inpath 'e-commerce-events/2019-Oct.csv' into table ecom_events_2019;
Loading data to table sks_c3_assignment.ecom_events_2019
OK
Time taken: 0.87 seconds
|hive> load data inpath 'e-commerce-events/2019-Nov.csv' into table ecom_events_2019;
Loading data to table sks_c3_assignment.ecom_events_2019
OK
Time taken: 0.637 seconds
```

Now let us switch to beeline and we can see the data has been loaded into the Hive table successfully in correct order.

```
|beeline -n -S -b -e | beeline -e://localhost:10000/default -n hadoop
Connecting to: jdbc:hive2://localhost:10000/default
Connected to: Apache Hive (version 2.3.6-0.1)
Driver: Hive JDBC (version 2.3.6-0.1)
Transferred bytes: 0
Transactions: 0
Tables: 0
Partitions: 0
Rows: 0
Memory usage: 0
File: /tmp/beeline-158844418211821.cdc3ce90-aa46-4bd2-8f18-e954cf9c9387
Beeline version 2.3.6-mzn-1 by Apache Hive
| jdbc:hive2://localhost:10000/default> show databases;
+-----+
| database_name |
+-----+
| default      |
| sks_c3_assignment |
+-----+
2 rows selected (1.358 seconds)

| jdbc:hive2://localhost:10000/default> use sks_c3_assignment;
+-----+
| default      |
+-----+
Semantic Analysis Completed
| jdbc:hive2://localhost:10000/default> show tables;
+-----+
| tab_name     |
+-----+
No rows selected (0.000 seconds)
| jdbc:hive2://localhost:10000/default> use sks_c3_assignment;
+-----+
| default      |
+-----+
Semantic Analysis Completed
| jdbc:hive2://localhost:10000/default> show tables;
+-----+
| tab_name     |
+-----+
No rows selected (0.000 seconds)
| jdbc:hive2://localhost:10000/default> show tables;
+-----+
| tab_name     |
+-----+
No rows selected (0.000 seconds)
| jdbc:hive2://localhost:10000/default> select * from ecom_events_2019 limit 5;
+-----+
| ecom_events_2019.event_type | ecom_events_2019.product_id | ecom_events_2019.category_id | ecom_events_2019.category_code | ecom_events_2019.brand | ecom_events_2019.price | ecom_events_2019.user_id | ecom_events_2019.user_session |
+-----+
| 2019-11-01 08:08:02 UTC | view | 5882432 | 148758000728559681 | | 0.32 | 562976648 | 00fa7d0c-cc99-4eb1-834f-3352774ad9241 | |
| 2019-11-01 08:08:09 UTC | view | 58844979 | 148758000651233537 | | 2.38 | 553279724 | 2867216-318a-434a-83cc-4375a2775a7fb |
| 2019-11-01 08:08:11 UTC | view | 58872646 | 148758000108639607 | | 2.22 | 553279724 | 57edc344-4544-4591-a897-5a87707f741 |
| 2019-11-01 08:08:11 UTC | cart | 58876812 | 148758000108295687 | | Jessnall | 3.16 | 664586666 | 18ec1395-0852-4b37-adce-d9944cd1d5f7 |
| 2019-11-01 08:08:12 UTC | remove_from_cart | 5826182 | 148758000748304900 | | Jessnall | 3.33 | 553279724 | 2867216-3150-4056-a1cc-ef00575a34f7b |
+-----+
5 rows selected (2.208 seconds)
0: jdbc:hive2://localhost:10000/default>
```

Now let us look at the questions that need to be answered.

Question 1:

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

Here we need to find the sum of the total price of all products which has an event as purchase and month is October. We can extract the month from the timestamp string and match it to 10th month.

Query:

