

# Map Reduce Project - BAN632

### 1) Python files (mapper and reducer) ← in project submission

**2) The commands for executing the Python application in Hadoop are as follows:**

```
# put data into correct directory
```

```
hdfs dfs -copyFromLocal CourseProjectData/ hdfs://msba-hadoop-name:9000/user/student29/
```

```
# run mapper
```

```
spark-submit NcdcRecordMapper.py CourseProjectData/*.gz
```

```
# run reducer
```

```
spark-submit --master local[*] NcdcRecordReducer.py
```

**3) The text file including Year and Temperature data created by you ← in project submission**

#### 4) The screenshot of the text file being created

```
# The file is created when running NcdcRecordReducer.py.
```

## First Screenshot: mapper

[illegible]

Second Screenshot: reducer (red markup shows reducer works), this is when the yr\_temp\_data.txt file is created.

```
24/05/01 11:22:28 INFO spark.SparkContext: Running Spark version 2.4.3
24/05/01 11:22:28 INFO spark.SparkContext: Submitted application: NDCReducer
24/05/01 11:22:28 INFO spark.SecurityManager: Changing view acls to: student29
24/05/01 11:22:28 INFO spark.SecurityManager: Changing modify acls to: student29
24/05/01 11:22:28 INFO spark.SecurityManager: Changing view acls groups to:
24/05/01 11:22:28 INFO spark.SecurityManager: Changing modify acls groups to:
24/05/01 11:22:28 INFO spark.SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(student29); groups with view permissions: Set(); users with modify permissions: Set(student29); groups with modify permissions: Set()
24/05/01 11:22:28 INFO util.Utils: Successfully started service 'sparkDriver' on port 44026.
24/05/01 11:22:28 INFO spark.SparkEnv: Registering MapOutputTracker
24/05/01 11:22:28 INFO spark.SparkEnv: Registering BlockManagerMaster
24/05/01 11:22:28 INFO storage.BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
24/05/01 11:22:28 INFO storage.BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
24/05/01 11:22:28 INFO storage.DiskBlockManager: Created local directory at /tmp/blockmgr-3d0dfd7a-59db-4155-bf1e-cc6447e6d739
24/05/01 11:22:28 INFO memory.MemoryStore: MemoryStore started with capacity 366.3 MB
24/05/01 11:22:28 INFO spark.SparkEnv: Registering OutputCommitCoordinator
24/05/01 11:22:29 INFO util.log: Logging initialized @2002ms
24/05/01 11:22:29 INFO server.Server: jetty-9.3.z-SNAPSHOT, build timestamp: unknown, git hash: unknown
24/05/01 11:22:29 INFO server.Server: Started @2068ms
24/05/01 11:22:29 INFO server.AbstractConnector: Started ServerConnector@68bbb52e[HTTP/1.1,[http/1.1]][0.0.0.0:4040]
24/05/01 11:22:29 INFO util.Utils: Successfully started service 'SparkUI' on port 4040.
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@4a94b07d[/jobs,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2bddcc24[/jobs/json,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@6d5a2dd7[/jobs/job,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@653cd5e3[/jobs/job/json,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@20c5898c[/stages,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@1c4d2cd6[/stages/json,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@7590c24c[/stages/stage,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@fdb446f[/stages/stage/json,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@7caa657a[/stages/pool,null,AVAILABLE,@Spark]
24/05/01 11:22:29 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@e4e7e2d[/stages/pool/json,null,AVAILABLE,@Spark]
```

5) the screenshot of the final Pig output showing the year and the highest and lowest temperatures

# load data from text file

temperature\_data = LOAD

'hdfs://msba-hadoop-name:9000/user/student29/output/aggregated\_temperature/yr\_temp\_data.txt' USING PigStorage(',') AS (year:int, temperature:int);

# group data by year

grouped\_data = GROUP temperature\_data BY year;

# extract highest and lowest temperatures

temperature\_stats = FOREACH grouped\_data { max\_temp = MAX(temperature\_data.temperature); min\_temp = MIN(temperature\_data.temperature); GENERATE group AS year, max\_temp AS max\_temperature, min\_temp AS min\_temperature; }

# MIN MAX TEMP OUTPUT AT THE BOTTOM OF SCREENSHOT

STORE temperature\_stats INTO 'output/temperature\_stats' USING PigStorage(',');  
DUMP temperature\_stats;

```
Success!
Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature  Outputs
Job_1714512003524_0176  1  1  3  3  3  2  2  2  2  grouped_data,temperature_data,temperature_stats  GROUP_BY,COMBINE  hdfs://msba-h
adpoo-name:9000/tmp/temp-2072864620/tmp-1484360426,

Input(s):
Successfully read 66319 records (561595 bytes) from: "hdfs://msba-hadoop-name:9000/user/student29/output/aggregated_temperature/part-000000"

Output(s):
Successfully stored 1 records (12 bytes) in: "hdfs://msba-hadoop-name:9000/tmp/temp-2072864620/tmp-1484360426"

Counters:
Total records written : 1
Total bytes written : 12
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
Job_1714512003524_0176

2024-05-01 12:00:08,564 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /127.0.0.1:8032
2024-05-01 12:00:08,567 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2024-05-01 12:00:08,592 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /127.0.0.1:8032
2024-05-01 12:00:08,594 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2024-05-01 12:00:08,600 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /127.0.0.1:8032
2024-05-01 12:00:08,603 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2024-05-01 12:00:08,621 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2024-05-01 12:00:08,624 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2024-05-01 12:00:08,624 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2024-05-01 12:00:08,638 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2024-05-01 12:00:08,638 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(1940,128,-61)
```

**6) the screenshot of the final Hive output showing the average year and temperature.**

# typical hive setup

ls -l | grep meta

mv metastore\_db metastore\_db.old

schematool -dbType derby -initSchema

hive

set hive.metastore.warehouse.dir;

# create table in a manner that I can load the text file into it

```
CREATE TABLE IF NOT EXISTS temperature_data (year INT, temperature INT) ROW
FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n'
STORED AS TEXTFILE;
```

# load text file into table

LOAD DATA INPATH

```
'hdfs://msba-hadoop-name:9000/user/student29/output/aggregated_temperature/yr_temp_data.txt
' INTO TABLE temperature_data;
```

# query for average temperature

```
SELECT year, AVG(temperature) AS avg_temperature FROM temperature_data WHERE year
IS NOT NULL AND temperature IS NOT NULL GROUP BY year;
```

```
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2024-05-01 12:14:40,241 Stage-1 map = 0%, reduce = 0%
2024-05-01 12:14:47,467 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 5.03 sec
2024-05-01 12:14:52,595 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 7.07 sec
MapReduce Total cumulative CPU time: 7 seconds 70 msec
Ended Job = job_1714512003524_0177
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 7.07 sec HDFS Read: 4060871 HDFS Write: 123 SUCCESS
Total MapReduce CPU Time Spent: 7 seconds 70 msec
OK
1940      56.413818061189104 |
Time taken: 21.548 seconds, Fetched: 1 row(s)
hive>
```