# Data Ingestion Using Scoop And Data Analysis Using Hive

## DATA INGESTION FROM RDS TO HDFS USING SQOOP

1. Sqoop import command *(See, sqoop/importFromRDSToHDFS.sqoop)*

   **Note**
   ✓ The table **Key_indicator_districtwise** available in RDS has some **NULL** values in some of the columns.
   ✓ Using the scoop command, the **NULL** values are replaced with **NA** for all String based columns and **\N** for all non-string based columns while importing data into HDFS.
   ✓ This is to make sure the **NULL** value is not written to the HDFS data.

2. Command to see the list of imported data *(See, hdfs/viewData.hdfs)*

# EXTERNAL TABLE CREATION IN HIVE AND LOADING INGESTED DATA

1. Command to create the external table *(See, hive/createHiveTable.hql)*

   **Note**
   ✓   A database named **India_Annual_Health_Survey_2012_13_DB** is created. All the tables pertaining to this project will be created in this database.
   ✓   An external table named **IAHS_2012_13** is created with 645 columns. This table will be used as a master repository of data.

2. Command to load the ingested data into the external table *(See, hive/loadDataInHiveTable.hql)*

3. Queries to verify that the ingestion is correctly accomplished

   3.1 Query to count the total number of rows of data fetched from RDS using MySQL Workbench and from Hive using Hue

   MySQL Workbench *(See, sql/verificationQuery1.sql)*
   Hue *(See, hive/verificationQuery1.hql)*

   3.2 Query to select the top 10 rows and first 8 columns of the data fetched from RDS using MySQL Workbench and from Hive using Hue

   MySQL Workbench *(See, sql/verificationQuery2.sql)*
   Hue *(See, hive/verificationQuery2.hql)*

**Note**

- ✓ The above listed 02 queries and their results across the RDBMS table **Key_indicator_districtwise** and the HIVE table **IAHS_2012_13** should show that the data is correctly imported from RDS to HDFS using sqoop.
- ✓ Later, the same imported data is correctly ingested into the HIVE table **IAHS_2012_13.**

## SUBSET SCHEMA CREATION IN HIVE TO SUPPORT ANALYSIS

1. Columns used in the subset schema
   *ID*
   *State_Name*
   *State_District_Name*
   *AA_Households_Total*
   *AA_Population_Total*
   *CC_Sex_Ratio_All_Ages_Total*
   *LL_Total_Fertility_Rate_Total*
   *YY_Under_Five_Mortality_Rate_U5MR_Total_Person*


2. Storage format used [Benchmark the performance before finalizing the storage format to be used. Create one schema using default format and one in any other format such as ORC for the columns to be used. Insert data into both the tables created. Compare the runtimes of the following queries and decide which format to be used.
   ✓ *select count(*) from <Table Name>;*
   ✓ *select State_Name, count(*) from <Table Name> group by State_Name;*
   ✓ *select * from <Table Name> where State_Name = 'Uttar Pradesh';]*


   **Note**
   In point **03** below,
   ✓ A subset table named **IAHS_2012_13_TEXT** is created with default TEXT format.
   ✓ The subset table contains selected **08** columns.
   ✓ The data is ingested into this table from the master table **IAHS_2012_13**
   In point **04** below,
   ✓ A subset table named **IAHS_2012_13_ORC** is created with ORC format.
   ✓ The subset table contains selected **08** columns.
   ✓ The data is ingested into this table from the master table **IAHS_2012_13**

In point **05** below,

- ✓ **03** sets of queries are executed against both the tables (reference, point 3 and 4) and their execution time is noted.
- ✓ On examining the execution time for all the **03** set of queries, it is observed that the queries executed on the table with ORC format has lower execution time in comparision to the execution time of queries executed on the table with default TEXT format.
- ✓ The difference in execution time of queries is marginal as the data set is small in size.
- ✓ The difference in execution time will increase for a voluminous production size data set.
- ✓ Based on the benchmarking performed for all the **03** queries, I have choosen the ORC format to be used for this project.
- ✓ Additionally, I have also used the compression algorithm SNAPPY with the ORC format as opposed to the non-compressed way of storing data with the default TEXT format.
  The data stored in compressed format saves on disk space which is again helpful when the size of the data set is voluminous.

