

Manideep Elasagaram

Net ID: is6496

## The Course Project

**Plagiarism and cheating will not be tolerated for this individual and independent project assignment. All submissions will undergo a thorough plagiarism check. Sharing your work with other students, copying others' work, or using AI tools such as ChatGPT is strictly prohibited. Any instances of plagiarism or cheating will result in a grade of 0 for the entire project submission.**

The course project consists of four parts. Parts 1 and 2 require the use of the project data zip file, while parts 3 and 4 require the use of a text file named “skyceiling.txt”.

Part 1:

The first part is to develop a Mapper and Reducer application to calculate the *range* (the difference between max and min values) of *sky ceiling height* (meters) for *each observation month* from NCDC records (note: 99999 indicates missing value, and [01459] indicate good quality value).

Step:1

chmod +x Mapper45.py

chmod +x Reducer45.py

```
student45@msba-hadoop ~ % 
[student45@msba-hadoop-name ManideepClassProject]$ pwd
/home/student45/ManideepClassProject
[student45@msba-hadoop-name ManideepClassProject]$ ls
011060-99999-1928.gz 014270-99999-1930.gz 028970-99999-1923.gz 029350-99999-1922.gz 030910-99999-1927.gz
011060-99999-1929.gz 023610-99999-1929.gz 028970-99999-1924.gz 029350-99999-1923.gz 030910-99999-1928.gz
011060-99999-1930.gz 023610-99999-1930.gz 028970-99999-1925.gz 029350-99999-1924.gz 032620-99999-1927.gz
012620-99999-1928.gz 028360-99999-1921.gz 028970-99999-1926.gz 029350-99999-1925.gz 033020-99999-1927.gz
012620-99999-1929.gz 028360-99999-1922.gz 029110-99999-1921.gz 029350-99999-1926.gz 034970-99999-1927.gz
012620-99999-1930.gz 028360-99999-1923.gz 029110-99999-1922.gz 029700-99999-1921.gz 038040-99999-1927.gz
014030-99999-1928.gz 028360-99999-1924.gz 029110-99999-1923.gz 029700-99999-1922.gz hadoop-streaming-2.7.3.jar
014030-99999-1929.gz 028360-99999-1925.gz 029110-99999-1924.gz 029700-99999-1923.gz Mapper45.py
014030-99999-1930.gz 028360-99999-1926.gz 029110-99999-1925.gz 029700-99999-1924.gz Reducer45.py
014270-99999-1928.gz 028970-99999-1921.gz 029110-99999-1926.gz 029700-99999-1925.gz
014270-99999-1929.gz 028970-99999-1922.gz 029350-99999-1921.gz 029700-99999-1926.gz
[student45@msba-hadoop-name ManideepClassProject]$ chmod +x Mapper45.py
[student45@msba-hadoop-name ManideepClassProject]$ chmod +x Reducer45.py
[student45@msba-hadoop-name ManideepClassProject]$ ls
011060-99999-1928.gz 014270-99999-1930.gz 028970-99999-1923.gz 029350-99999-1922.gz 030910-99999-1927.gz
011060-99999-1929.gz 023610-99999-1929.gz 028970-99999-1924.gz 029350-99999-1923.gz 030910-99999-1928.gz
011060-99999-1930.gz 023610-99999-1930.gz 028970-99999-1925.gz 029350-99999-1924.gz 032620-99999-1927.gz
012620-99999-1928.gz 028360-99999-1921.gz 028970-99999-1926.gz 029350-99999-1925.gz 033020-99999-1927.gz
012620-99999-1929.gz 028360-99999-1922.gz 029110-99999-1921.gz 029350-99999-1926.gz 034970-99999-1927.gz
012620-99999-1930.gz 028360-99999-1923.gz 029110-99999-1922.gz 029700-99999-1921.gz 038040-99999-1927.gz
014030-99999-1928.gz 028360-99999-1924.gz 029110-99999-1923.gz 029700-99999-1922.gz hadoop-streaming-2.7.3.jar
014030-99999-1929.gz 028360-99999-1925.gz 029110-99999-1924.gz 029700-99999-1923.gz Mapper45.py
014030-99999-1930.gz 028360-99999-1926.gz 029110-99999-1925.gz 029700-99999-1924.gz Reducer45.py
014270-99999-1928.gz 028970-99999-1921.gz 029110-99999-1926.gz 029700-99999-1925.gz
014270-99999-1929.gz 028970-99999-1922.gz 029350-99999-1921.gz 029700-99999-1926.gz
[student45@msba-hadoop-name ManideepClassProject]$
```

## Step:2

Unzipping the files:

```
gunzip *-99999-*.gz
```

```
[student45@msba-hadoop-name ManideepClassProject]$ gunzip *-99999-*.gz
[student45@msba-hadoop-name ManideepClassProject]$ ls
011060-99999-1928 014270-99999-1928 028360-99999-1925 029110-99999-1922 029350-99999-1925 030910-99999-1928
011060-99999-1929 014270-99999-1929 028360-99999-1926 029110-99999-1923 029350-99999-1926 032620-99999-1927
011060-99999-1930 014270-99999-1930 028970-99999-1921 029110-99999-1924 029700-99999-1921 033020-99999-1927
012620-99999-1928 023610-99999-1929 028970-99999-1922 029110-99999-1925 029700-99999-1922 034970-99999-1927
012620-99999-1929 023610-99999-1930 028970-99999-1923 029110-99999-1926 029700-99999-1923 038040-99999-1927
012620-99999-1930 023610-99999-1921 028970-99999-1924 029350-99999-1921 029700-99999-1924 hadoop-streaming-2.7.3.jar
014030-99999-1928 028360-99999-1922 028970-99999-1925 029350-99999-1922 029700-99999-1925 Mapper45.py
014030-99999-1929 028360-99999-1923 028970-99999-1926 029350-99999-1923 029700-99999-1926 Reducer45.py
014030-99999-1930 028360-99999-1924 029110-99999-1921 029350-99999-1924 030910-99999-1927
[student45@msba-hadoop-name ManideepClassProject]$
```

