# MapReduce Assignment

Group Case Study
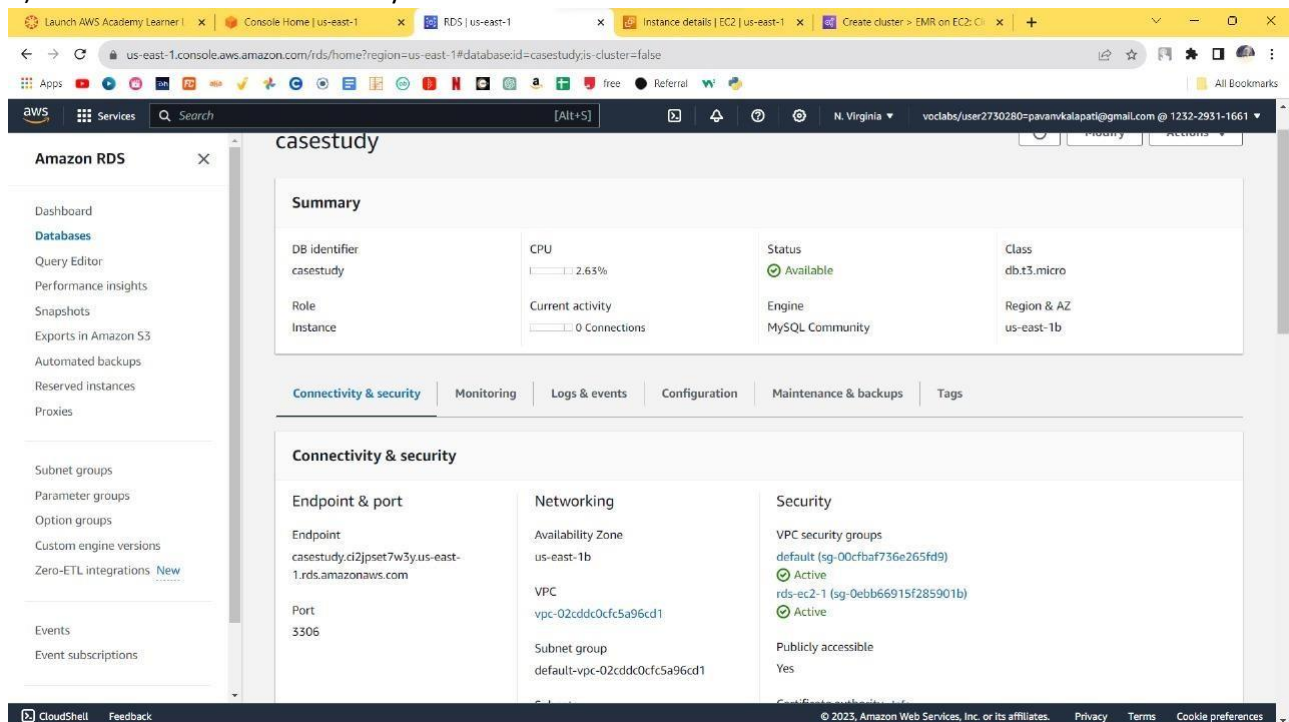
## Prepared By

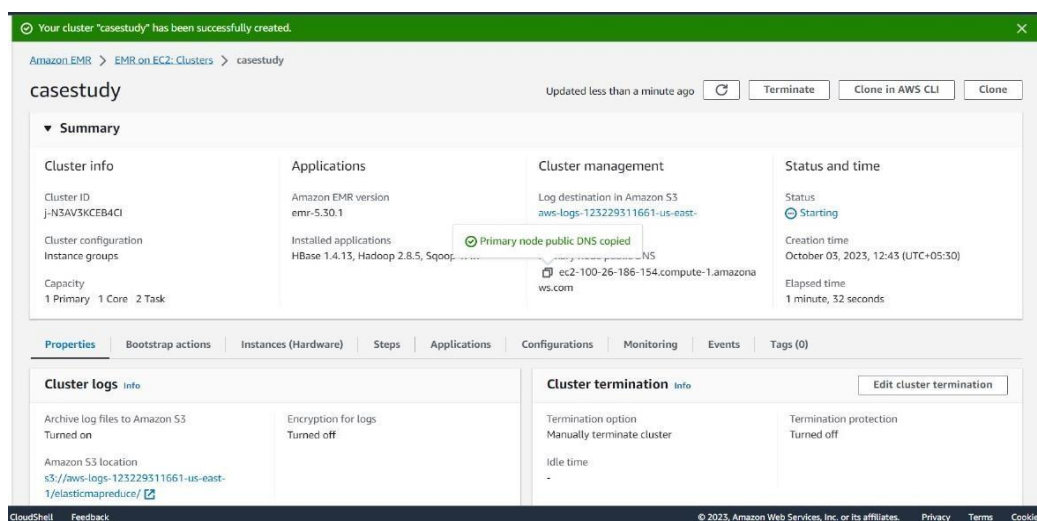ATHENI MADHU SUDHAN