3. Create and insert command for the default format *(See, hive/createAndInsertDefaultFormat.hql)*

4. Create and insert command for the formats such as ORC *(See, hive/createAndInsertORCFormat.hql)*

5. Screenshot of runtimes against each query given above for the default format as well as for the formats such as ORC

**TEXT FORMAT**

*SELECT COUNT(\*) FROM*
*India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_text;*

Time Taken: **66.117 seconds**

## ORC FORMAT

*SELECT COUNT(\*) FROM*
*India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_orc;*

**Time Taken: *60.657 seconds***

## TEXT FORMAT

SELECT State_Name, COUNT(*) FROM
India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_text GROUP
BY State_Name;

Time Taken: **185.335 seconds**

# ORC FORMAT

*SELECT State_Name, COUNT(*) FROM India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_orc GROUP BY State_Name;*

Time Taken: *144.896 seconds*

**TEXT FORMAT**

*SELECT \* FROM India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_text WHERE State_Name = "Uttar Pradesh";*

Time Taken: **47.038 seconds**

# ORC FORMAT

*SELECT * FROM*
*India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_orc WHERE*
*State_Name = "Uttar Pradesh";*

Time Taken: **46.913 seconds**

6. Create and insert command for the partition table for analyses 1 & 2. The partition table should be created using the table created above. *(See, hive/createAndInsertORCFormatPartitioned.hql)*

   **Note**
   For analyses 1 and 2,
   ✔ A partitioned table named **IAHS_2012_13_PARTITIONED_ORC_FORMAT** is created with ORC format.
   ✔ The data into this table is ingested from the master table **IAHS_2012_13.**
   ✔ This table will be used only for writing queries for analyses 1 and 2.

   For analyses 3, 4 and 5, the non-partitioned ORC format table **IAHS_2012_13_ORC** will be used.

# QUERY ANALYSIS, RESULT AND CHART

1. State wise child mortality rate

   **Query**
   *SELECT State_Name,*
   *ROUND(AVG(YY_Under_Five_Mortality_Rate_U5MR_Total_Person),2) AS*
   *State_Wise_Average_Child_Mortality_Rate*
   *FROM*
   *India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_partitioned_*
   *orc_format GROUP BY State_Name;*