## Step:3

Creating new directory to save input files, copying the input file to the directory:

```
hdfs dfs -mkdir /home/45student45/input_project
```

```
hdfs dfs -copyFromLocal /home/student45/ManideepClassProject/*-99999-*
/home/45student45/input_project/
```

```
[student45@msba-hadoop ~] + *
[student45@msba-hadoop-name ManideepClassProject]$ hdfs dfs -mkdir /home/45student45/input_project/
[student45@msba-hadoop-name ManideepClassProject]$ hdfs dfs -copyFromLocal /home/student45/ManideepClassProject/*-99999-* /home/45student45/input_project/
[student45@msba-hadoop-name ManideepClassProject]$ hdfs -ls /home/45student45/input_project/
Found 59 items:
-rw-r--r-- 5 student45 supergroup 18946 2023-05-05 22:34 /home/45student45/input_project/011060-99999-1928
-rw-r--r-- 5 student45 supergroup 19934 2023-05-05 22:34 /home/45student45/input_project/011060-99999-1929
-rw-r--r-- 5 student45 supergroup 19971 2023-05-05 22:34 /home/45student45/input_project/011060-99999-1930
-rw-r--r-- 5 student45 supergroup 8726 2023-05-05 22:34 /home/45student45/input_project/012620-99999-1928
-rw-r--r-- 5 student45 supergroup 9657 2023-05-05 22:34 /home/45student45/input_project/012620-99999-1929
-rw-r--r-- 5 student45 supergroup 10854 2023-05-05 22:34 /home/45student45/input_project/012620-99999-1930
-rw-r--r-- 5 student45 supergroup 13933 2023-05-05 22:34 /home/45student45/input_project/014030-99999-1928
-rw-r--r-- 5 student45 supergroup 59908 2023-05-05 22:34 /home/45student45/input_project/014030-99999-1929
-rw-r--r-- 5 student45 supergroup 119256 2023-05-05 22:34 /home/45student45/input_project/014030-99999-1930
-rw-r--r-- 5 student45 supergroup 9948 2023-05-05 22:34 /home/45student45/input_project/014270-99999-1928
-rw-r--r-- 5 student45 supergroup 9942 2023-05-05 22:34 /home/45student45/input_project/014270-99999-1929
-rw-r--r-- 5 student45 supergroup 9936 2023-05-05 22:34 /home/45student45/input_project/014270-99999-1930
-rw-r--r-- 5 student45 supergroup 37058 2023-05-05 22:34 /home/45student45/input_project/023610-99999-1928
-rw-r--r-- 5 student45 supergroup 37058 2023-05-05 22:34 /home/45student45/input_project/023610-99999-1929
-rw-r--r-- 5 student45 supergroup 151504 2023-05-05 22:34 /home/45student45/input_project/023610-99999-1931
-rw-r--r-- 5 student45 supergroup 152226 2023-05-05 22:34 /home/45student45/input_project/028360-99999-1922
-rw-r--r-- 5 student45 supergroup 151749 2023-05-05 22:34 /home/45student45/input_project/028360-99999-1923
-rw-r--r-- 5 student45 supergroup 153831 2023-05-05 22:34 /home/45student45/input_project/028360-99999-1924
-rw-r--r-- 5 student45 supergroup 150831 2023-05-05 22:34 /home/45student45/input_project/028360-99999-1925
-rw-r--r-- 5 student45 supergroup 151838 2023-05-05 22:34 /home/45student45/input_project/028360-99999-1926
-rw-r--r-- 5 student45 supergroup 147469 2023-05-05 22:34 /home/45student45/input_project/028970-99999-1921
-rw-r--r-- 5 student45 supergroup 148648 2023-05-05 22:34 /home/45student45/input_project/028970-99999-1922
-rw-r--r-- 5 student45 supergroup 151594 2023-05-05 22:34 /home/45student45/input_project/028970-99999-1923
-rw-r--r-- 5 student45 supergroup 152911 2023-05-05 22:34 /home/45student45/input_project/028970-99999-1924
-rw-r--r-- 5 student45 supergroup 151436 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1925
-rw-r--r-- 5 student45 supergroup 119046 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1926
-rw-r--r-- 5 student45 supergroup 150955 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1921
-rw-r--r-- 5 student45 supergroup 150762 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1922
-rw-r--r-- 5 student45 supergroup 151936 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1923
-rw-r--r-- 5 student45 supergroup 136425 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1924
-rw-r--r-- 5 student45 supergroup 148719 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1925
-rw-r--r-- 5 student45 supergroup 148629 2023-05-05 22:34 /home/45student45/input_project/029110-99999-1926
-rw-r--r-- 5 student45 supergroup 151947 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1921
-rw-r--r-- 5 student45 supergroup 149122 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1922
-rw-r--r-- 5 student45 supergroup 148971 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1923
-rw-r--r-- 5 student45 supergroup 149467 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1924
-rw-r--r-- 5 student45 supergroup 152108 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1925
-rw-r--r-- 5 student45 supergroup 152487 2023-05-05 22:34 /home/45student45/input_project/029350-99999-1926
-rw-r--r-- 5 student45 supergroup 152081 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1921
-rw-r--r-- 5 student45 supergroup 152708 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1922
-rw-r--r-- 5 student45 supergroup 152707 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1923
-rw-r--r-- 5 student45 supergroup 152707 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1923
-rw-r--r-- 5 student45 supergroup 152754 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1924
-rw-r--r-- 5 student45 supergroup 152978 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1925
-rw-r--r-- 5 student45 supergroup 152522 2023-05-05 22:34 /home/45student45/input_project/029700-99999-1926
-rw-r--r-- 5 student45 supergroup 30057 2023-05-05 22:34 /home/45student45/input_project/030910-99999-1927
```