ABDULLAH QANEEH
NEETU SINGH

# Task 1

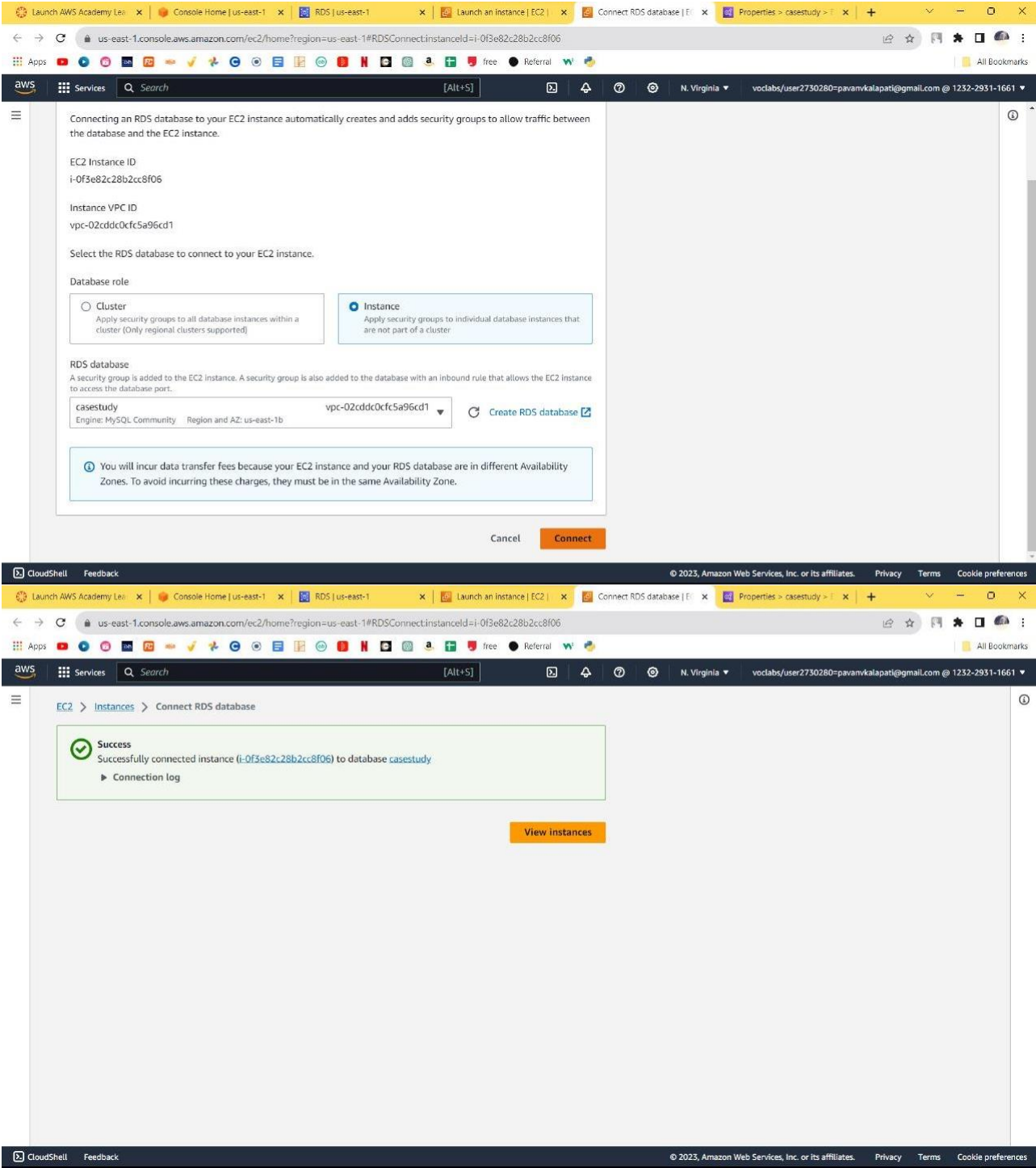## Create an RDS instance in your AWS account and upload the data to the RDS instance

