# MapReduce - Programming Assignment

## Data Ingestion Task

BY
Dharan R
Sanjana Dutta
Punith Vadla

# Ingestion_Task

>>sudo -i

>>mkdir HBase

>>ls

>>cd HBase/

>>hbase shell

>>status 'detailed'

>>version

>>table_help


>>create "trip_data","TLC"


**Sqoop requires MySQL's JDBC driver to be installed in order to talk to the MySQL database engine. Download the driver and install it in /usr/lib/sqoop/lib.**
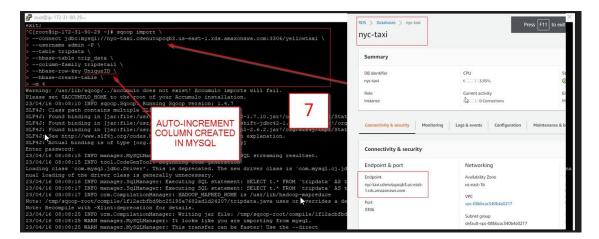
**at root run following commands:**


>>wget http://dev.mysql.com/get/Downloads/Connector-J/mysql-connector-java-5.1.38.tar.gz

>>tar -xvzf mysql-connector-java-5.1.38.tar.gz

>>sudo cp mysql-connector-java-5.1.38/mysql-connector-java-5.1.38-bin.jar /usr/lib/sqoop/lib/


**Get the tables available**

**Now that you have a table created, use Sqoop to get a list of tables that are available in the database.**


**In the following command, use the –P option to be prompted to enter a password. Replace "db-instance-endpoint" with the RDS endpoint retrieved in Step 7, and replace "yourdatabase" with the name of your database in which the "pv_aggregates" table was created.**

>>sqoop list-tables --connect jdbc:mysql://assignmentdb.cqwodvedrqu4.us-east-1.rds.amazonaws.com:3306/assignmentdb --username admin -P

```
 sqoop import \
> --connect  "jdbc:mysql://assignmentdb.cqwodvedrqu4.us-east-1.rds.amazonaws.com:3306/assignmentdb" \
> --table "trip_data" \
> --hbase-table "trip_data" \
> --hbase-row-key "ID" \
> --column-family "TLC" \
> --hbase-create-table \
> --username admin -P
```

in hbase run this command to check count of rows:

count 'trip_data', INTERVAL => 100000