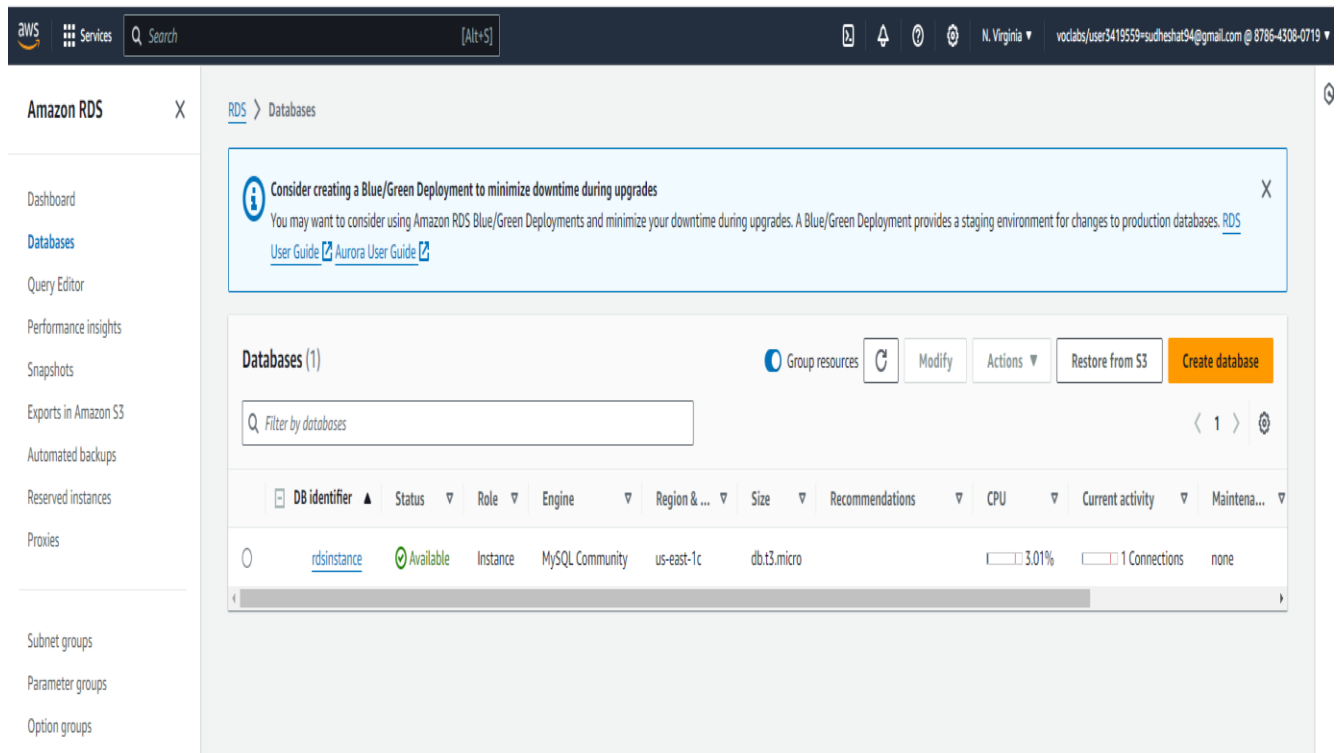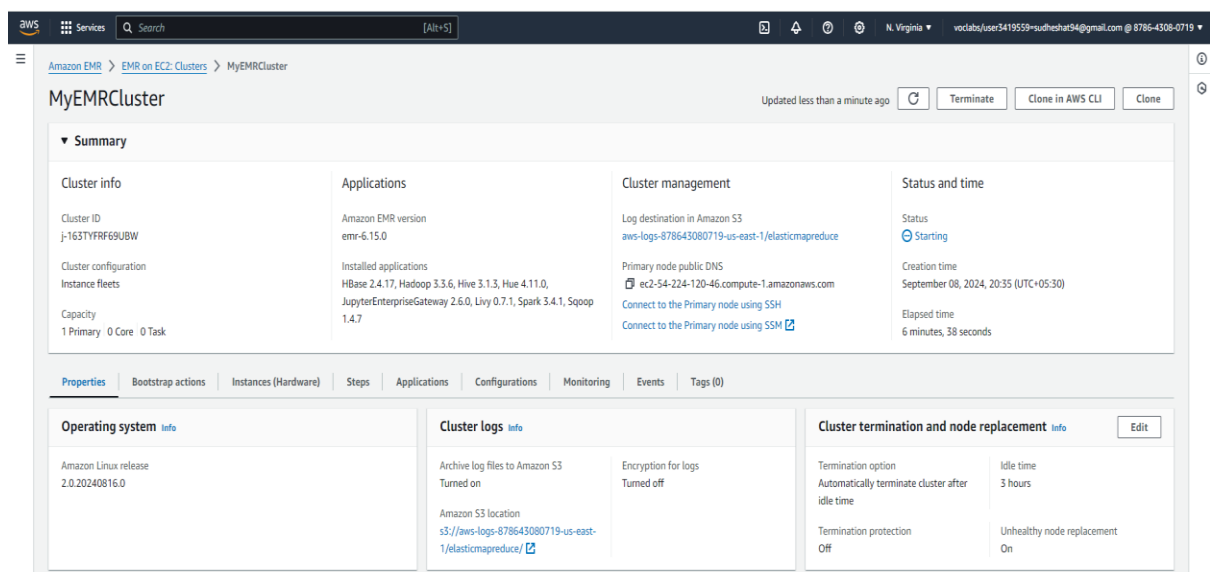# MapReduce Programing Assignment

