

1. Finish mapper.py and test the code

```
Last login: Mon Nov 16 11:05:17 on ttys000
(base) qmr@qumingruodeMBP ~ % /Users/qmr/Desktop/cmu_semester1/cloud
zsh: permission denied: /Users/qmr/Desktop/cmu_semester1/cloud
(base) qmr@qumingruodeMBP ~ % cd /Users/qmr/Desktop/cmu_semester1/cloud
(base) qmr@qumingruodeMBP cloud % cd temperature
(base) qmr@qumingruodeMBP temperature % cat data/1901 | python mapper.py
1901 -78
1901 -72
1901 -94
1901 -61
1901 -56
1901 -28
1901 -67
1901 -33
1901 -28
1901 -33
1901 -44
1901 -39
1901 0
1901 6
1901 0
1901 6
1901 6
1901 -11
1901 -33
1901 -50
1901 -44
1901 -28
1901 -33
1901 -33
1901 -50
1901 -33
1901 -28
1901 -44
1901 -44
1901 -44
1901 -39
1901 -50
1901 -44
1901 -39
1901 -33
1901 -22
1901 0
1901 -6
1901 -17
1901 -44
1901 -39
1901 -33
1901 -6
1901 17
1901 22
1901 28
1901 28
1901 11
1901 -17
1901 -28
1901 -56
1901 -44
1901 -44
1901 -67
1901 -44
1901 -39
1901 -22
```

2. Finish reducer.py and test the code, reduce on local machine

```
BrokenPipeError: [Errno 32] Broken pipe
(base) qmr@qumingruodeMacBook-Pro:~/temperature % cat data/1902 | python mapper.py | sort -k1,1 | python reducer.py
19020101      28
19020102      11
19020103     -61
19020104      28
19020105      22
19020106      11
19020107      28
19020108      28
19020109      28
19020110      33
19020111       6
19020112     -11
19020113      -6
19020114     -22
19020115     -11
19020116      28
19020117       0
19020118     -17
19020119      -6
19020120      22
19020121      11
19020122      -6
19020123       6
19020124      28
19020125      22
19020126      17
19020127      11
19020128      -6
19020129      -6
19020130     -33
19020131     -17
19020132      22
19020133      17
19020134     -17
19020135     -50
19020136     -11
19020137     -11
19020138       0
19020139     -44
19020140     -28
19020141     -61
19020142       0
19020143     -11
19020144     -33
19020145      -6
19020146      11
19020147      22
19020148       0
19020149     117
19020150     -22
19020151      -6
19020152       0
19020153       0
19020154       0
19020155       0
19020156       6
19020157     -6
19020158       6
19020159      17
19020160      22
```

3. Create a cluster on Data Proc and upload the files we need

bucket id: **dataproc-staging-us-central1-147654780204-veglhg7c**

The screenshot shows the Google Cloud Platform Storage Bucket details page for 'dataproc-staging-us-central1-147654780204-veglhg7c'. The bucket is located in 'us-central1 (Iowa)' with a 'Standard' storage class. It has 'Public access' set to 'Subject to object ACLs' and 'Protection' set to 'None'. The 'OBJECTS' tab is selected, displaying five files: '.DS_Store', 'data/', 'mapper.py', 'reducer.py', and '1902'. The 'CONFIGURATION', 'PERMISSIONS', 'PROTECTION', and 'LIFECYCLE' tabs are also present.

4. Open SSH

The screenshot shows the Google Cloud Platform Cluster details page for 'cluster-508f'. The cluster is a 'Dataproc Cluster' with a 'Running' status. The 'VM INSTANCES' tab is selected, listing three instances: 'cluster-508f-m' (Master), 'cluster-508f-w-0' (Worker), and 'cluster-508f-w-1' (Worker). An SSH context menu is open over the first instance, with the option 'Open in browser window' highlighted by a red circle.

```

cluster-508f - VM Instances - X mingruoqu123@cluster-508f-m + 
← → C ssh.cloud.google.com/projects/hadoop-mapreduce-332010/zones/us-central1-c/ins

Connected, host fingerprint: ssh-rsa 0 31:C3:D6:A5:85:E9:35:66:1D:D6:3C:42:7C:C4
:3A:DB:9E:D1:09:75:7C:8B:13:25:86:20:8E:CD:3E:D6:38:89
Linux cluster-508f-m 5.10.0-0.bpo.8-amd64 #1 SMP Debian 5.10.46-4~bpo10+1 (2021-08-07) x86_64

The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
mingruoqu123@cluster-508f-m:~$ gsutil ls
gs://dataproc-staging-us-central1-147654780204-veglhg7c/
gs://dataproc-temp-us-central1-147654780204-fxdwdnbo/
mingruoqu123@cluster-508f-m:~$ 