1) We've generated an RDS instance using the Learner Lab and initiated an EMR cluster.

2) We obtained the necessary files for this task, namely yellow_tripdata_2017-01.csv and yellow_tripdata_201702.csv.

3) We've established a connection from EMR to the RDS instance, created a table, and imported records from the CSV files.

4) The data has been successfully loaded into the RDS instance we created.



EMR Cluster:

Connecting RDS with EMR EC2 instance:



Connecting an RDS database to your EC2 instance automatically creates and adds security groups to allow traffic between the database and the EC2 instance.

EC2 Instance ID
i-0f3e82c28b2cc8f06

Instance VPC ID
vpc-02cddc0cfc5a96cd1

Select the RDS database to connect to your EC2 instance.

Database role

○ Cluster
Apply security groups to all database instances within a cluster (Only regional clusters supported)

● Instance
Apply security groups to individual database instances that are not part of a cluster

RDS database
A security group is added to the EC2 instance. A security group is also added to the database with an inbound rule that allows the EC2 instance to access the database port.

casestudy                              vpc-02cddc0cfc5a96cd1
Engine: MySQL Community   Region and AZ: us-east-1b          ↻  Create RDS database

ⓘ You will incur data transfer fees because your EC2 instance and your RDS database are in different Availability Zones. To avoid incurring these charges, they must be in the same Availability Zone.

Cancel    Connect

EC2 > Instances > Connect RDS database

✓ Success
Successfully connected instance (i-0f3e82c28b2cc8f06) to database casestudy
▶ Connection log

View instances

Connecting EMR instance with PuTTy and then downloading **yellow_tripdata_2017-01.csv** & **yellow_tripdata_2017-02.csv** :