   **Screenshot of the result**

# Chart
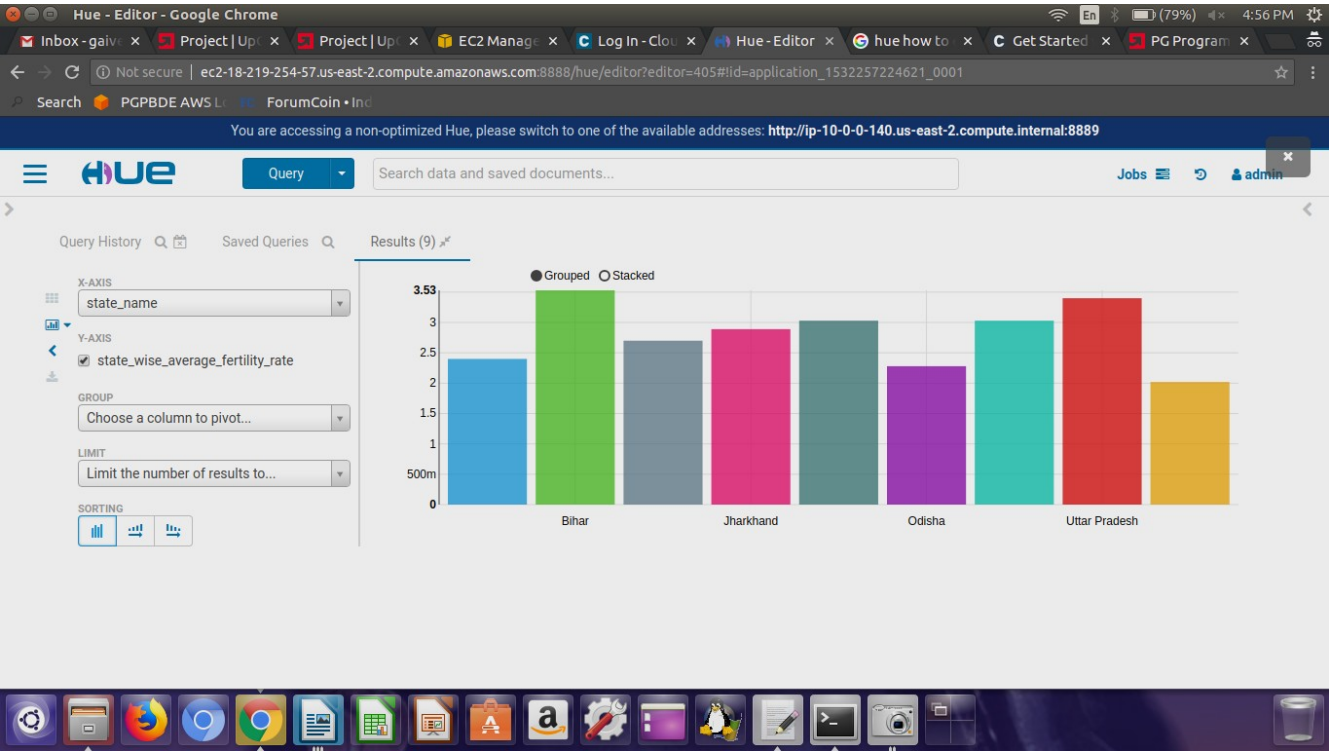
2. State wise fertility rate

**Query**

*SELECT State_Name,*
*ROUND(AVG(LL_Total_Fertility_Rate_Total),2) AS*
*State_Wise_Average_Fertility_Rate*
*FROM*
*India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_partitioned_*
*orc_format GROUP BY State_Name;*

**Screenshot of the result**



| | state_name | state_wise_average_fertility_rate |
|---|---|---|
| 1 | Assam | 2.4 |
| 2 | Bihar | 3.53 |
| 3 | Chhattisgarh | 2.7 |
| 4 | Jharkhand | 2.89 |
| 5 | Madhya Pradesh | 3.03 |
| 6 | Odisha | 2.28 |
| 7 | Rajasthan | 3.03 |
| 8 | Uttar Pradesh | 3.4 |
| 9 | Uttarakhand | 2.02 |

