# Data Ingestion from the RDS to HDFS using Sqoop

1. **Sqoop command used for importing table from RDS to HDFS**

- Run Sqoop command to import "**member_score**" table from RDS to HDFS.

  sqoop import --connect jdbc:mysql://upgradawsrds1.cyaielc9bmnf.us-east-
  1.rds.amazonaws.com/cred_financials_data \
  --username upgraduser \
  --password upgraduser \
  --table member_score \
  --null-string 'NA' \
  --null-non-string '\\N' \
  --delete-target-dir \
  --target-dir '/ccfd_capstone_project/member_score' \
  -m 1

- Run Sqoop command to import "**card_member**"  table from RDS to HDFS.

  sqoop import --connect jdbc:mysql://upgradawsrds1.cyaielc9bmnf.us-east-
  1.rds.amazonaws.com/cred_financials_data \
  --username upgraduser \
  --password upgraduser \
  --table card_member \
  --null-string 'NA' \
  --null-non-string '\\N' \
  --delete-target-dir \
  --target-dir '/ccfd_capstone_project/card_member' \
  -m 1

2. **<Command to see the list of imported data in HDFS>**

- Create external table "**card_member_ext**" to hold data from card_member table in RDS.

  CREATE EXTERNAL TABLE IF NOT EXISTS CARD_MEMBER_EXT(`CARD_ID`
  STRING,`MEMBER_ID` STRING,`MEMBER_JOINING_DT` TIMESTAMP,`CARD_PURCHASE_DT`
  STRING,`COUNTRY` STRING,`CITY` STRING)
  ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LOCATION
  '/ccfd_capstone_project/card_member';

- Create external table "**member_score_ext**" to hold data from member_score table in RDS.
  CREATE EXTERNAL TABLE IF NOT EXISTS MEMBER_SCORE_EXT(
  `MEMBER_ID` STRING,
  `SCORE` INT)
  ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
  LOCATION '/ccfd_capstone_project/member_score';

- Create "**card_member_orc**" table. For better performance.
  CREATE TABLE IF NOT EXISTS CARD_MEMBER_ORC(
  `CARD_ID` STRING,
  `MEMBER_ID` STRING,
  `MEMBER_JOINING_DT` TIMESTAMP,
  `CARD_PURCHASE_DT` STRING,
  `COUNTRY` STRING,
  `CITY` STRING)
  STORED AS ORC
  TBLPROPERTIES ("orc.compress"="SNAPPY");

- Create "**member_score_orc**" table. For better performance.
  CREATE TABLE IF NOT EXISTS MEMBER_SCORE_ORC(
  `MEMBER_ID` STRING,
  `SCORE` INT) STORED AS ORC
  TBLPROPERTIES ("orc.compress"="SNAPPY");

- Load data into "**card_member_orc**" table from "**card_member_ext**" table.
  INSERT OVERWRITE TABLE CARD_MEMBER_ORC
  SELECT CARD_ID, MEMBER_ID, MEMBER_JOINING_DT, CARD_PURCHASE_DT, COUNTRY,
  CITY FROM CARD_MEMBER_EXT;

- Load data into "**member_score_orc**" table from "**member_score_ext**" table.
  INSERT OVERWRITE TABLE MEMBER_SCORE_ORC
  SELECT MEMBER_ID, SCORE FROM MEMBER_SCORE_EXT;

- Verify data in "**card_member_orc**" table.
  SELECT * FROM CARD_MEMBER_ORC LIMIT 10;

- Verify data in "**member_score_orc**" table.
  SELECT * FROM MEMBER_SCORE_ORC LIMIT 10;

### 3. Screenshot of the imported data

- Run Sqoop command to import "**member_score**" table from RDS to HDFS.





- Run Sqoop command to import "**card_member**"  table from RDS to HDFS.

- Verify data in "**card_member_orc**" table.

```
hive> SELECT * FROM CARD_MEMBER_ORC LIMIT 10;
OK
340028465709212 009250698176266 2012-02-08 06:04:13      05/13   United States   Barberton
340054675199675 835873341185231 2017-03-10 09:24:44      03/17   United States   Fort Dodge
340082915339645 512969555857346 2014-02-15 06:30:30      07/14   United States   Graham
340134186926007 887711945571282 2012-02-05 01:21:58      02/13   United States   Dix Hills
340265728490548 680324265406190 2014-03-29 07:49:14      11/14   United States   Rancho Cucamonga
340268219434811 929799084911715 2012-07-08 02:46:08      08/12   United States   San Francisco
340379737226464 089615510858348 2010-03-10 00:06:42      09/10   United States   Clinton
340383645652108 181180599313885 2012-02-24 05:32:44      10/16   United States   West New York
340803866934451 417664728506297 2015-05-21 04:30:45      08/17   United States   Beaverton
340889618969736 459292914761635 2013-04-23 08:40:11      11/15   United States   West Palm Beach
Time taken: 0.111 seconds, Fetched: 10 row(s)
hive>
```

- Verify data in "**member_score_orc**" table.

```
hive> SELECT * FROM MEMBER_SCORE_ORC LIMIT 10;
OK
000037495066290 339
000117826301530 289
001147922084344 393
001314074991813 225
001739553947511 642
003761426295463 413
004494068832701 217
006836124210484 504
006991872634058 697
007955566230397 372
Time taken: 0.098 seconds, Fetched: 10 row(s)
hive>
```