## Step:4

## Displaying contents of the sample file:

```
hdfs dfs -cat /home/45student45/input_project/011060-99999-1928
```

### Step:5

### Executing python file using the jar file

hadoop jar hadoop-streaming-2.7.3.jar \

-file /home/student45/ManideepClassProject/Mapper45.py \

```
-mapper /home/student45/ManideepClassProject/Mapper45.py \
```

```
-file /home/student45/ManideepClassProject/Reducer45.py \
```

-reducer /home/student45/ManideepClassProject/Reducer45.py`

-input /home/45student45/input\_project/ \

-output /home/45student45/output project

```

student45@msba-hadoop: ~ + 
[student45@msba-hadoop: ~] hadoop jar hadoop-streaming-2.7.3.jar \
> -file /home/student5/ManideepClassProject/Mapper40.py \
> -mapper "/home/student5/ManideepClassProject/Mapper40.py" \
> -reducer "/home/student5/ManideepClassProject/Reducer40.py" \
> -input "/home/45student5/Input_project" \
> -output "/home/45student5/Output_project" \
23/05/05 22:38:41 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/student45/ManideepClassProject/Mapper40.py, /home/student5/ManideepClassProject/Reducer40.py, /tmp/hadoop-usjar52137676643084419784/] [] /tmp/streamjob569799392846532175.jar tmpDir=null
23/05/05 22:38:42 INFO client.MiniDFSPosix: Connecting to ResourceManager at /127.0.0.1:8083
23/05/05 22:38:43 INFO mapred.FileInputFormat: Total input files to process : 50
23/05/05 22:38:43 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
23/05/05 22:38:43 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1669743171306_3688
23/05/05 22:38:43 INFO impl.YarnClientImpl: Submitted application: application_1669743171306_3688/
23/05/05 22:38:43 INFO mapreduce.Job: Job job_1669743171306_3688 running in uber mode : false
23/05/05 22:38:43 INFO mapreduce.Job: user code starting in uber mode
23/05/05 22:38:43 INFO mapreduce.Job: map 4% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 12% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 18% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 24% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 30% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 36% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 42% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 48% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 54% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 60% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 66% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 72% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 78% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 84% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 90% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 96% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 100% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: Job job_1669743171306_3688 running in uber mode : false
23/05/05 22:38:43 INFO mapreduce.Job: user code starting in uber mode
23/05/05 22:38:43 INFO mapreduce.Job: map 4% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 12% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 18% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 24% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 30% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 36% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 42% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 48% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 54% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 60% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 66% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 72% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 78% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 84% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 90% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 96% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: map 100% reduce 0%
23/05/05 22:38:43 INFO mapreduce.Job: Job job_1669743171306_3688 completed successfully
23/05/05 22:38:43 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=37527
    FILE: Number of bytes written=16468183
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=37527
    HDFS: Number of bytes written=99
    HDFS: Number of read operations=153
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=0
  Job Counters
    Launched map tasks=50

```

```

student45@msba-hadoop: ~ + 
[student45@msba-hadoop: ~] hadoop jar hadoop-streaming-2.7.3.jar \
> -file /home/student5/ManideepClassProject/Mapper40.py \
> -mapper "/home/student5/ManideepClassProject/Mapper40.py" \
> -reducer "/home/student5/ManideepClassProject/Reducer40.py" \
> -input "/home/45student5/Input_project" \
> -output "/home/45student5/Output_project" \
23/05/05 22:39:58 INFO mapreduce.Job: map 86% reduce 27%
23/05/05 22:39:58 INFO mapreduce.Job: map 88% reduce 27%
23/05/05 22:39:58 INFO mapreduce.Job: map 90% reduce 27%
23/05/05 22:39:58 INFO mapreduce.Job: map 92% reduce 30%
23/05/05 22:39:57 INFO mapreduce.Job: map 98% reduce 30%
23/05/05 22:39:57 INFO mapreduce.Job: map 100% reduce 30%
23/05/05 22:40:08 INFO mapreduce.Job: map 100% reduce 100%
23/05/05 22:40:08 INFO mapreduce.Job: Job job_1669743171306_3688 completed successfully
23/05/05 22:40:08 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=37527
    FILE: Number of bytes written=16468183
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes written=99
    HDFS: Number of read operations=153
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=0
  Job Counters
    Launched map tasks=50
    Launched reduce tasks=1
    Data-local map tasks=50
    Total time spent by all maps in occupied slots (ms)=292666
    Total time spent by all mappers in occupied slots (ms)=292666
    Total time spent by all map reduce tasks (ms)=52834
    Total time spent by all reduce tasks (ms)=2838
    Total vcore-milliseconds taken by all map tasks=292464
    Total vcore-milliseconds taken by all reduce tasks=52834
    Total mapbytes-milliseconds taken by all map tasks=29983136
    Total mapbytes-milliseconds taken by all reduce tasks=53282816
  Map-Reduce Framework
    Map input records=36404
    Map output records=36404
    Map output bytes=156911
    Map output materialized bytes=37821
    Input split bytes=9680
    Combine output records=8
    Reduce input groups=13
    Reduce output groups=13
    Reduce input records=37821
    Reduce input records=3411
    Reduce output records=12
    Spilled records=2222
    Shuffled Maps =59
    Failed Shuffles=0
    Merged Map output=59
    Co-located map tasks=11
    CPU time spent (ms)=35880
    Physical memory (bytes) snapshot=17714833560
    Virtual memory (bytes) snapshot=184148977328
    Total committed heap usage (bytes)=15611199488
  Shuffle Errors
    File Input errors=0
    Bytes Read=517759
    Bytes Written=99
  File Output Format Counters
    Bytes Written=99
23/05/05 22:40:08 INFO streaming.StreamJob: Output directory: /home/45student5/output_project/