```
select sum(price) as OctPurchase from ecom_events_2019 where
event_type = 'purchase' and month(event_time) = 10;
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select sum(price) as OctPurchase from ecom_events_2019 where event_type = 'purchase' and month(event_time) = 10;
INFO : Compiling command(queryId=hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885): select sum(price) as OctPurchase from ecom_events_2019 where event_type = 'purchase' and month(event_time) = 10
INFO : Semantic Analysis Completed
INFO : Executing command(queryId=hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885): select sum(price) as OctPurchase from ecom_events_2019 where event_type = 'purchase' and month(event_time) = 10
INFO : Completed compiling command(queryId=hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885); Time taken: 0.890 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885): select sum(price) as OctPurchase from ecom_events_2019 where event_type = 'purchase' and month(event_time) = 10
INFO : Query ID = hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885
INFO : Total job = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Tez session hasn't been created yet. Opening session
INFO : Dag name: select sum(price) as OctPurchase from e..10(Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0004)
INFO : Map 1: -/- Reducer 2: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1
INFO : Map 1: 1(+1)/2 Reducer 2: 0(+1)/1
INFO : Map 1: 2/2 Reducer 2: 0(+1)/1
INFO : Map 1: 2/2 Reducer 2: 0/1
INFO : Map 1: 2/2 Reducer 2: 0/1
INFO : Map 1: 0/1
INFO : Completed executing command(queryId=hive_20201003102601_8d64140a-56df-445c-a812-b441d8d9d885); Time taken: 32.612 seconds
INFO : OK
+-----+
| octpurchase |
+-----+
| 1211538.4295325726 |
+-----+
1 row selected (33.634 seconds)
0: jdbc:hive2://localhost:10000/default>
```

Partitioning:

As we can see most of the queries so far and also the forthcoming queries are focussed on the purchase event_type. So let us partition the data which has been purchased.

```
create table if not exists purchases_2019 (event_time string,
product_id string, category_id string, category_code string, brand
string, price float, user_id bigint, user_session string) partitioned
by (event_type string) row format delimited fields terminated by ','
lines terminated by '\n';
```

```
0: jdbc:hive2://localhost:10000/default> create table if not exists purchases_2019 (event_time string, product_id string, category_id string, category_code string, brand string, price float, user_id
bigint, user_session string) partitioned by (event_type string) row format delimited fields terminated by ',' lines terminated by '\n';
INFO : Compiling command(queryId=hive_20201003104012_1f7a3bf9-44ea-4546-a3c3-9d7e9c8203d0): create table if not exists purchases_2019 (event_time string, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string) partitioned by (event_type string) row format delimited fields terminated by ',' lines terminated by '\n'
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
INFO : Completed compiling command(queryId=hive_20201003104012_1f7a3bf9-44ea-4546-a3c3-9d7e9c8203d0); Time taken: 0.018 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003104012_1f7a3bf9-44ea-4546-a3c3-9d7e9c8203d0): create table if not exists purchases_2019 (event_time string, product_id string, category_id string, category_code string, brand string, price float, user_id bigint, user_session string) partitioned by (event_type string) row format delimited fields terminated by ',' lines terminated by '\n'
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20201003104012_1f7a3bf9-44ea-4546-a3c3-9d7e9c8203d0); Time taken: 0.147 seconds
INFO : OK
No rows affected (0.183 seconds)
```

```
insert into table purchases_2019 partition (event_type = 'purchase')
select event_time, product_id, category_id, category_code, brand,
```


Question 2:

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

Here since we need to show data/month let us group by based on month.

Query:

```
select sum(price) as PurchaseTotal, month(event_time) as PurchaseMonth from purchases_2019 group by month(event_time);
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select sum(price) as PurchaseTotal, month(event_time) as PurchaseMonth from purchases_2019 group by month(event_time);
INFO : Compiling command(queryId=hive_20201003105939_c3723880-93aa-4945-ab34-4bd9281b22c7); select sum(price) as PurchaseTotal, month(event_time) as PurchaseMonth from purchases_2019 group by month(event_time)
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema fieldSchemas:[FieldSchema(name:purchasetotal, type:double, comment:null), FieldSchema(name:purchasemonth, type:int, comment:null)], properties:null
INFO : Completed compiling command(queryId=hive_20201003105939_c3723880-93aa-4945-ab34-4bd9281b22c7); Time taken: 0.137 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003105939_c3723880-93aa-4945-ab34-4bd9281b22c7): select sum(price) as PurchaseTotal, month(event_time) as PurchaseMonth from purchases_2019 group by month(event_time)
INFO : Query ID = hive_20201003105939_c3723880-93aa-4945-ab34-4bd9281b22c7
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select sum(price) as Pur...month(event_time)(Stage=1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0007)
INFO : Map 1: 0/2 Reducer 2: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1
INFO : Map 1: 1(+1)/2 Reducer 2: 0/1
INFO : Map 1: 1(+1)/2 Reducer 2: 0/1
INFO : Map 1: 2/2 Reducer 2: 0(+1)/1
INFO : Map 1: 2/2 Reducer 2: 1/1
INFO : Completed executing command(queryId=hive_20201003105939_c3723880-93aa-4945-ab34-4bd9281b22c7); Time taken: 14.271 seconds
INFO : OK
+-----+-----+
| purchasetotal | purchasemonth |
+-----+-----+
| 1211538.4295325726 | 10 |
| 1531616.8991247676 | 11 |
+-----+
2 rows selected (14.449 seconds)
0: jdbc:hive2://localhost:10000/default>
```

Question 3:

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

Here we need to find the difference in revenue between 2 months. We will use the below subquery with a cartesian product between 2 tables for this purpose.

NOTE: In order to run below query we need to set `hive.strict.checks.cartesian.product` property to false.

Query:

```
select (OctRev.Rev - NovRev.Rev) as RevenueDiff
from
(select sum(price) as Rev from purchases_2019 where
month(event_time) = 10) as OctRev
join
(select sum(price) as Rev from purchases_2019 where
month(event_time) = 11) as NovRev;
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select (OctRev.Rev - NovRev.Rev) as RevenueDiff from
... . . . . . . . . . . . . . . . . . . . . . . . . . . . . . > (select sum(price) as Rev from purchases_2019 where month(event_time) = 10) as OctRev
... . . . . . . . . . . . . . . . . . . . . . . . . . . . . . > join
... . . . . . . . . . . . . . . . . . . . . . . . . . . . . . > (select sum(price) as Rev from purchases_2019 where month(event_time) = 11) as NovRev;
INFO : Compiling command(queryId=hive_20201003115048_81b7c406-7452-4485-8a69-566ebfdf9499): select (OctRev.Rev - NovRev.Rev) as RevenueDiff from
(select sum(price) as Rev from purchases_2019 where month(event_time) = 10) as OctRev
join
(select sum(price) as Rev from purchases_2019 where month(event_time) = 11) as NovRev
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:revvenuediff, type:double, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20201003115048_81b7c406-7452-4485-8a69-566ebfdf9499); Time taken: 0.245 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003115048_81b7c406-7452-4485-8a69-566ebfdf9499): select (OctRev.Rev - NovRev.Rev) as RevenueDiff from
(select sum(price) as Rev from purchases_2019 where month(event_time) = 10) as OctRev
join
(select sum(price) as Rev from purchases_2019 where month(event_time) = 11) as NovRev
INFO : Query ID = hive_20201003115048_81b7c406-7452-4485-8a69-566ebfdf9499
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select (OctRev.Rev - NovRev.Rev) as...NovRev(Stage-1)
INFO : Setting tez.task.scale.memory.reserve-fraction to 0.30000001192092896
INFO : Tez session was closed. Reopening...
INFO : Session re-established.
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0010)

INFO : Map 1: --/-- Map 3: --/-- Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0/2 Map 3: 0/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0/2 Map 3: 0/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0/2 Map 3: 0/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 0(+2)/2 Map 3: 0(+1)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 1(+1)/2 Map 3: 0(+2)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 2/2 Map 3: 0(+2)/2 Reducer 2: 0/1 Reducer 4: 0/1
INFO : Map 1: 2/2 Map 3: 0(+2)/2 Reducer 2: 0(+1)/1 Reducer 4: 0/1
INFO : Map 1: 2/2 Map 3: 1(+1)/2 Reducer 2: 0(+1)/1 Reducer 4: 0(+1)/1
INFO : Map 1: 2/2 Map 3: 2/2 Reducer 2: 0(+1)/1 Reducer 4: 1/1
INFO : Map 1: 2/2 Map 3: 2/2 Reducer 2: 1/1 Reducer 4: 1/1
INFO : Completed executing command(queryId=hive_20201003115048_81b7c406-7452-4485-8a69-566ebfdf9499); Time taken: 31.095 seconds
INFO : OK
+-----+
| revvenuediff |
+-----+
| -319478.469592195 |
+-----+
1 row selected (31.404 seconds)
0: jdbc:hive2://localhost:10000/default>
```

The minus value indicates that the revenue in November was more than the revenue in October and there was an 3,19,478 increase in revenue in November when compared to October

Time taken: 31 seconds

Further Optimizations:

Also we can see that we need data based on a month. So let us create a new table and split the above partition into 2 buckets based on month(event_time).

For this in order to perform bucketing let us create another table ecom_events_month_2019 with month(event_time) as a separate column.

We can see the data populated in the final optimized table parq_buck_purchase.