**Connecting RDS with EMR Instance:**

**Hostname:** casestudy.ci2jpset7w3y.us-east-1.íds.amazonaws.com Command:

mysql -h case-study-dbb.ck4jzoqb1yn7.us-east-1.rds.amazonaws.com -P 3306 -u root -p

```sql
create database
yellow_taxi;
 use yellow_taxi;
Creating table—
CREATE TABLE taxi
(
vendorID INT,
tpep_pickup_datetime DATETIME,
tpep_dropoff_datetime DATETIME,
passenger_count INT, trip_distance
DOUBLE, puLocationID INT,
doLocationID INT, rateCodeID INT,
store_and_fwd_flag VARCHAR(255),
payment_type INT, fare_amount
DOUBLE, extra DOUBLE, mta_tax
DOUBLE, improvement_surcharge
DOUBLE, tip_amount DOUBLE,

tolls_amount DOUBLE, total_amount
DOUBLE, congestion_Surcharge
DOUBLE, airport_fee DOUBLE
);
```

```
o your MySQL server version for the right syntax to use near 'decribe table taxi' at line 1
MySQL [yellow_taxi]> describe taxi;
+----------------------+-------------+------+-----+---------+-------+
| Field                | Type        | Null | Key | Default | Extra |
+----------------------+-------------+------+-----+---------+-------+
| vendorID             | int         | YES  |     | NULL    |       |
| tpep_pickup_datetime | datetime    | YES  |     | NULL    |       |
| tpep_dropoff_datetime| datetime    | YES  |     | NULL    |       |
| passenger_count      | int         | YES  |     | NULL    |       |
| trip_distance        | double      | YES  |     | NULL    |       |
| puLocationID         | int         | YES  |     | NULL    |       |
| doLocationID         | int         | YES  |     | NULL    |       |
| rateCodeID           | int         | YES  |     | NULL    |       |
| store_and_fwd_flag   | varchar(255)| YES  |     | NULL    |       |
| payment_type         | int         | YES  |     | NULL    |       |
| fare_amount          | double      | YES  |     | NULL    |       |
| extra                | double      | YES  |     | NULL    |       |
| mta_tax              | double      | YES  |     | NULL    |       |
| improvement_surcharge| double      | YES  |     | NULL    |       |
| tip_amount           | double      | YES  |     | NULL    |       |
| tolls_amount         | double      | YES  |     | NULL    |       |
| total_amount         | double      | YES  |     | NULL    |       |
| congestion_Surcharge | double      | YES  |     | NULL    |       |
| airport_fee          | double      | YES  |     | NULL    |       |
+----------------------+-------------+------+-----+---------+-------+
19 rows in set (0.00 sec)
```

**Load Data into Above table**

LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'

INTO TABLE taxi

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;


LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-02.csv'

INTO TABLE taxi

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;

```
Copyright (c) 2000, 2018, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MySQL [(none)]> use yellow_taxi;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
MySQL [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'
    -> INTO TABLE taxi
    -> FIELDS TERMINATED BY ','
    -> LINES TERMINATED BY '\n'
    -> IGNORE 1 LINES;
ERROR 2 (HY000): File '/home/hadoop/yellow_tripdata_2017-01.csv' not found (Errcode: 2)
MySQL [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/ec2-user/yellow_tripdata_2017-01.csv' INTO
 TABLE taxi FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 LINES;
Query OK, 9710820 row
MySQL [yellow_taxi]>
MySQL [yellow_taxi]>
MySQL [yellow_taxi]>
MySQL [yellow_taxi]>                                      hadoop/yellow_tripdata_2017-01.csv' INTO TA
MySQL [yellow_taxi]> RMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 LINES;017-01.csv' INTO
MySQL [yellow_taxi]> Ctrl-C -- exit!
Aborted
[ec2-user@ip-172-31-84-154 ~]$ mysql -h casestudy.ci2jpset7w3y.us-east-1.rds.amazonaws.com -P 3306 -u root -p
Enter password:
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 23
Server version: 8.0.33 Source distribution

Copyright (c) 2000, 2018, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MySQL [(none)]> use yellow_taxi;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
MySQL [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/ec2-user/yellow_tripdata_2017-02.csv'
    -> into table taxi
    -> fields terminated by ','
    -> lines terminated by '\n'
    -> ignore 1 lines;
Query OK, 9169775 rows affected, 65535 warnings (2 min 56.44 sec)
Records: 9169775  Deleted: 0  Skipped: 0  Warnings: 27509325

MySQL [yellow_taxi]>
```

# TASK 2

**Task 2. Use Sqoop command to ingest the data from RDS into the HBase Table.**

**The following steps were followed for data ingestion from RDS into HBase table**

1) Load the data to RDS instance
   Following screen short for reference from previous step (i.e. Task 1)





2) Exit form RDS and load the table data in to the hbase.
   Created the hbase table : hbasetaxi, set the column-family as cf. Copied the data from the RDS taxi table

   **Sqoop import command:**
   sqoop import \
     --connect "jdbc:mysql://casestudy.ci2jpset7w3y.us-east-1.rds.amazonaws.com:3306/yellow_taxi" \
     --username root \
     --password 123456789 \
     --table taxi \
     --columns
   "vendorID,tpep_pickup_datetime,tpep_dropoff_datetime,passenger_count,trip_distance,puLocationID,doLoc

ationID,rateCodeID,store_and_fwd_flag,payment_type,fare_amount,extra,mta_tax,improvement_surcharge,t
ip_amount,tolls_amount,total_amount,congestion_Surcharge,airport_fee" \
  --hbase-create-table \
  --hbase-table hbasetaxi \
  --column-family trip_details \
  --hbase-row-key "vendorID,tpep_pickup_datetime,tpep_dropoff_datetime" \
  --split-by tpep_dropoff_datetime \
  -m 8