```

## Step:6

Displaying the final results:

```
hdfs dfs -ls /home/45student45/output_project/
```

```
hdfs dfs -cat /home/45student45/output_project/part-00000
```

```
[student45@msba-hadoop-name ManideepClassProject]$ hdfs dfs -ls /home/45student45/output_project/
Found 2 items
-rw-r--r-- 5 student45 supergroup          0 2023-05-05 22:39 /home/45student45/output_project/_SUCCESS
-rw-r--r-- 5 student45 supergroup        99 2023-05-05 22:39 /home/45student45/output_project/part-00000
[student45@msba-hadoop-name ManideepClassProject]$ hdfs dfs -cat /home/45student45/output_project/part-00000
1      21985
2      21985
3      21985
4      21985
5      21985
6      21985
7      21985
8      21985
9      21985
10     21985
11     21985
12     21985
[student45@msba-hadoop-name ManideepClassProject]$
```

## Part 2:

The second part is to develop a python application that can be implemented in PySpark to calculate the *average visibility distance* (meters) for *each USAF weather station ID* from NCDC records (note: 999999 indicates missing value, and [01459] indicate good quality value).

### Step:1

Unzipping sample files and creating a new directory:

```
gunzip *-99999-* .gz
```

```
hdfs dfs -mkdir /user/hadoop/input_2
```

```
Last login: Fri May  5 23:00:27 on ttys003
(base) manideepelasagaram@Manideeps-MacBook-Pro ~ % cd Desktop
(base) manideepelasagaram@Manideeps-MacBook-Pro Desktop % ssh -i BAN632.pem hadoop@ec2-34-205-90-229.compute-1.amazonaws.com
--| --|_
 _| (   /   Amazon Linux 2 AMI
---\---|---|_

https://aws.amazon.com/amazon-linux-2/
6 package(s) needed for security, out of 18 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEE MMMMMMM          MMMMMMM RRRRRRRRRRRRRR
E:::::::::::E M::::::M      M::::::M R:::::::::::R
EE:::::E:::::E M:::::M      M:::::M R:::::RRRRR:::::R
E:::E     EEEE M:::::M      M:::::M R:::::R    R:::::R
E:::E     M:::::M:M      M:::M:M M:::::M R:::::R    R:::::R
E:::::EEEEEEEEE M:::::M M:::::M M:::::M R:::::RRRRR:::::R
E:::::::::::E M:::::M M:::::M:M      M:::::M R:::::::::::R
E:::::EEEEEEEEE M:::::M M:::::M      M:::::M R:::::RRRRR:::::R
E:::E     M:::::M M:::::M      M:::::M R:::::R    R:::::R
E:::::E     EEEE M:::::M      M:::::M R:::::R    R:::::R
EE:::::E:::::E M:::::M      M:::::M R:::::R    R:::::R
E:::::::::::E M:::::M      M:::::M R:::::R    R:::::R
EEEEEEEEEEEEEEEEE MMMMMMM          MMMMMMM RRRRRRR    RRRRRR

[hadoop@ip-172-31-4-169 ~]$ gunzip *-99999-* .gz
[hadoop@ip-172-31-4-169 ~]$ hdfs dfs -mkdir /user/hadoop/input_2
```

### Step:2

Copying sample files from local to Hadoop:

```
hdfs dfs -copyFromLocal *-99999-* /user/hadoop/input_2
```

```
[hadoop@ip-172-31-4-169 ~]$ hdfs dfs -copyFromLocal *-99999-* /user/hadoop/input_2
```