```
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls /user/hive/warehouse/sks_c3_assignment.db/pq_buck_purchase
Found 4 items
drwxrwxrwt  - hadoop hadoop          0 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/pq_buck_purchase/event_type=cart
drwxrwxrwt  - hadoop hadoop          0 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/pq_buck_purchase/event_type=purchase
drwxrwxrwt  - hadoop hadoop          0 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/pq_buck_purchase/event_type=remove_from_cart
drwxrwxrwt  - hadoop hadoop          0 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/pq_buck_purchase/event_type=view
[hadoop@ip-172-31-18-181 ~]$
```

In the above table we can see 4 partitions created dynamically for each event type. Now let us check the buckets in each partition.

```
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=cart
Found 2 items
-rwxrwxrwt 1 hadoop hadoop 89313368 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=cart/000000_0
-rwxrwxrwt 1 hadoop hadoop 93813729 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=cart/000001_0
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=purchase
Found 2 items
-rwxrwxrwt 1 hadoop hadoop 6106454 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=purchase/000000_0
-rwxrwxrwt 1 hadoop hadoop 10210823 2020-10-03 12:52 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=purchase/000001_0
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=remove_from_cart
Found 2 items
-rwxrwxrwt 1 hadoop hadoop 50038840 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=remove_from_cart/000000_0
-rwxrwxrwt 1 hadoop hadoop 63278928 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=remove_from_cart/000001_0
[hadoop@ip-172-31-18-181 ~]$ hadoop fs -ls /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=view
Found 2 items
-rwxrwxrwt 1 hadoop hadoop 146096722 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=view/000000_0
-rwxrwxrwt 1 hadoop hadoop 161165126 2020-10-03 12:53 /user/hive/warehouse/sks_c3_assignment.db/parq_buck_purchase/event_type=view/000001_0
[hadoop@ip-172-31-18-181 ~]$
```

As you can see above each partition is further divided into 2 buckets based on 2 months.

Now let us try running the 3rd question in the newly optimized table and check the time taken. We need to add condition here since the previous query was run on an already filtered table with event type as purchase.

```
select (OctRev.Rev - NovRev.Rev) as RevenueDiff
from
(select sum(price) as Rev from parq_buck_purchase where
month(event_time) = 10 and event_type = 'purchase') as OctRev
join
(select sum(price) as Rev from parq_buck_purchase where
month(event_time) = 11 and event_type = 'purchase') as NovRev;
```

```
0: jdbc:hive2://localhost:10000/default> select (OctRev.Rev - NovRev.Rev) as RevenueDiff
...> from
...> (select sum(price) as Rev from parq_buck_purchase where month(event_time) = 10 and event_type = 'purchase') as OctRev
...> join
...> (select sum(price) as Rev from parq_buck_purchase where month(event_time) = 11 and event_type = 'purchase') as NovRev
INFO : Compiling command(queryId=hive_20201003130835_cfb8a59c-7462-433f-9075-8afa5002daa4): select (OctRev.Rev - NovRev.Rev) as RevenueDiff
from
(select sum(price) as Rev from parq_buck_purchase where month(event_time) = 10 and event_type = 'purchase') as OctRev
join
(select sum(price) as Rev from parq_buck_purchase where month(event_time) = 11 and event_type = 'purchase') as NovRev
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldschemas:[FieldSchema(name:revenueendiff, type:double, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20201003130835_cfb8a59c-7462-433f-9075-8afa5002daa4); Time taken: 0.684 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003130835_cfb8a59c-7462-433f-9075-8afa5002daa4): select (OctRev.Rev - NovRev.Rev) as RevenueDiff
from
(select sum(price) as Rev from parq_buck_purchase where month(event_time) = 10 and event_type = 'purchase') as OctRev
join
(select sum(price) as Rev from parq_buck_purchase where month(event_time) = 11 and event_type = 'purchase') as NovRev
INFO : Query ID = hive_20201003130835_cfb8a59c-7462-433f-9075-8afa5002daa4
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select (OctRev.Rev - NovRev.Rev) as ...NovRev(Stage-1)
INFO : Setting tez.task.scale.memory.reserve-fraction to 0.3000000192092896
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0014)

INFO : Map 1: 0/1      Map 3: 0/1      Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0/1      Map 3: 0/1      Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0/1
INFO : Map 1: 0(+1)/1  Map 3: 0(+1)/1  Reducer 2: 0/1  Reducer 4: 0(+1)/1
INFO : Map 1: 0/1      Map 3: 1/1      Reducer 2: 0/1  Reducer 4: 0(+1)/1
INFO : Map 1: 0/1      Map 3: 1/1      Reducer 2: 0/1  Reducer 4: 1/1
INFO : Map 1: 1/1      Map 3: 1/1      Reducer 2: 1/1  Reducer 4: 1/1
INFO : Completed executing command(queryId=hive_20201003130835_cfb8a59c-7462-433f-9075-8afa5002daa4); Time taken: 17.143 seconds
INFO : OK
+-----+
| revenueendiff |
+-----+
| -319478.469592195 |
+-----+
1 row selected (17.859 seconds)
0: jdbc:hive2://localhost:10000/default> ■
```

Time taken : 17.86 seconds

As you can see above there is a 13 second difference between both queries even though they fetch the same data. The final optimized table with partitions and buckets stored as parquet fetches data faster than a table with just partitions.

Now let us proceed with answering rest of the questions with our most optimized table `parq_buck_purchase`

Question 4:
Find distinct categories of products.

Query:

```
select distinct(category_code) as categories from parq_buck_purchase;
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select distinct(category_code) as categories from parq_buck_purchase;
INFO : Compiling command[queryId=hive_20201003131437_c879da2c-965a-472b-8939-33b7bfaf87cf]: select distinct(category_code) as categories from parq_buck_purchase
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldsSchemas:[FieldSchema(name:categories, type:string, comment:null)], properties:null)
INFO : Completed compiling command[queryId=hive_20201003131437_c879da2c-965a-472b-8939-33b7bfaf87cf]; Time taken: 0.115 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command[queryId=hive_20201003131437_c879da2c-965a-472b-8939-33b7bfaf87cf]: select distinct(category_code) as categories from parq_buck_purchase
INFO : Query ID = hive_20201003131437_c879da2c-965a-472b-8939-33b7bfaf87cf
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select distinct(categories) from parq_buck_purchase(Stage-1)
INFO : Tez session was closed. Reopening...
INFO : Session re-established.
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0015)

INFO : Map 1: -/- Reducer 2: 0/1
INFO : Map 1: 0/6 Reducer 2: 0/1
INFO : Map 1: 0/6 Reducer 2: 0/1
INFO : Map 1: 0/6 Reducer 2: 0/1
INFO : Map 1: 0(+1)/6 Reducer 2: 0/1
INFO : Map 1: 0(+3)/6 Reducer 2: 0/1
INFO : Map 1: 1(+2)/6 Reducer 2: 0/1
INFO : Map 1: 1(+3)/6 Reducer 2: 0/1
INFO : Map 1: 2(+2)/6 Reducer 2: 0/1
INFO : Map 1: 3(+2)/6 Reducer 2: 0/1
INFO : Map 1: 3(+3)/6 Reducer 2: 0/1
INFO : Map 1: 3(+2)/6 Reducer 2: 0/1
INFO : Map 1: 4(+2)/6 Reducer 2: 0(+1)/1
INFO : Map 1: 6/6 Reducer 2: 0(+1)/1
INFO : Map 1: 6/6 Reducer 2: 1/1
INFO : Completed executing command[queryId=hive_20201003131437_c879da2c-965a-472b-8939-33b7bfaf87cf]; Time taken: 31.051 seconds
INFO : OK
+-----+
|      categories      |
+-----+
| 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.cartrige          |
+-----+
12 rows selected (31.229 seconds)
0: jdbc:hive2://localhost:10000/default>
```

As you can see above there are 11 distinct categories and there are products with no category also.

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

Query:

