**TASK-2:**

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

1) First login into the EMR instance using Hadoop and switch to root user using the following command:

```
sudo -i
```



2) Complete the initial steps of setup by running the following commands for setting up MySQL connector

```
wget https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz
tar -xvf mysql-connector-java-8.0.25.tar.gz
cd mysql-connector-java-8.0.25/
sudo cp mysql-connector-java-8.0.25.jar /usr/lib/sqoop/lib/
```

3) Start HBase shell to create a table named 'trip_data_hbase' with a column family 'cf1'

```
hbase shell
create 'trip_log_hbase', 'cf1'
exit
```



4) To ingest data from MySQL RDS to HBase table, we run the following commands:

(a) Create a 'sqoop_command.sh' file and insert the following code in it:

```
vim sqoop_command.sh
```



```
sqoop import --connect jdbc:mysql://mydbinstance.c9zffzdmradr.us-east-
1.rds.amazonaws.com:3306/taxi_records --username admin --password 20021960 --table
trip_log --hbase-table trip_log_hbase --column-family cf1 --hbase-create-table --
hbase-row-key tpep_pickup_datetime,tpep_dropoff_datetime --hbase-bulkload --split-
by payment_type
```

This command facilitates efficient data migration from MySQL to HBase by:

Connecting to the specified MySQL database and table:

```
sqoop import --connect jdbc:mysql://mydbinstance.c9zffzdmradr.us-east-
1.rds.amazonaws.com:3306/taxi_records
```

Creating the necessary HBase table and column family:

```
--username admin --password 20021960 --table trip_log --hbase-table
trip_log_hbase --column-family cf1 --hbase-create-table
```

Using composite row keys for unique row identification:

```
--hbase-row-key tpep_pickup_datetime,tpep_dropoff_datetime
```

Employing bulk load for optimal performance:

```
--hbase-bulkload
```

Utilizing parallel processing to speed up the import:

```
--split-by payment_type
```

(b) After saving the 'sqoop_command.sh' run the following command:

```
chmod +x sqoop_command.sh
./sqoop_command.sh
```

5) Check the count of the Hbase table

```
hbase shell
count 'trip_log_hbase'
```

```
hbase(main):001:0> list
TABLE
trip_log_hbase
1 row(s)
Took 0.9007 seconds
=> ["trip_log_hbase"]
hbase(main):002:0> count 'trip_log_hbase'
Current count: 1000, row: 2017-01-01 00:17:38
Current count: 2000, row: 2017-01-01 00:34:19
Current count: 3000, row: 2017-01-01 00:51:00
Current count: 4000, row: 2017-01-01 01:07:40
Current count: 5000, row: 2017-01-01 01:24:20
Current count: 6000, row: 2017-01-01 01:41:04

Current count: 3953000, row: 2017-02-28 14:53:41
Current count: 3954000, row: 2017-02-28 15:13:29
Current count: 3955000, row: 2017-02-28 15:33:52
Current count: 3956000, row: 2017-02-28 15:54:19
Current count: 3957000, row: 2017-02-28 16:14:52
Current count: 3958000, row: 2017-02-28 16:36:06
Current count: 3959000, row: 2017-02-28 16:57:31
Current count: 3960000, row: 2017-02-28 17:18:26
Current count: 3961000, row: 2017-02-28 17:39:13
Current count: 3962000, row: 2017-02-28 18:00:06
Current count: 3963000, row: 2017-02-28 18:20:13
Current count: 3964000, row: 2017-02-28 18:40:29
Current count: 3965000, row: 2017-02-28 19:00:51
Current count: 3966000, row: 2017-02-28 19:21:24
Current count: 3967000, row: 2017-02-28 19:42:43
Current count: 3968000, row: 2017-02-28 20:03:47
Current count: 3969000, row: 2017-02-28 20:25:52
Current count: 3970000, row: 2017-02-28 20:48:00
Current count: 3971000, row: 2017-02-28 21:09:48
Current count: 3972000, row: 2017-02-28 21:32:14
Current count: 3973000, row: 2017-02-28 21:53:47
Current count: 3974000, row: 2017-02-28 22:15:51
Current count: 3975000, row: 2017-02-28 22:39:27
Current count: 3976000, row: 2017-02-28 23:04:25
Current count: 3977000, row: 2017-02-28 23:29:03
Current count: 3978000, row: 2017-02-28 23:57:31
3978081 row(s)
Took 408.6991 seconds
=> 3978081
hbase(main):003:0> |
```