### Step:3

Displaying the copied files:

```
hdfs dfs -ls /user/hadoop/input_2
```

```
[hadoop@ip-172-31-4-169 ~]$ hdfs dfs -ls /user/hadoop/input_2
Found 50 items
-rw-r--r-- 1 hadoop hdfsadmingroup 18946 2023-05-06 06:07 /user/hadoop/input_2/011060-99999-1928
-rw-r--r-- 1 hadoop hdfsadmingroup 19934 2023-05-06 06:07 /user/hadoop/input_2/011060-99999-1929
-rw-r--r-- 1 hadoop hdfsadmingroup 19971 2023-05-06 06:07 /user/hadoop/input_2/011060-99999-1930
-rw-r--r-- 1 hadoop hdfsadmingroup 8726 2023-05-06 06:07 /user/hadoop/input_2/012620-99999-1928
-rw-r--r-- 1 hadoop hdfsadmingroup 9657 2023-05-06 06:07 /user/hadoop/input_2/012620-99999-1929
-rw-r--r-- 1 hadoop hdfsadmingroup 10054 2023-05-06 06:07 /user/hadoop/input_2/012620-99999-1930
-rw-r--r-- 1 hadoop hdfsadmingroup 13033 2023-05-06 06:07 /user/hadoop/input_2/014030-99999-1928
-rw-r--r-- 1 hadoop hdfsadmingroup 59900 2023-05-06 06:07 /user/hadoop/input_2/014030-99999-1929
-rw-r--r-- 1 hadoop hdfsadmingroup 119256 2023-05-06 06:07 /user/hadoop/input_2/014030-99999-1930
-rw-r--r-- 1 hadoop hdfsadmingroup 9948 2023-05-06 06:07 /user/hadoop/input_2/014270-99999-1928
-rw-r--r-- 1 hadoop hdfsadmingroup 9442 2023-05-06 06:07 /user/hadoop/input_2/014270-99999-1929
-rw-r--r-- 1 hadoop hdfsadmingroup 9998 2023-05-06 06:07 /user/hadoop/input_2/014270-99999-1930
-rw-r--r-- 1 hadoop hdfsadmingroup 37353 2023-05-06 06:07 /user/hadoop/input_2/023610-99999-1929
-rw-r--r-- 1 hadoop hdfsadmingroup 87958 2023-05-06 06:07 /user/hadoop/input_2/023610-99999-1930
-rw-r--r-- 1 hadoop hdfsadmingroup 151540 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1921
-rw-r--r-- 1 hadoop hdfsadmingroup 152226 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1922
-rw-r--r-- 1 hadoop hdfsadmingroup 151749 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1923
-rw-r--r-- 1 hadoop hdfsadmingroup 153031 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1924
-rw-r--r-- 1 hadoop hdfsadmingroup 150831 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1925
-rw-r--r-- 1 hadoop hdfsadmingroup 151838 2023-05-06 06:07 /user/hadoop/input_2/028360-99999-1926
-rw-r--r-- 1 hadoop hdfsadmingroup 147469 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1921
-rw-r--r-- 1 hadoop hdfsadmingroup 148648 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1922
-rw-r--r-- 1 hadoop hdfsadmingroup 151594 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1923
-rw-r--r-- 1 hadoop hdfsadmingroup 152011 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1924
-rw-r--r-- 1 hadoop hdfsadmingroup 151430 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1925
-rw-r--r-- 1 hadoop hdfsadmingroup 119414 2023-05-06 06:07 /user/hadoop/input_2/028970-99999-1926
-rw-r--r-- 1 hadoop hdfsadmingroup 150855 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1921
-rw-r--r-- 1 hadoop hdfsadmingroup 150762 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1922
-rw-r--r-- 1 hadoop hdfsadmingroup 151938 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1923
-rw-r--r-- 1 hadoop hdfsadmingroup 136425 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1924
-rw-r--r-- 1 hadoop hdfsadmingroup 148719 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1925
-rw-r--r-- 1 hadoop hdfsadmingroup 148629 2023-05-06 06:07 /user/hadoop/input_2/029110-99999-1926
-rw-r--r-- 1 hadoop hdfsadmingroup 151947 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1921
-rw-r--r-- 1 hadoop hdfsadmingroup 149122 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1922
-rw-r--r-- 1 hadoop hdfsadmingroup 148971 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1923
-rw-r--r-- 1 hadoop hdfsadmingroup 149467 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1924
-rw-r--r-- 1 hadoop hdfsadmingroup 152100 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1925
-rw-r--r-- 1 hadoop hdfsadmingroup 152487 2023-05-06 06:07 /user/hadoop/input_2/029350-99999-1926
-rw-r--r-- 1 hadoop hdfsadmingroup 152001 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1921
-rw-r--r-- 1 hadoop hdfsadmingroup 151532 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1922
-rw-r--r-- 1 hadoop hdfsadmingroup 152747 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1923
-rw-r--r-- 1 hadoop hdfsadmingroup 152754 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1924
-rw-r--r-- 1 hadoop hdfsadmingroup 152978 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1925
-rw-r--r-- 1 hadoop hdfsadmingroup 152522 2023-05-06 06:07 /user/hadoop/input_2/029700-99999-1926
-rw-r--r-- 1 hadoop hdfsadmingroup 30057 2023-05-06 06:07 /user/hadoop/input_2/030910-99999-1927
-rw-r--r-- 1 hadoop hdfsadmingroup 34042 2023-05-06 06:07 /user/hadoop/input_2/030910-99999-1928
-rw-r--r-- 1 hadoop hdfsadmingroup 33359 2023-05-06 06:07 /user/hadoop/input_2/032620-99999-1927
-rw-r--r-- 1 hadoop hdfsadmingroup 36616 2023-05-06 06:07 /user/hadoop/input_2/033020-99999-1927
-rw-r--r-- 1 hadoop hdfsadmingroup 35171 2023-05-06 06:07 /user/hadoop/input_2/034970-99999-1927
-rw-r--r-- 1 hadoop hdfsadmingroup 36622 2023-05-06 06:07 /user/hadoop/input_2/038040-99999-1927
```