```
select count(product_id) as ProdCount, category_code from  
parq_buck_purchase group by category_code;
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select count(product_id) as ProdCount, category_code from parq_buck_purchase group by category_code;  
INFO : Compiling command(queryId=hive_20201003131808_97f90673-8d38-44b4-bb29-bfc986b8f277); select count(product_id) as ProdCount, category_code from parq_buck_purchase group by category_code  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:prodcount, type:bigint, comment:null), FieldSchema(name:category_code, type:string, comment:null)], properties:null)  
INFO : Completed compiling command(queryId=hive_20201003131808_97f90673-8d38-44b4-bb29-bfc986b8f277); Time taken: 0.263 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20201003131808_97f90673-8d38-44b4-bb29-bfc986b8f277); select count(product_id) as ProdCount, category_code from parq_buck_purchase group by category_code  
INFO : Total jobs = 1  
INFO : Launching Job 1 out of 1  
INFO : Starting task [Stage-1:MAPRED] in serial mode  
INFO : Session is already open  
INFO : Dag name: select count(product_id) as ...category_code(Stage-1)  
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0015)  
INFO : Map 1: 0/6 Reducer 2: 0/1  
INFO : Map 1: 0/6 Reducer 2: 0/1  
INFO : Map 1: 0/6 Reducer 2: 0/1  
INFO : Map 1: 0(+1)/6 Reducer 2: 0/1  
INFO : Map 1: 0(+2)/6 Reducer 2: 0/1  
INFO : Map 1: 0(+3)/6 Reducer 2: 0/1  
INFO : Map 1: 1(+3)/6 Reducer 2: 0/1  
INFO : Map 1: 2(+2)/6 Reducer 2: 0/1  
INFO : Map 1: 3(+1)/6 Reducer 2: 0/1  
INFO : Map 1: 3(+2)/6 Reducer 2: 0/1  
INFO : Map 1: 3(+3)/6 Reducer 2: 0/1  
INFO : Map 1: 4(+2)/6 Reducer 2: 0/1  
INFO : Map 1: 4(+2)/6 Reducer 2: 0(+1)/1  
INFO : Map 1: 5(+1)/6 Reducer 2: 0(+1)/1  
INFO : Map 1: 6/6 Reducer 2: 0(+1)/1  
INFO : Map 1: 6/6 Reducer 2: 1/1  
INFO : Completed executing command(queryId=hive_20201003131808_97f90673-8d38-44b4-bb29-bfc986b8f277); Time taken: 24.353 seconds  
INFO : OK  
+-----+  
| prodcount | category_code |  
+-----+  
| 8594995 |  
| 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 |  
| 398 | furniture.living_room.chair |  
| 2 | sport.diving |  
| 26722 | stationery.cartidge |  
+-----+  
12 rows selected (24.667 seconds)  
0: jdbc:hive2://localhost:10000/default> ||
```

You can see above the highest number of products have no category assigned.

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

Query:

```
Select brand from (select sum(price) as Sales, brand from  
purchases_2019 group by brand order by Sales desc limit 1) as  
TopBrands;
```

Let us use the purchases_2019 table for the forthcoming queries since we need data only which are purchased.

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales desc limit 1) as TopBrands;
INFO : Compiling command(queryId=hive_20201003140926_3201b806-18ad-4ae3-b55b-58228660bb94); Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales de
sc limit 1) as TopBrands
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:brand, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20201003140926_3201b806-18ad-4ae3-b55b-58228660bb94); Time taken: 0.104 seconds
INFO : Lock mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003140926_3201b806-18ad-4ae3-b55b-58228660bb94); Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales de
sc limit 1) as TopBrands
INFO : Query ID = hive_20201003140926_3201b806-18ad-4ae3-b55b-58228660bb94
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: Select brand from (select sum(pr...TopBrands(Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0019)

INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 1(+1)/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 1/1 Reducer 3: 0/1
INFO : Completed executing command(queryId=hive_20201003140926_3201b806-18ad-4ae3-b55b-58228660bb94); Time taken: 13.676 seconds
INFO : OK
+-----+
| brand |
+-----+
|       |
|       |
|       |
+-----+
1 row selected (13.803 seconds)
0: jdbc:hive2://localhost:10000/default>
```

As you can see above the highest sales happened in products which did not contain any brand. So let us see which is the next highest (i.e.) top among available brands.

```
Select brand from (select sum(price) as Sales, brand from
purchases_2019 group by brand order by Sales desc limit 2) as
TopBrands;
```

```
0: jdbc:hive2://localhost:10000/default> Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales desc limit 2) as TopBrands;
INFO : Compiling command(queryId=hive_20201003140811_02ff6bb1-df06-4b39-a5c3-2aec7737122f); Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales de
sc limit 2) as TopBrands
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:brand, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20201003140811_02ff6bb1-df06-4b39-a5c3-2aec7737122f); Time taken: 0.115 seconds
INFO : Lock mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003140811_02ff6bb1-df06-4b39-a5c3-2aec7737122f); Select brand from (select sum(price) as Sales, brand from purchases_2019 group by brand order by Sales de
sc limit 2) as TopBrands
INFO : Query ID = hive_20201003140811_02ff6bb1-df06-4b39-a5c3-2aec7737122f
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: Select brand from (select sum(pr...TopBrands(Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0019)

INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 1(+1)/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 1/1 Reducer 3: 0(+1)/1
INFO : Completed executing command(queryId=hive_20201003140811_02ff6bb1-df06-4b39-a5c3-2aec7737122f); Time taken: 14.046 seconds
INFO : OK
+-----+
| brand |
+-----+
|       |
|       |
| runail |
+-----+
2 rows selected (14.187 seconds)
0: jdbc:hive2://localhost:10000/default>
```

Runail is the brand with the highest sales amount of all available brands.