Screenshots for reference:

# TASK 2

After the mapreduce

Lets check for the hbasetaxi by running

Scan 'hbasetaxi'

In hbase shell

Bulk import data from two files in the dataset on your EMR cluster to your HBase Table using the relevant codes.

1) In this task, our objective is to import data from two CSV files, namely "yellow_tripdata_2017-03.csv" and "yellow_tripdata_2017-04.csv," into an HBase table.

2) We will initiate a new HBase table specifically for this task.

Note: Since the datasets do not contain a primary key, we will need to modify these datasets by adding a primary column ID. Given the large size of the dataset, we will split the files and work with them in smaller portions for further processing.

**Create new HBase table –**

```
hbase(main):001:0> create 'trip_data_batch', 'cf'
0 row(s) in 1.6310 seconds

=> Hbase::Table - trip_data_batch
hbase(main):002:0>
```

**Batch Insert command and Execution –**

```
[root@ip-172-31-44-105 hadoop]# python file/batch_insert.py
```

**HBase Table Records After Import –**

```
hbase(main):004:0> scan 'trip_data_batch'
ROW                          COLUMN+CELL
 1                           column=cf:DOLocationID, timestamp=1683444168419, value=42
 1                           column=cf:PULocationID, timestamp=1683444168419, value=231
 1                           column=cf:RatecodeID, timestamp=1683444168419, value=1
 1                           column=cf:VendorID, timestamp=1683444168419, value=1
 1                           column=cf:airport_fee, timestamp=1683444168419, value=
 1                           column=cf:congestion_surcharge, timestamp=1683444168419, value=
 1                           column=cf:extra, timestamp=1683444168419, value=0.5
 1                           column=cf:fare_amount, timestamp=1683444168419, value=30.5
 1                           column=cf:improvement_surcharge, timestamp=1683444168419, value=0.3
 1                           column=cf:mta_tax, timestamp=1683444168419, value=0.5
 1                           column=cf:passenger_count, timestamp=1683444168419, value=1
 1                           column=cf:payment_type, timestamp=1683444168419, value=1
 1                           column=cf:store_and_fwd_flag, timestamp=1683444168419, value=N
 1                           column=cf:tip_amount, timestamp=1683444168419, value=6
 1                           column=cf:tolls_amount, timestamp=1683444168419, value=0
 1                           column=cf:total_amount, timestamp=1683444168419, value=37.8
 1                           column=cf:tpep_dropoff_datetime, timestamp=1683444168419, value=01-03-2017 00:59
 1                           column=cf:tpep_pickup_datetime, timestamp=1683444168419, value=01-03-2017 00:38
 1                           column=cf:trip_distance, timestamp=1683444168419, value=10.5
 10                          column=cf:DOLocationID, timestamp=1683444168469, value=82
 10                          column=cf:PULocationID, timestamp=1683444168469, value=82
 10                          column=cf:RatecodeID, timestamp=1683444168469, value=1
 10                          column=cf:VendorID, timestamp=1683444168469, value=1
```

```
 98                          column=cf:tpep_pickup_datetime, timestamp=1683444168752, value=01-03-2017 00:08
 98                          column=cf:trip_distance, timestamp=1683444168752, value=7.2
 99                          column=cf:DOLocationID, timestamp=1683444168757, value=170
 99                          column=cf:PULocationID, timestamp=1683444168757, value=186
 99                          column=cf:RatecodeID, timestamp=1683444168757, value=1
 99                          column=cf:VendorID, timestamp=1683444168757, value=2
 99                          column=cf:airport_fee, timestamp=1683444168757, value=
 99                          column=cf:congestion_surcharge, timestamp=1683444168757, value=
 99                          column=cf:extra, timestamp=1683444168757, value=0.5
 99                          column=cf:fare_amount, timestamp=1683444168757, value=5.5
 99                          column=cf:improvement_surcharge, timestamp=1683444168757, value=0.3
 99                          column=cf:mta_tax, timestamp=1683444168757, value=0.5
 99                          column=cf:passenger_count, timestamp=1683444168757, value=1
 99                          column=cf:payment_type, timestamp=1683444168757, value=2
 99                          column=cf:store_and_fwd_flag, timestamp=1683444168757, value=N
 99                          column=cf:tip_amount, timestamp=1683444168757, value=0
 99                          column=cf:tolls_amount, timestamp=1683444168757, value=0
 99                          column=cf:total_amount, timestamp=1683444168757, value=6.8
 99                          column=cf:tpep_dropoff_datetime, timestamp=1683444168757, value=01-03-2017 00:31
 99                          column=cf:tpep_pickup_datetime, timestamp=1683444168757, value=01-03-2017 00:25
 99                          column=cf:trip_distance, timestamp=1683444168757, value=0.86
198 row(s) in 1.3820 seconds
```