### Step:4

Executing the python file:

```
spark-submit --master yarn AVD_Sparkfile.py
```

```
[hadoop@ip-172-31-4-169 ~]$ spark-submit --master yarn AVO_SparkHdfs.py
23/05/06 06:20:18 INFO SparkContext: Running Spark version 2.4.8-amzn-2
23/05/06 06:20:18 INFO SparkContext: Submitted application: SparkVisibility
23/05/06 06:20:18 INFO SecurityManager: Changing view acls to: hedop
23/05/06 06:20:18 INFO SecurityManager: Changing view acls groups to:
23/05/06 06:20:18 INFO SecurityManager: Changing exec acls groups to:
23/05/06 06:20:18 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(hadoop); groups with view permissions: Set(); users with modify permissions: Set(hadoop); groups with m
23/05/06 06:20:18 INFO Util: Successfully started service 'sparkDriver' on port 42599.
23/05/06 06:20:18 INFO SparkEnv: Registering BlockManagerMaster
23/05/06 06:20:18 INFO SparkEnv: Registering BlockManagerMaster
23/05/06 06:20:18 INFO SparkEnv: Registering BlockManagerMaster
23/05/06 06:20:18 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
23/05/06 06:20:11 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-3c0e9b9-9428-43cf-a6fc-0b9e956d5b37
23/05/06 06:20:11 INFO MemoryStore: MemoryStore started with capacity 912.3 MB
23/05/06 06:20:11 INFO MemoryStore: BlockManagerId(mapper_1, ip-172-31-4-169, 0)
23/05/06 06:20:11 INFO Util: Successfully started service 'SparkUI' on port 4048.
23/05/06 06:20:11 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-4-169.ec2.internal:4048
23/05/06 06:20:11 INFO Util: Using 100 prelocated executors (minExecutors: 0). Set spark.dynamicallocation.preallocateExecutors to 'false' disable executor preallocation.
23/05/06 06:20:11 INFO Util: Using 100 prelocated executors (minExecutors: 0). Set spark.dynamicallocation.preallocateExecutors to 'false' disable executor preallocation.
23/05/06 06:20:11 INFO Client: Requesting a new application from cluster with 1 NodeManagers
23/05/06 06:20:12 INFO Configuration: resource-types.xml not found
23/05/06 06:20:12 INFO ResourceUtil: Unable to find 'resource-types.xml'.
23/05/06 06:20:12 INFO ResourceUtil: Adding resource type - name = memory_mb, units = Mi, type = COUNTABLE
23/05/06 06:20:12 INFO ResourceUtil: Adding resource type - name = vcores, units = ., type = COUNTABLE
23/05/06 06:20:12 INFO Client: Verifying our application has not requested more than the maximum memory capability of the cluster (12288 MB per container)
23/05/06 06:20:12 INFO Client: Will allocate AM container with 896 MB memory including 384 MB overhead
23/05/06 06:20:12 INFO Client: Setting up the launch environment for our AM container
23/05/06 06:20:12 INFO Client: Setting up the launch environment for our AM container
23/05/06 06:20:12 INFO Client: Waiting for AM container to launch, falling back to uploading libraries under SPARK_HOME.
23/05/06 06:20:12 INFO Client: Uploading resource file:/mnt/tmp/spark-374d1ed-d286-4c18-835c-0ef719f9fad6/_spark_lbs_74283283464776533.zip -> hdfs://ip-172-31-4-169.ec2.internal:8020/user/hadoop/.sparkStaging/application_1683352956191_0001/_spark_lbs_74283283464776533.zip
23/05/06 06:20:24 INFO Client: Uploading resource file:/etc/hadoop/conf/slaves.conf -> hdfs://ip-172-31-4-169.ec2.internal:8020/user/hadoop/.sparkStaging/application_1683352956191_0001/hud-defaults.conf
23/05/06 06:20:24 INFO Client: Uploading resource file:/usr/lib/python/lib/pyspark.zip -> hdfs://ip-172-31-4-169.ec2.internal:8020/user/hadoop/.sparkStaging/application_1683352956191_0001/pyspark.zip
23/05/06 06:20:24 INFO Client: Uploading resource file:/usr/lib/python/lib/pyspark/lib/py4j-0.10.7-src.zip -> hdfs://ip-172-31-4-169.ec2.internal:8020/user/hadoop/.sparkStaging/application_1683352956191_0001/py4j-0.10.7-src.zip
23/05/06 06:20:24 INFO Client: Uploading resource file:/mnt/tmp/spark-374d1ed-d286-4c18-835c-0ef719f9fad6/_spark_conf_593527854/27859413.zip -> hdfs://ip-172-31-4-169.ec2.internal:8020/user/hadoop/.sparkStaging/application_1683352956191_0001/_spark_conf_593527854/27859413.zip
23/05/06 06:20:24 INFO SecurityManager: Changing view acls to: hedop
23/05/06 06:20:24 INFO SecurityManager: Changing exec acls groups to:
23/05/06 06:20:24 INFO SecurityManager: Changing modif acls groups to:
23/05/06 06:20:24 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(hadoop); groups with view permissions: Set(); users with modify permissions: Set(hadoop); groups with m
23/05/06 06:20:24 INFO Client: Submitting application application_1683352956191_0001 to ResourceManager
23/05/06 06:20:26 INFO YarnClientImpl: Submitted application application_1683352956191_0001
23/05/06 06:20:26 INFO SchedulerExtensionServices: Starting Yarn extension services with app application_1683352956191_0001 and attemptId None
23/05/06 06:20:26 INFO Client: Application report for application application_1683352956191_0001 (state: ACCEPTED)
23/05/06 06:20:26 INFO Client:
  client token: N/A
  diagnostics: AM container is launched, waiting for AM container to Register with RM
  ApplicationMaster host: N/A
  ApplicationMaster port: -1
  queue: default
  start time: 168335402405630
  final status: UNDEFINED
  tracking URL: http://ip-172-31-4-169.ec2.internal:20888/proxy/application_1683352956191_0001/
  user: hadoop

23/05/06 06:20:27 INFO Client: Application report for application application_1683352956191_0001 (state: ACCEPTED)
23/05/06 06:20:27 INFO Client: Application report for application application_1683352956191_0001 (state: ACCEPTED)
23/05/06 06:20:29 INFO Client: Application report for application application_1683352956191_0001 (state: ACCEPTED)
23/05/06 06:20:30 INFO Client: Application report for application application_1683352956191_0001 (state: RUNNING)
23/05/06 06:20:30 INFO Client:
  client token: N/A
  diagnostics: N/A
```