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

We need to first get the details of each brand on their revenue for October and similarly for November, then compare which rows have November Total sales greater than October total sales and fetch those brands distinctly.

Here also we use purchases_2019 since we need only purchase data .

Query:

```
select distinct(novdet.brand) as brands
from
(Select sum(price) as OctSales, brand from purchases_2019 group by
month(event_time), brand having month(event_time) = 10) as OctDet
Inner Join
(select sum(price) as NovSales, brand from purchases_2019 group by
month(event_time), brand having month(event_time) = 11) as NovDet
On (OctDet.brand = NovDet.brand)
Where OctDet.OctSales < NovDet.NovSales;
```

As you can see in screenshots below out of a total of 219 brands 153 brands increased their sales from October to November.

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select distinct(novdet.brand) as brands
+-----+-----+
| brands |-----+
| airnails |-----+
| art-visage |-----+
| artex |-----+
| aura |-----+
```

0: jdbc:hive2://localhost:10000/default> select distinct(novdet.brand) as brands
.....> from
.....> (Select sum(price) as OctSales, brand from purchases_2019 group by
.....> month(event_time), brand having month(event_time) = 10) as OctDet
.....> inner join
.....> (select sum(price) as NovSales, brand from purchases_2019 group by
.....> month(event_time), brand having month(event_time) = 11) as NovDet
.....> On (OctDet.brand = NovDet.brand)
.....> Where OctDet.OctSales < NovDet.NovSales;
INFO : Compiling command(queryId=hive_20201003155034_9e96a774-421c-4054-b05a-b598c9ba6ef6): select distinct(novdet.brand) as brands
from
|(Select sum(price) as OctSales, brand from purchases_2019 group by
|month(event_time), brand having month(event_time) = 10) as OctDet
|inner join
|(select sum(price) as NovSales, brand from purchases_2019 group by
|month(event_time), brand having month(event_time) = 11) as NovDet
|On (OctDet.brand = NovDet.brand)
|Where OctDet.OctSales < NovDet.NovSales
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:brands, type:string, comment:null)], properties:null)
INFO : Completed compiling command(queryId=hive_20201003155034_9e96a774-421c-4054-b05a-b598c9ba6ef6); Time taken: 0.167 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003155034_9e96a774-421c-4054-b05a-b598c9ba6ef6): select distinct(novdet.brand) as brands
from
(Select sum(price) as OctSales, brand from purchases_2019 group by
month(event_time), brand having month(event_time) = 10) as OctDet
inner join
(select sum(price) as NovSales, brand from purchases_2019 group by
month(event_time), brand having month(event_time) = 11) as NovDet
On (OctDet.brand = NovDet.brand)
Where OctDet.OctSales < NovDet.NovSales
INFO : Query ID = hive_20201003155034_9e96a774-421c-4054-b05a-b598c9ba6ef6
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select distinct(novdet.br...NovDet.NovSales(Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0023)
INFO : Map 1: 0/2 Map 5: 0/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0/2 Map 5: 0/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0/2 Map 5: 0/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 0(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 0(+1)/2 Map 5: 1(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0/1
INFO : Map 1: 1(+1)/2 Map 5: 1(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0(+1)/1
INFO : Map 1: 1(+1)/2 Map 5: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 0(+1)/1
INFO : Map 1: 1(+1)/2 Map 5: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 1/1
INFO : Map 1: 2/2 Map 5: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1 Reducer 4: 0/1 Reducer 6: 1/1
INFO : Map 1: 2/2 Map 5: 2/2 Reducer 2: 1/1 Reducer 3: 0(+1)/1 Reducer 4: 0/1 Reducer 6: 1/1
INFO : Map 1: 2/2 Map 5: 2/2 Reducer 2: 1/1 Reducer 3: 1/1 Reducer 4: 0(+1)/1 Reducer 6: 1/1
INFO : Map 1: 2/2 Map 5: 2/2 Reducer 2: 1/1 Reducer 3: 1/1 Reducer 4: 1/1 Reducer 6: 1/1
INFO : Completed executing command(queryId=hive_20201003155034_9e96a774-421c-4054-b05a-b598c9ba6ef6); Time taken: 24.47 seconds
INFO : OK

