**Exp. No : 3**

# Map Reduce program to process Weather dataset

1. Download Weather dataset.



2. Create mapper.py program

3. Create reducer.py

```
  GNU nano 7.2                         reducer.py                        Modified
#!/usr/bin/env python
from operator import itemgetter
import sys
current_month = None
current_max = 0
month = None
for line in sys.stdin:
        line = line.strip()
        month, daily_max = line.split('\t', 1)
        try:
                daily_max = float(daily_max)
        except ValueError:
                continue
        if current_month == month:
                if daily_max > current_max:
                        current_max = daily_max
        else:
                if current_month:
                        print ('%s\t%s' % (current_month, current_max))
                current_max = daily_max
                current_month = month
if current_month == month:
        print ('%s\t%s' % (current_month, current_max))

^G Help        ^O Write Out   ^W Where Is    ^K Cut         ^T Execute
^X Exit        ^R Read File   ^\ Replace     ^U Paste       ^J Justify
```

4. Run the Map reduce program using Hadoop Streaming.

```
keerthi@fedora:-$ hadoop jar $HADOOP_STREAMING -input /exp3/dataset.txt -output /exp3/output3 -mapper ~/exp3/mapper.py -reducer ~/exp3/re
ducer.py
packageJobJar: [/tmp/hadoop-unjar7638139713067505238/] [] /tmp/streamjob2412380005828032716.jar tmpDir=null
2024-10-20 12:07:28,334 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2024-10-20 12:07:28,493 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2024-10-20 12:07:28,951 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/keerthi/.staging/
job_1729431042099_0006
2024-10-20 12:07:29,267 INFO mapred.FileInputFormat: Total input files to process : 1
2024-10-20 12:07:29,846 INFO mapreduce.JobSubmitter: number of splits:2
2024-10-20 12:07:30,203 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1729431042099_0006
2024-10-20 12:07:30,203 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-10-20 12:07:30,431 INFO conf.Configuration: resource-types.xml not found
2024-10-20 12:07:30,432 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2024-10-20 12:07:30,874 INFO impl.YarnClientImpl: Submitted application application_1729431042099_0006
2024-10-20 12:07:30,969 INFO mapreduce.Job: The url to track the job: http://fedora:8088/proxy/application_1729431042099_0006/
2024-10-20 12:07:30,972 INFO mapreduce.Job: Running job: job_1729431042099_0006
2024-10-20 12:07:42,453 INFO mapreduce.Job: Job job_1729431042099_0006 running in uber mode : false
2024-10-20 12:07:42,457 INFO mapreduce.Job:  map 0% reduce 0%
2024-10-20 12:07:54,898 INFO mapreduce.Job:  map 100% reduce 0%
2024-10-20 12:08:02,002 INFO mapreduce.Job:  map 100% reduce 100%
2024-10-20 12:08:05,045 INFO mapreduce.Job: Job job_1729431042099_0006 completed successfully
2024-10-20 12:08:05,377 INFO mapreduce.Job: Counters: 54
        File System Counters
                FILE: Number of bytes read=102094
                FILE: Number of bytes written=1039165
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=83480
                HDFS: Number of bytes written=96
                HDFS: Number of read operations=11
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
```

```
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
Job Counters
        Launched map tasks=2
        Launched reduce tasks=1
        Data-local map tasks=2
        Total time spent by all maps in occupied slots (ms)=20527
        Total time spent by all reduces in occupied slots (ms)=4945
        Total time spent by all map tasks (ms)=20527
        Total time spent by all reduce tasks (ms)=4945
        Total vcore-milliseconds taken by all map tasks=20527
        Total vcore-milliseconds taken by all reduce tasks=4945
        Total megabyte-milliseconds taken by all map tasks=21019648
        Total megabyte-milliseconds taken by all reduce tasks=5063680
Map-Reduce Framework
        Map input records=365
        Map output records=10220
        Map output bytes=81648
        Map output materialized bytes=102100
        Input split bytes=180
        Combine input records=0
        Combine output records=0
        Reduce input groups=12
        Reduce shuffle bytes=102100
        Reduce input records=10220
        Reduce output records=12
        Spilled Records=20440
        Shuffled Maps =2
        Failed Shuffles=0
        Merged Map outputs=2
        GC time elapsed (ms)=513
        CPU time spent (ms)=5780
```

```
        Input split bytes=180
        Combine input records=0
        Combine output records=0
        Reduce input groups=12
        Reduce shuffle bytes=102100
        Reduce input records=10220
        Reduce output records=12
        Spilled Records=20440
        Shuffled Maps =2
        Failed Shuffles=0
        Merged Map outputs=2
        GC time elapsed (ms)=513
        CPU time spent (ms)=5780
        Physical memory (bytes) snapshot=873496576
        Virtual memory (bytes) snapshot=7734071296
        Total committed heap usage (bytes)=693108736
        Peak Map Physical memory (bytes)=319991808
        Peak Map Virtual memory (bytes)=2576502784
        Peak Reduce Physical memory (bytes)=234762240
        Peak Reduce Virtual memory (bytes)=2582114304
Shuffle Errors
        BAD_ID=0
        CONNECTION=0
        IO_ERROR=0
        WRONG_LENGTH=0
        WRONG_MAP=0
        WRONG_REDUCE=0
File Input Format Counters
        Bytes Read=83300
File Output Format Counters
        Bytes Written=96
2024-10-20 12:08:05,377 INFO streaming.StreamJob: Output directory: /exp3/output3
keerthi@fedora:~$
```

**Output :**

```
keerthi@fedora:~$ hdfs dfs -cat /exp3/output3/part-00000
01      26.5
02      26.6
03      29.1
04      30.8
05      31.1
06      33.6
07      38.5
08      40.2
09      36.5
10      36.9
11      27.6
12      25.9
```