### Step:5

### Displaying the output files:

```
hdfs dfs -ls /user/hadoop/
```

## Step:6

Displaying the contents of the output files combined:

```
hdfs dfs -cat /user/hadoop/output_2/output_q2.txt/part-000*
```

```
[hadoop@ip-172-31-4-169 ~]$ hdfs dfs -cat /user/hadoop/output_2/output_q2.txt/part-000*
('023610', 37068.553459119496)
('029350', 0.0)
('014270', 17137.426900584796)
('034970', 5803.20197044335)
('012620', 26542.331288343557)
('033020', 12318.483412322275)
('029110', 0.0)
('028970', 0.0)
('032620', 8316.497461928933)
('029700', 0.0)
('014030', 33686.024844720494)
('030910', 11362.198391420912)
('028360', 0.0)
('011060', 24848.672566371682)
('038040', 14158.064516129032)
[hadoop@ip-172-31-4-169 ~]$ █
```

### Part 3:

The third part is to load the text file into Pig and get the range of sky ceiling height for each USAF weather station ID.

Step:1 records\_ch = LOAD 'pig/skyceiling\_data.txt'

AS (USAF\_stationID:chararray, ceiling\_height:int);

```
grunt> records_ch = LOAD 'pig/skyceiling_data.txt'  
->   AS (USAF_stationID:chararray, ceiling_height:int);  
2023-05-04 19:54:03,912 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum  
grunt> ■
```

### Step:2 DUMP records\_ch;

```
(#38848, 22880)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 458)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 788)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 3688)  
(#38848, 1238)  
(#38848, 1238)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 1238)  
(#38848, 22880)  
(#38848, 788)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 22880)  
(#38848, 22880)  
grunt: ■
```

Step:3 DESCRIBE records\_ch;

```
[8888:16:22888]
[grunt> DESCRIBE records_ch;
records_ch: {USAF_stationID: chararray,ceiling_height: int}
grunt>
```

Step:4 filtered records ch = FILTER records ch BY ceiling height != 99999;

DUMP filtered records ch;

Step:5 grouped records ch = GROUP filtered records ch BY USAF\_stationID:

DUMP grouped records ch;

Step:6 DESCRIBE grouped records ch;

```
|grunt> DESCRIBE grouped_records_ch;
grouped_records_ch: {group: chararray, filtered_records_ch: {{(USAF_stationID: chararray,ceiling_height: int)}}
grunt> |
```

Step:7 min ceiling height = FOREACH grouped records ch GENERATE group

MIN(filtered records ch.ceiling height):

```
grunt> min_ceiling_height = FOREACH grouped_records_ch GENERATE group
>>> MIN(filtered_records_ch.ceiling_height);
grunt>
```

## DUMP max\_ceiling\_height;

```

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.9.0 0.17.0 student45 2023-05-04 20:00:15 2023-05-04 20:00:16 GROUP_BY,FILTER

Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_local1033654270_0004 1 1 n/a n/a n/a n/a n/a n/a filtered_records_ch,grouped_records_ch,min_ceiling_height,records_ch GROUP_BY,COMBINER file:/tmp/temp-863784154/tmp11

Input(s):
Successfully read 3391 records from: "file:///home/student45/pig/skycelling_data.txt"

Output(s):
Successfully stored 10 records in: "file:/tmp/temp-863784154/tmp1102737492"

Counters:
Total records written : 10
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1033654270_0004

2023-05-04 20:00:16,034 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:00:16,036 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:00:16,036 [main] INFO org.apache.hive.hbase.metrics.HBaseMetrics - Cannot initialize HBase Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:00:16,038 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2023-05-04 20:00:16,039 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2023-05-04 20:00:16,040 [main] WARN org.apache.pig.data.SchemaUtil - schema@local has already been initialized
2023-05-04 20:00:16,044 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2023-05-04 20:00:16,054 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1

(81104,15)
(81105,15)
(81405,15)
(814270,15)
(82360,15)
(83020,15)
(83260,240)
(833020,15)
(834970,15)
(83860,60)
grants: []

```

Step:8 max\_ceiling\_height = FOREACH grouped\_records\_ch GENERATE group,

MAX(filtered\_records\_ch.ceiling\_height);

```

grunt> max_ceiling_height = FOREACH grouped_records_ch GENERATE group,
>>   MAX(filtered_records_ch.ceiling_height);

```

DUMP max\_ceiling\_height;

```

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.9.0 0.17.0 student45 2023-05-04 20:01:34 2023-05-04 20:01:34 GROUP_BY,FILTER

Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_local172565997_0005 1 n/a n/a n/a n/a n/a n/a n/a filtered_records_ch,grouped_records_ch,max_ceiling_height,records_ch GROUP_BY,COMBINER file:/tmp/temp-863784154/tmp156975719.

Input(s):
Successfully read 3391 records from: "file:///home/student45/pig/skycelling_data.txt"

Output(s):
Successfully stored 10 records in: "file:/tmp/temp-863784154/tmp156975719"

Counters:
Total records written : 10
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local172565997_0005

2023-05-04 20:01:34,566 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:01:34,569 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:01:34,569 [main] INFO org.apache.hive.hbase.metrics.HBaseMetrics - Cannot initialize HBase Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:01:34,572 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2023-05-04 20:01:34,572 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2023-05-04 20:01:34,585 [main] INFO org.apache.pig.data.SchemaUtil - schema@local has already been initialized
2023-05-04 20:01:34,585 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2023-05-04 20:01:34,585 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1

(81104,22800)
(81105,22800)
(81405,22800)
(814270,22800)
(82360,22800)
(83020,22800)
(83260,22800)
(833020,22800)
(834970,22800)
(83860,22800)
grants: []