# Chart

3. Does high fertility correlate with high child mortality?

**Query**

*SELECT State_Name,*
*CORR(YY_Under_Five_Mortality_Rate_U5MR_Total_Person,*
*LL_Total_Fertility_Rate_Total)*
*FROM India_Annual_Health_Survey_2012_13_DB.iahs_2012_13_orc*
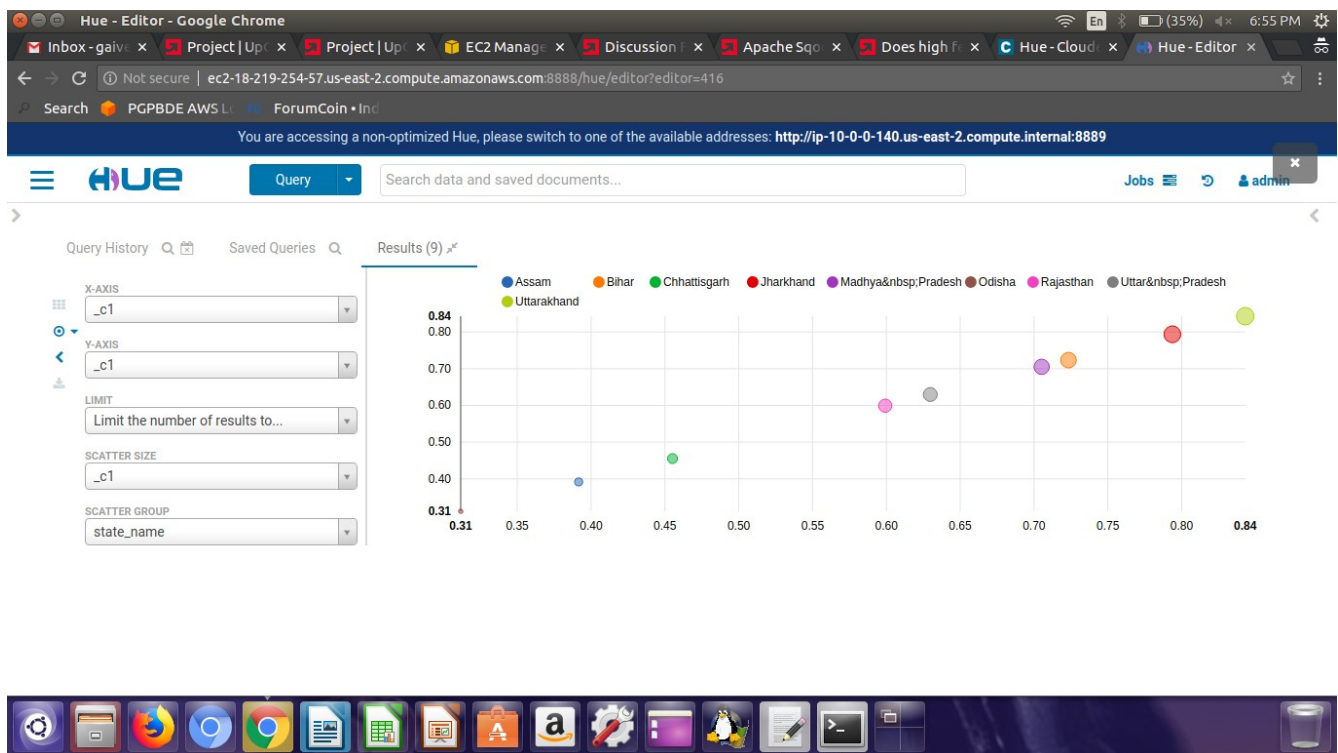*GROUP BY State_Name;*

**Screenshot of the result**

# Chart



## Note

✔ Based on the analysis of the output, we see a
positive slope in the scatter plot above as all the
correlation co-efficient lie in the range of 0.3 to
0.8