## Taks 1:

Go to AWS console and search for RDS and the click on create cluster. Select the options as per our requirements and click on create. The below RDS cluster with Database DemoDB is created.



Then went back to Console and searched for EMR. Created EMR cluster as per requirements. Below EMR cluster was created.

Open the EMR cluster using putty from local machine. And download below 2 files using wget commands.

Wget https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-01.csv

Wget https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-02.csv

Connect to RDS instance using below command.

**mysql -h rdsinstance.ckadi12cwuvi.us-east-1.rds.amazonaws.com -P 3306 -u admin -p**

**And give the password.**

Use demoDB database and run below create table query to create the table.

**create table taxi_data (**

**VendorID INT**

**,tpep_pickup_datetime DATETIME DEFAULT CURRENT_TIMESTAMP**

**,tpep_dropoff_datetime DATETIME DEFAULT CURRENT_TIMESTAMP**

**,passenger_count INT**

**,trip_distance FLOAT**

**,RatecodeID INT**

**,store_and_fwd_flag VARCHAR (255)**

**,PULocationID INT**

**,DOLocationID INT**

**,payment_type INT**

**,fare_amount FLOAT**

**,extra FLOAT**

**,mta_tax FLOAT**

**,tip_amount FLOAT**

**,tolls_amount FLOAT**

**,improvement_surcharge FLOAT**

**,total_amount FLOAT**

**,congestion_surcharge FLOAT**

**,airport_fee FLOAT);**

Taxi_data table is created.

```
hadoop@ip-172-31-47-49:~

[hadoop@ip-172-31-47-49 ~]$ mysql -h rdsinstance.ckadil2cwuvi.us-east-1.rds.amazonaws.com -P 3306 -u admin -p
Enter password:
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 100
Server version: 8.0.35 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)]> show databases;
+--------------------+
| Database           |
+--------------------+
| demoDB             |
| information_schema |
| mysql              |
| performance_schema |
| sys                |
+--------------------+
5 rows in set (0.00 sec)

MySQL [(none)]> use demoDB;
Database changed
MySQL [demoDB]> show tables;
Empty set (0.00 sec)

MySQL [demoDB]> create table taxi_data (
    -> VendorID INT
    -> ,tpep_pickup_datetime DATETIME DEFAULT CURRENT_TIMESTAMP
    -> ,tpep_dropoff_datetime DATETIME DEFAULT CURRENT_TIMESTAMP
    -> ,passenger_count INT
    -> ,trip_distance FLOAT
    -> ,RatecodeID INT
    -> ,store_and_fwd_flag VARCHAR (255)
    -> ,PULocationID INT
    -> ,DOLocationID INT
    -> ,payment_type INT
    -> ,fare_amount FLOAT
    -> ,extra FLOAT
    -> ,mta_tax FLOAT
    -> ,tip_amount FLOAT
    -> ,tolls_amount FLOAT
    -> ,improvement_surcharge FLOAT
    -> ,total_amount FLOAT
    -> ,congestion_surcharge FLOAT
    -> ,airport_fee FLOAT);
Query OK, 0 rows affected (0.03 sec)

MySQL [demoDB]> show tables;
+------------------+
| Tables_in_demoDB |
+------------------+
| taxi_data        |
+------------------+
1 row in set (0.00 sec)

MySQL [demoDB]> 
```

Loading the data from 2 files: Run below 2 commands to load data to RDS table taxi_data.

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

INTO TABLE taxi_data

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;

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

INTO TABLE taxi_data

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;