# connecting to HBase server and opening the table
**HBase Table Records After Import –**
**Code –**

**1) Create HBase Table**

```
create 'trip_data_batch', 'cf'
```

**2) Ingest Batch data code –**

**batch_ingest.py**

```
import csv
import happybase
import glob
# Define the HBase table and column family
table_name = 'trip_data_batch'
column_family = 'cf'

# connecting to HBase server and opening the table
connection = happybase.Connection(host='ec2-3-80-189-243.compute-1.amazonaws.com')
table = connection.table(table_name)
csv_dir_path = '/home/hadoop/files/file'
csv_files = glob.glob(csv_dir_path + '/*.csv')

# open CSV file and read the data
for file in csv_files:
        with open(file, 'r') as csvfile:
                reader = csv.DictReader(csvfile)
                for row in reader:
# defining row key and column values for each row
                        row_key = row['ID']
                        column_values = {
                                f'{column_family}:VendorID': row['VendorID'],
                                f'{column_family}:tpep_pickup_datetime':
row['tpep_pickup_datetime'],
                                f'{column_family}:tpep_dropoff_datetime':
row['tpep_dropoff_datetime'],
                                f'{column_family}:passenger_count': row['passenger_count'],
                                f'{column_family}:trip_distance': row['trip_distance'],
                                f'{column_family}:RatecodeID': row['RatecodeID'],
                                f'{column_family}:store_and_fwd_flag': row['store_and_fwd_flag'],
                                f'{column_family}:PULocationID': row['PULocationID'],
                                f'{column_family}:DOLocationID': row['DOLocationID'],
                                f'{column_family}:payment_type': row['payment_type'],
                                f'{column_family}:fare_amount': row['fare_amount'],
                                f'{column_family}:extra': row['extra'],
                                f'{column_family}:mta_tax': row['mta_tax'],
```

```
                    f'{column_family}:tip_amount': row['tip_amount'],
                    f'{column_family}:tolls_amount': row['tolls_amount'],
                    f'{column_family}:improvement_surcharge':
row['improvement_surcharge'],
                    f'{column_family}:total_amount': row['total_amount'],
                    f'{column_family}:congestion_surcharge':
row['congestion_surcharge'],
                    f'{column_family}:airport_fee': row['airport_fee']
                    }
# inserting the row into the HBase table
                    table.put(row_key.encode('utf-8'), column_values)
# Close the connection
connection.close()
```

**Task 4.** Write MapReduce codes to perform the tasks using the files you've downloaded on your EMR Instance:

a) Which vendors have the most trips, and what is the total revenue generated by that vendor?