```

Step:9 fheight\_range = FOREACH grouped\_records\_ch GENERATE group AS USAF\_stationID, (max\_ceiling\_height - min\_ceiling\_height) AS range\_ch;

```
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.9.0 0.17.0 student42 2023-05-04 20:03:43 2023-05-04 20:03:44 GROUP_BY

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime Alias Feature Outputs
job_local363916344_0002 1 n/a n/a n/a n/a n/a n/a Range_height_end,grouped_records,records GROUP_BY,COMBINER file:/tmp/temp1415349613/tmp-109594995

Input(s):
Successfully read 3391 records from: "file:///home/student42/pig/skyceiling_data.txt"

Output(s):
Successfully stored 10 records in: "file:/tmp/temp1415349613/tmp-109594995"

Counters:
Total records written : 10
Total bytes written : 8
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local363916344_0002

2023-05-04 20:03:44,294 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:03:44,295 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:03:44,296 [main] INFO org.apache.hadoop.metrics.jvm.JvmMetrics - Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
2023-05-04 20:03:44,297 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2023-05-04 20:03:44,298 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2023-05-04 20:03:44,299 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input blocks (128) to process : 1
2023-05-04 20:03:44,300 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2023-05-04 20:03:44,311 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2023-05-04 20:03:44,311 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input blocks (128) to process : 1
2023-05-04 20:03:44,311 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
(01060,21985)
(012620,21985)
(014030,21985)
(015050,21985)
(023610,21985)
(032620,21760)
(033020,21985)
(034970,21985)
(038040,21940)
grunt> █
```

Part 4:

The fourth part is to load the text file into Hive and get the average sky ceiling height for each USAF weather station ID.

Step:1 DROP TABLE IF EXISTS records45;

```
[hive> DROP TABLE IF EXISTS records45;
OK
Time taken: 6.291 seconds
hive> ]
```

Step:2 CREATE TABLE records45 (USAF\_stationID STRING, ceiling\_height INT)

ROW FORMAT DELIMITED

FIELDS TERMINATED BY '\t';

```
hive> CREATE TABLE records45 (USAF_stationID STRING, ceiling_height INT)
      > ROW FORMAT DELIMITED
      [   >   FIELDS TERMINATED BY '\t';
OK
Time taken: 0.56 seconds
hive> ]
```

Step:3 LOAD DATA LOCAL INPATH 'pig/skyceiling\_data.txt'

OVERWRITE INTO TABLE records45;

```
hive> LOAD DATA LOCAL INPATH 'pig/skyceiling_data.txt'
[   > OVERWRITE INTO TABLE records45;
Loading data to table default.records45
OK
Time taken: 0.531 seconds
hive> ]
```

Step:4 SELECT USAF\_stationID, AVG(ceiling\_height)

FROM records45

WHERE ceiling\_height != 9999

GROUP BY USAF\_stationID;

```

hive> SELECT USAF_stationID, AVG(ceiling_height)
    > FROM records45
    > WHERE ceiling_height != 9999
    |> GROUP BY USAF_stationID;
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: client45_20230504194858_b9f6c831-afc5-4bd7-9ebc-013fbcc0b81cd
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=<n>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<n>
Starting Job = job_1669743171306_3287, Tracking URL: http://msba-hadoop-name:8088/proxy/application_1669743171306_3287/
Kill Command = /usr/local/hadoop/bin/hadoop job -kill job_1669743171306_3287
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2023-05-04 19:49:05,378 Stage-1 map = 0%, reduce = 0%
2023-05-04 19:49:17,805 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.35 sec
2023-05-04 19:49:17,805 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 6.55 sec
MapReduce Total cumulative CPU time: 5 seconds 550 msec
Ended Job = job_1669743171306_3287
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 5.55 sec HDFS Read: 49855 HDFS Write: 460 SUCCESS
Total MapReduce CPU Time Spent: 5 seconds 550 msec
OK
031010 10047.007377947101
032620 9879.468885104383
034090 9564.369282294651
034270 7448.4851724137936
032610 11340.65266579974
030910 11135.77388952381
032620 8984.336734693878
033020 8474.796511627987
034970 10384.42857142857
038040 10921.725888324872
Time taken: 20.636 seconds, Fetched: 10 row(s)
hive> 

```

You need to turn in:

**1) Part 1:**

- a. *if you are using JAVA to develop the Mapper and Reducer applications:* the three java files (mapper, reducer and main);
- b. *if you are using Hadoop streaming jar and developing two python programs (mapper python file and reducer python file):* the two python files (mapper and reducer);
- c. *if you are using mrjob library and developing one python program with two functions:* the python file (with the mapper and reducer functions);

**2) Part 2:** the python program you developed;

**3)** the commands from converting java files into a Jar file to running the Jar file in Hadoop, or the commands to execute the python files in Hadoop and in Spark;

**4)** the step by step commands and screenshots of solutions from all the parts;

The original dataset for this project is available on Blackboard.