```
| balocore
| batiste
| beautix
| beauty-free
| beautyblender
| beautyskin
| benovy
| biobea
| bior
| blixz
| bluesky
| bodyskin
| bpw-style
| browxenna
| candy
| carnex
| chi
| colfin
| concept
| cessim
| cessomorfi
| cristalinas
| cutrin
| de.lux
| deource
| depilflax
| dizao
| domix
| ecraft
| ecolab
| eppenisa
| ellizavecca
| ellips
| elskin
| enjpy
| entity
| eos
| extrel
| extrelare
| f.o.x
| farronits
| farzona
| fedua
| finnish
| fly
| foanie
| frencor
| freeshubble
| gewol
| glysolid
| godofroy
| grace
| gratal
| gray
| happyfons
| haruyana
| ignobeautey
| ingorden
| inn
| insight
| italk
| italex
| jaguer
| jax
| jessnail
| joico
| kaxai
| kawall
| kapous
| kayoro
| kess
| kersey
| kime
| kinetics
| kiss
| kcostar
| kostica
| kozif
| konad
```

```
| kcostar
| kostice
| kostif
| kond
| konomoka
| laboratorium
| lador
| ledskin
| latolinil
| levissime
| live
| liemall
| likeito
| lish
| lovely
| lowence
| manu
| marathon
| markell
| marutaka-foot
+-----+
| brands
+-----+
| maseur
| matrix
| mavala
| metzger
| mil
| misskin
| missha
| moyou
| nagoraku
| nefertiti
| neader
| nirel
| nitrile
| ophid
| orly
| osmo
| ovalle
| pilar
| polarus
| protafil
| proxene
| protokeratin
| provoc
| rasse
| reflectoil
| ross
| roundoff
| roundl
| s.core
| sancto
| seavirina
| shary
| shis
| skilicity
| skililite
| smart
| solo
| solomeya
| sophin
| stileks
| stting
| supertan
| swarovski
| tectic
| treaclemoon
| trind
| uno
| uskusi
| versaclara
| vilenta
| yokd
| yu-r
| zeitun
```

```
153 rows selected (24.663 seconds)
0: jdbc:hive2://localhost:10000/default> ■
```

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 from (select sum(price) as UserRevenue, user_id from purchases_2019 group by user_id order by UserRevenue desc limit 10) as TopUsers;
```

Screenshot:

```
0: jdbc:hive2://localhost:10000/default> select user_id from (select sum(price) as UserRevenue, user_id from purchases_2019 group by user_id order by UserRevenue desc limit 10) as TopUsers;
INFO : Compiling command(queryId=hive_20201003140654_c3ce6e06-7a01-4851-925e-a925e6842e18): select user_id from (select sum(price) as UserRevenue, user_id from purchases_2019 group by user_id order by UserRevenue desc limit 10) as TopUsers
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema fieldSchemas:[FieldsSchema(name:user_id, type:bigint, comment:null)], properties:null
INFO : Completed compiling command(queryId=hive_20201003140654_c3ce6e06-7a01-4851-925e-a925e6842e18); Time taken: 0.111 seconds
INFO : Concurrency monitor is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20201003140654_c3ce6e06-7a01-4851-925e-a925e6842e18): select user_id from (select sum(price) as UserRevenue, user_id from purchases_2019 group by user_id order by UserRevenue desc limit 10) as TopUsers
INFO : Query ID = hive_20201003140654_c3ce6e06-7a01-4851-925e-a925e6842e18
INFO : Total jobs = 1
INFO : Launching Job 1 out of 1
INFO : Starting task [Stage-1:MAPRED] in serial mode
INFO : Session is already open
INFO : Dag name: select user_id from (select sum(price)...TopUsers(Stage-1)
INFO : Status: Running (Executing on YARN cluster with App id application_1601718673301_0019)

INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+1)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+2)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 0(+4)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 1(-1)/2 Reducer 2: 0/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 0(+1)/1 Reducer 3: 0/1
INFO : Map 1: 2/2 Reducer 2: 1/1 Reducer 3: 0(+1)/1
INFO : Map 1: 2/2 Reducer 2: 1/1 Reducer 3: 1/1
INFO : Completed executing command(queryId=hive_20201003140654_c3ce6e06-7a01-4851-925e-a925e6842e18); Time taken: 15.292 seconds
INFO : OK
+-----+
| user_id |
+-----+
| 557790271 |
| 159318419 |
| 562167663 |
| 531969924 |
| 557859743 |
| 522138011 |
| 561592895 |
| 431980134 |
| 566576888 |
| 521347289 |
+-----+
10 rows selected (15.418 seconds)
0: jdbc:hive2://localhost:10000/default>
```

Above are the top 10 users.

Step 3: Cleaning up

We need to drop the tables first then the database.

Query:

```
drop table ecom_events_2019;
drop table ecom_events_month_2019;
drop table parq_buck_purchase;
drop table purchases_2019;
```