Code File  reference: mrtask_a.py

Code execution screenshot:

```
[hadoop@ip-172-31-77-47 ~]$ python mrtask_a.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_a.hadoop.20230312.203403.902112
Running step 1 of 2...
reading from STDIN
Running step 2 of 2...
job output is in /tmp/mrtask_a.hadoop.20230312.203403.902112/output
Streaming final output from /tmp/mrtask_a.hadoop.20230312.203403.902112/output...
5583181 "VeriFone Inc."
Removing temp directory /tmp/mrtask_a.hadoop.20230312.203403.902112...
[hadoop@ip-172-31-77-47 ~]$ 
```

**Comment**: VeriFone Inc. holds the highest number of trips in the yellow_tripdata_2017-03.csv dataset, totaling **5,583,181 trips.**

Code File reference:

mrtask_a_TC.py    Code    execution

screenshot:

```
/usr/bin/python3: can't open file 'mrtask_a_TC': [Errno 2] No such file or directory
[hadoop@ip-172-31-77-47 ~]$ python mrtask_a_TC.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_a_TC.hadoop.20230312.211218.095998
Running step 1 of 1...
reading from STDIN
job output is in /tmp/mrtask_a_TC.hadoop.20230312.211218.095998/output
Streaming final output from /tmp/mrtask_a_TC.hadoop.20230312.211218.095998/output...
"Creative Mobile Technologies"  75347398.64700934
"VeriFone Inc." 91682368.32536966
Removing temp directory /tmp/mrtask_a_TC.hadoop.20230312.211218.095998...
[hadoop@ip-172-31-77-47 ~]$ 
```

Comment: VeriFone Inc. has generated a total revenue of **91,682,368.32.**

Code File reference: mrtask_b.py

Code execution screenshot:

```
MR.csv  mrtask_a.py  mrtask_a_10.py  mrtask_a - TotalCharges.py  mrtask_b.py  mrtas
[hadoop@ip-172-31-77-47 ~]$ python mrtask_b.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_b.hadoop.20230312.215459.904500
Running step 1 of 2...
reading from STDIN
Running step 2 of 2...
job output is in /tmp/mrtask_b.hadoop.20230312.215459.904500/output
Streaming final output from /tmp/mrtask_b.hadoop.20230312.215459.904500/output...
13307409.48001407       "132"
Removing temp directory /tmp/mrtask_b.hadoop.20230312.215459.904500...
[hadoop@ip-172-31-77-47 ~]$ []
```

Remarks: Location 132 generated most revenue in the file tested.

What are the different payment types used by customers and their count? The final results should be in a sorted format.

Code File reference: mrtask_c.py

Code execution screenshot:

```
Removing temp directory /tmp/mrtask_c.hadoop.20230312.222000.37101...
[hadoop@ip-172-31-77-47 ~]$ python mrtask_c.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_c.hadoop.20230312.222032.987646
Running step 1 of 2...
reading from STDIN
Running step 2 of 2...
job output is in /tmp/mrtask_c.hadoop.20230312.222032.987646/output
Streaming final output from /tmp/mrtask_c.hadoop.20230312.222032.987646/output...
14999    "4"
53815    "3"
3231928  "2"
6994699  "1"
Removing temp directory /tmp/mrtask_c.hadoop.20230312.222032.987646...
[hadoop@ip-172-31-77-47 ~]$ []
```

Remark: Credit card is mostly used followed by cash, then no charge and lowest one observed is dispute

1= Credit card, 2= Cash, 3= No charge,4= Dispute, 5= Unknown, 6= Voided trip

What is the average trip time for different pickup locations?

Code File reference: mrtask_d.py

Code execution screenshot:

```
[hadoop@ip-172-31-77-47 ~]$ python mrtask_d.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_d.hadoop.20230312.183025.917044
Running step 1 of 1...
reading from STDIN
job output is in /tmp/mrtask_d.hadoop.20230312.183025.917044/output
Streaming final output from /tmp/mrtask_d.hadoop.20230312.183025.917044/output...
"1"     255689.0
"10"    3250964.5
"100"   64877691.0
"101"   7566.333333333333
"102"   77081.66666666667
"105"   4461.666666666667
"106"   637954.0
"107"   71554280.33333333
"108"   13074.333333333334
"109"   28439.333333333332
"11"    18669.5
"111"   4654.0
"112"   1945240.6666666667
"113"   52092151.666666664
"114"   47185916.333333336
"115"   3313.0
"116"   4167355.6666666665
"117"   2352.0
"118"   358.3333333333333
"119"   96017.0
"12"    2766627.5
"120"   14736.333333333334
"121"   20760.333333333332
"122"   21882.666666666668
"123"   50655.333333333336
"124"   46644.333333333336
"125"   20483051.666666668
"126"   40420.333333333336
"127"   349983.3333333333
"128"   19702.0
"129"   1769505.6666666667
"13"    54531881.0
"130"   611231.6666666666
"131"   16631.666666666668
"132"   210844927.66666666
"133"   200180.66666666666
"134"   215742.33333333334
```