4. Find top 2 districts per state with the highest population per household

**Query**
*SELECT*
*tmp_table.State_Name,*
*tmp_table.State_district_name,*
*tmp_table.Population_Per_House_Hold*
*FROM (*
*SELECT*
*State_Name,*
*State_district_name,*
*(AA_Population_Total/AA_Households_Total) AS*
*Population_Per_House_Hold,*
*RANK() OVER (PARTITION BY State_Name ORDER BY*
*(AA_Population_Total/AA_Households_Total) DESC) AS Rank*
*FROM iahs_2012_13_orc*
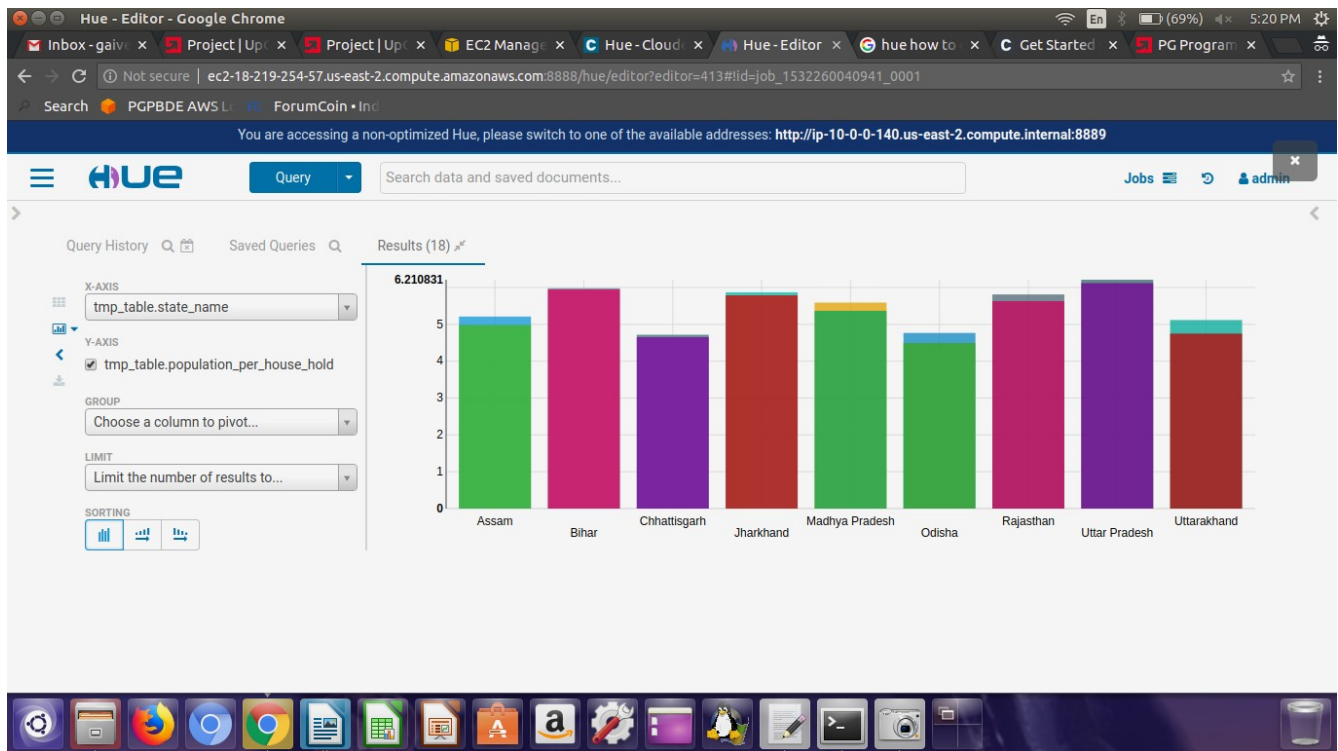*) tmp_table WHERE Rank < 3;*

**Screenshot of the result**



| | tmp_table.state_name | tmp_table.state_district_name | tmp_table.population_per_ho |
|---|---|---|---|
| 1 | Assam | Dhemaji | 5.2103445894620535 |
| 2 | Assam | Marigaon | 4.978445126406547 |
| 3 | Bihar | Gopalganj | 5.979195301761839 |
| 4 | Bihar | Nawada | 5.944978455419291 |
| 5 | Chhattisgarh | Durg | 4.716408016844732 |
| 6 | Chhattisgarh | Rajnandgaon | 4.651162790697675 |
| 7 | Jharkhand | Kodarma | 5.868167462952465 |
| 8 | Jharkhand | Giridih | 5.787106964805766 |
| 9 | Madhya Pradesh | Jhabua | 5.5903925014645575 |
| 10 | Madhya Pradesh | Sehore | 5.366774132372464 |
| 11 | Odisha | Bhadrak | 4.765950743055191 |
| 12 | Odisha | Jajapur | 4.494145867839397 |
| 13 | Rajasthan | Dhaulpur | 5.810972222222222 |
| 14 | Rajasthan | Barmer | 5.629192111322455 |
| 15 | | | |

# Chart

5. Find top 2 districts per state with the lowest sex ratios

**Query**
SELECT
tmp_table.State_Name,
tmp_table.State_district_name,
tmp_table.CC_Sex_Ratio_All_Ages_Total
FROM (
SELECT
State_Name,
State_district_name,
CC_Sex_Ratio_All_Ages_Total,
RANK() OVER (PARTITION BY State_Name ORDER BY
CC_Sex_Ratio_All_Ages_Total ASC) AS Rank
FROM iahs_2012_13_orc
) tmp_table WHERE Rank < 3;

**Screenshot of the result**

# Chart