# The Course Project- Jasvitha Buggana – student59- js3225

The course project includes 3 parts. The first part is to develop a Python application to retrieve Year and Temperature from original NCDC records (i.e., the dataset we are using for this class) and then write the Year and Temperature data into a text file. The second part is to load the text file into Pig and get the highest and lowest temperatures for each year. The third part is to load the text file into Hive and get the average temperature for each year.

**Project Steps and Outputs:**

Created two Python files: Project_Map.py and Project_Reduce.py

Download CourseProjectData.zip and transfer the files using winscp.

- hdfs dfs -copyFromLocal /home/student59/CourseProjectData /home/59student59/CourseProjectData
- hadoop jar /home/student59/hadoop-streaming-2.9.0.jar \input /home/59student59/CourseProjectData \output /home/59student59/Project_js3225 \mapper /home/student59/Project_map.py \reducer /home/student59/Project_Reduce.py \file /home/student59/Project_map.py \file /home/student59/Project_Reduce.py

- hdfs dfs -ls /home/59student59/Project_js3225/

- hdfs dfs -copyToLocal /home/59student59/Project_js3225/part-00000 /home/student59/js3225_output.txt
- js_3225_output file will be generated.



- **In Pig:** pig -x local
- records = LOAD 'js3225_output.txt' AS (year:chararray, temperature:int);
- grouped_records = GROUP records BY year;

- maxtemp = FOREACH grouped_records GENERATE group, MAX(records.temperature);
- DUMP maxtemp;

- mintemp = FOREACH grouped_records GENERATE group, MIN(records.temperature);
- DUMP mintemp;

- **In hive:**
- DROP TABLE IF exists js3225table;
- CREATE TABLE js3225table (year STRING, temperature INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';

- LOAD DATA LOCAL INPATH 'js3225_output.txt' OVERWRITE INTO TABLE js3225table;
- SELECT year, AVG(temperature)

  FROM js3225table

  GROUP BY year;

```
student59@msba-hadoop-name:~
SLF4J: Found binding in [jar:file:/usr/local/hive-2.3.2/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/local/hadoop-2.9.0/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]

Logging initialized using configuration in jar:file:/usr/local/hive-2.3.2/lib/hive-common-2.3.2.jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hi
ve 1.X releases.
hive> DROP TABLE IF exists js3225table;
OK
Time taken: 4.861 seconds
hive> CREATE TABLE js3225table (year STRING, temperature INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
OK
Time taken: 0.869 seconds
hive> LOAD DATA LOCAL INPATH 'js3225_output.txt' OVERWRITE INTO TABLE js3225table;
Loading data to table default.js3225table
OK
Time taken: 0.893 seconds
hive> SELECT year, AVG(temperature)
    > FROM js3225table
    > GROUP BY year;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or
using Hive 1.X releases.
Query ID = student59_20240505130753_0c075c30-c4ca-440c-b5fa-405510405b9d
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1714512003524_1037, Tracking URL = http://msba-hadoop-name:8088/proxy/application_1714512003524_1037/
Kill Command = /usr/local/hadoop/bin/hadoop job  -kill job_1714512003524_1037
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2024-05-05 13:08:20,308 Stage-1 map = 0%,  reduce = 0%
2024-05-05 13:08:26,513 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 3.05 sec
2024-05-05 13:08:32,669 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 5.18 sec
MapReduce Total cumulative CPU time: 5 seconds 180 msec
Ended Job = job_1714512003524_1037
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1   Cumulative CPU: 5.18 sec   HDFS Read: 1192610 HDFS Write: 122 SUCCESS
Total MapReduce CPU Time Spent: 5 seconds 180 msec
OK
1940    74.91863717027938
Time taken: 40.66 seconds, Fetched: 1 row(s)
hive>
```

You need to turn in 1) Python files (mapper and reducer), 2) the commands for executing the Python application in Hadoop, 3) the text file including Year and Temperature data created by you, 4) the screenshot of the text file being created, 5) the screenshot of the final Pig output showing the year and the highest and lowest temperatures, and 6) the screenshot of the final Hive output showing the year and average temperature.

The original dataset for this project is available on Canvas.