Remarks: Running map reduce job on yellow_tripdata_2017-03.csv. Above screenshot has average trip time in seconds for each location.

Calculate the average tips to revenue ratio of the drivers for different locations in sorted format.

Code File reference: mrtask_e.py

```
[hadoop@ip-172-31-77-47 ~]$ python mrtask_e.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_e.hadoop.20230312.214218.070937
Running step 1 of 2...
reading from STDIN
Running step 2 of 2...
job output is in /tmp/mrtask_e.hadoop.20230312.214218.070937/output
Streaming final output from /tmp/mrtask_e.hadoop.20230312.214218.070937/output...
0.0     "27"
0.0     "59"
0.02618897967735177     "245"
0.05555555555555555     "118"
0.08324022346368715     "30"
0.08330668372241766     "44"
0.08333333333333333     "46"
0.08413395744299974     "117"
0.11013215859030838     "5"
0.13513927150290786     "214"
0.13664929969115955     "156"
0.15384615384615385     "187"
0.15791949398506777     "183"
0.16539883643787795     "221"
0.1664588528678304      "84"
0.18001762032128074     "206"
0.1966803775847734      "176"
0.199110563232387       "184"
0.2115120004910365      "139"
0.21793842034805888     "115"
0.23677482792527046     "204"
0.3277439150190755      "253"
0.3423711643513236      "111"
0.3727399722441718      "172"
0.3993444695654762      "150"
0.43193567854033804     "109"
0.45540935672514626     "105"
0.4941239450709766      "222"
0.4965950145681869      "58"
0.5082618573388116      "15"
0.552106455325866       "240"
0.6292976213757296      "122"
0.671368303560592       "248"
0.6782093351324366      "154"
0.6834950879863939      "32"
0.688921042420514       "3"
0.7185910067954439      "254"
0.7281153939519145      "2"
0.7288072293656345      "64"
0.7310509617742412      "205"
```

Comment: The screenshot above was taken during the testing of the code designed to calculate the average tip-to-revenue ratio for drivers in various locations.

Furthermore, the analysis aims to understand how revenue changes over time by calculating the average trip revenue per month. This analysis is carried out by considering the time of day (day versus night) and the day of the week (weekday versus weekend) as factors for comparison.

Code File reference: mrtask_f.py

Code execution screenshot:

```
[hadoop@ip-172-31-77-47 ~]$ python mrtask_f.py < yellow_tripdata_2017-03.csv
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/mrtask_f.hadoop.20230313.112115.944790
Running step 1 of 1...
reading from STDIN
job output is in /tmp/mrtask_f.hadoop.20230313.112115.944790/output
Streaming final output from /tmp/mrtask_f.hadoop.20230313.112115.944790/output..
"March" 33405953.39662997
"day"   46632269.931141086
"night" 5426591.436077364
"weekday"       17745196.803211816
"weekend"       6116198.478531879
Removing temp directory /tmp/mrtask_f.hadoop.20230313.112115.944790...
[hadoop@ip-172-31-77-47 ~]$
```

Observations/Findings:

1. The dataset predominantly consists of data for the month of March.

2. The average revenue for trips in March can be calculated.

3. On average, trip revenue during daytime hours exceeds that of nighttime hours. We assume nighttime to be between 11 PM and 5 AM.

4. The dataset suggests that the average trip revenue during weekdays surpasses that of weekends, with Saturday and Sunday considered as weekends.