```

5. Moving the python code to Local Cluster (Namenode): create a new folder named '14848_temperature' and copy the uploaded file into the new folder through 'gsutil cp -r'

```

mingruoqu123@cluster-508f-m:~$ mkdir 14848_temperature
mingruoqu123@cluster-508f-m:~$ ls
14848_temperature
mingruoqu123@cluster-508f-m:~/14848_temperature$ cd 14848_temperature
mingruoqu123@cluster-508f-m:~/14848_temperature$ gsutil ls gs://dataproc-staging-us-central1-147654780204-veglhg7c/
-bash: gsutil: command not found
mingruoqu123@cluster-508f-m:~/14848_temperature$ gsutil ls gs://dataproc-staging-us-central1-147654780204-veglhg7c/
gs://dataproc-staging-us-central1-147654780204-veglhg7c/google-cloud-dataproc-metainfo/
gs://dataproc-staging-us-central1-147654780204-veglhg7c/notebooks/
gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/
mingruoqu123@cluster-508f-m:~/14848_temperature$ gsutil ls gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/
gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/.DS_Store
gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/mapper.py
gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/reducer.py
gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/data/
mingruoqu123@cluster-508f-m:~/14848_temperature$ gsutil cp -r gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/.DS_Store...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/data/1901...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/data/1902...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/mapper.py...
/ [4 files][ 1.7 MiB/ 1.7 MiB]
==> NOTE: You are performing a sequence of gsutil operations that may
run significantly faster if you instead use gsutil -m cp ... Please
see the -m section under "gsutil help options" for further information
about when gsutil -m can be advantageous.

Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/reducer.py...
/ [5 files][ 1.7 MiB/ 1.7 MiB]
Operation completed over 5 objects/1.7 MiB.
mingruoqu123@cluster-508f-m:~/14848_temperature$ 

```

Verify it is the python file we uploaded

```

mingruoqu123@cluster-508f-m:~/14848_temperature$ gsutil cp -r gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/ .
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/.DS_Store...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/data/1901...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/data/1902...
Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/mapper.py...
-[4 files][ 1.7 MiB/ 1.7 MiB]
==> NOTE: You are performing a sequence of gsutil operations that may
run significantly faster if you instead use gsutil -m cp ... Please
see the -m section under "gsutil help options" for further information
about when gsutil -m can be advantageous.

Copying gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/reducer.py...
-[5 files][ 1.7 MiB/ 1.7 MiB]
Operation completed over 5 objects/1.7 MiB.
mingruoqu123@cluster-508f-m:~/14848_temperature$ ls temperature
data mapper.py reducer.py
mingruoqu123@cluster-508f-m:~/14848_temperature$ cd temperature/
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ ls
data mapper.py reducer.py
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ cat mapper.py
import sys
valid_quality = [0,1,4,5,9]
valid_temperature = [9999]
for line in sys.stdin:
    line = line.strip()
    date = line[15:23]
    temperature = int(line[87:92])
    quality = int(line[92])
    if (temperature not in valid_temperature and quality in valid_quality):
        print('%s\t%d' % (date, temperature))
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ cat reducer.py

```

6. create data folder

```

print('%s\t%d' % (current_date, current_temperature))
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /
Found 3 items
drwxrwxrwt  - hadoop          0 2021-11-15 22:33 /tmp
drwxrwxrwt  - hadoop          0 2021-11-15 22:33 /user
drwx-wx-wx  - hive            0 2021-11-15 22:33 /var
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -put data/ /
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /
Found 4 items
drwxr-xr-x  - mingruoqu123 hadoop          0 2021-11-16 00:47 /data
drwxrwxrwt  - hdfs           hadoop          0 2021-11-15 22:33 /tmp
drwxrwxrwt  - hdfs           hadoop          0 2021-11-15 22:33 /user
drwx-wx-wx  - hive            hadoop          0 2021-11-15 22:33 /var
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /data
Found 2 items
-rw-r--r--  2 mingruoqu123 hadoop  888190 2021-11-16 00:47 /data/1901
-rw-r--r--  2 mingruoqu123 hadoop  888978 2021-11-16 00:47 /data/1902
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ █

```

7. the execution of Hadoop MapReduce Job in the terminal

```

mingruoqu123@cluster-508f:~/14848_temperature/temperature$ hadoop jar /usr/lib/hadoop/hadoop-streaming.jar -file mapper.py -mapper 'python mapper.py' -file reducer.py -reducer 'python reducer.py' -input /input -output /outputFolder5
2021-11-16 02:08:35.950 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [mapper.py, reducer.py] [/usr/lib/hadoop/hadoop-streaming-3.2.2.jar] /tmp/streamjob839252088802662895.jar tmpDir=null
2021-11-16 02:08:36.761 INFO client.RMProxy: Connecting to ResourceManager at cluster-508f-m/10.128.0.4:8032
2021-11-16 02:08:37.027 INFO client.AHSProxy: Connecting to Application History server at cluster-508f-m/10.128.0.4:10200
2021-11-16 02:08:37.520 INFO client.AHSProxy: Connecting to Application History server at cluster-508f-m/10.128.0.4:10200
2021-11-16 02:08:37.701 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/mingruoqu123/.staging/job_1637015590244_0003
2021-11-16 02:08:38.077 INFO mapred.FileInputFormat: Total input files to process : 2
2021-11-16 02:08:38.146 INFO mapreduce.JobSubmitter: number of splits:22
2021-11-16 02:08:38.288 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1637015590244_0003
2021-11-16 02:08:38.298 INFO mapreduce.JobSubmitter: Executing token with tokens: []
2021-11-16 02:08:38.335 INFO mapreduce.Job: Job configuration: conf.Configuration@535333
2021-11-16 02:08:38.536 INFO mapreduce.Job: User Warnings: Unable to find 'mapred-site-types.xml'.
2021-11-16 02:08:38.536 INFO mapreduce.Job: Impl: YarnClientImpl@600: Submitted application application_1637015590244_0003
2021-11-16 02:08:38.634 INFO mapreduce.Job: The url to track the job: http://cluster-508f-m:8088/proxy/application_1637015590244_0003/
2021-11-16 02:08:38.636 INFO mapreduce.Job: Running job: job_1637015590244_0003
2021-11-16 02:08:46.742 INFO mapreduce.Job: Job: job_1637015590244_0003 running in uber mode : false
2021-11-16 02:08:46.743 INFO mapreduce.Job: map 0% reduce 0%
2021-11-16 02:08:54.851 INFO mapreduce.Job: map 14% reduce 0%
2021-11-16 02:08:55.857 INFO mapreduce.Job: map 18% reduce 0%
2021-11-16 02:08:56.893 INFO mapreduce.Job: map 32% reduce 0%
2021-11-16 02:08:58.903 INFO mapreduce.Job: map 36% reduce 0%
2021-11-16 02:09:00.914 INFO mapreduce.Job: map 45% reduce 0%
2021-11-16 02:09:03.931 INFO mapreduce.Job: map 64% reduce 0%
2021-11-16 02:09:04.931 INFO mapreduce.Job: map 68% reduce 0%
2021-11-16 02:09:04.931 INFO mapreduce.Job: map 72% reduce 0%
2021-11-16 02:09:08.359 INFO mapreduce.Job: map 82% reduce 0%
2021-11-16 02:09:09.964 INFO mapreduce.Job: map 86% reduce 0%
2021-11-16 02:09:10.969 INFO mapreduce.Job: map 95% reduce 0%
2021-11-16 02:09:11.974 INFO mapreduce.Job: map 100% reduce 0%
2021-11-16 02:09:18.013 INFO mapreduce.Job: map 100% reduce 29%
2021-11-16 02:09:20.026 INFO mapreduce.Job: map 100% reduce 71%
2021-11-16 02:09:22.036 INFO mapreduce.Job: map 100% reduce 100%
2021-11-16 02:09:24.053 INFO mapreduce.Job: Job: job_1637015590244_0003 completed successfully
2021-11-16 02:09:24.141 INFO mapreduce.Job: Counters: 55
File System Counters
  FILE: Number of bytes read=192588
  FILE: Number of bytes written=7590655
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  FILE: Number of large write operations=0
  HDFS: Number of bytes read=1860914
  HDFS: Number of bytes written=9111
  HDFS: Number of read operations=101
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=21
  HDFS: Number of bytes read erasure-coded=0
Job Counters
  Killed map tasks=1
  Launched map tasks=22
  Launched reduce tasks=7
  Data-local map tasks=22

HDFS: Number of large read operations=0
HDFS: Number of write operations=21
HDFS: Number of bytes read erasure-coded=0
Job Counters
  Killed map tasks=1
  Launched map tasks=22
  Launched reduce tasks=7
  Data-local map tasks=22
Total time spent by all maps in occupied slots (ms)=409891812
Total time spent by all reduces in occupied slots (ms)=104779200
Total time spent by all map tasks (ms)=129877
Total time spent by all reduce tasks (ms)=3200
Total vcore-milliseconds taken by all map tasks=129877
Total vcore-milliseconds taken by all reduce tasks=3200
Total megabyte-milliseconds taken by all map tasks=409891812
Total megabyte-milliseconds taken by all reduce tasks=104779200
Map-Reduce Framework
  Map input records=13130
  Map output records=13129
  Map output bytes=166288
  Map output materialized bytes=193470
  Input split bytes=1826
  Combine input records=0
  Combine output records=0
  Reduce input groups=730
  Reduce shuffle bytes=193470
  Reduce input records=13129
  Reduce output records=730
  Spill Record bytes=2528
  Shuffled Maps =154
  Failed Shuffles=0
  Merged Map outputs=154
  GC time elapsed (ms)=4327
  CPU time spent (ms)=23080
  Physical memory (bytes) snapshot=12989292544
  Virtual memory (bytes) snapshot=128415109120
  Total committed heap usage (bytes)=12656836608
  Peak Map Physical memory (bytes)=553037824
  Peak Map Virtual memory (bytes)=4429299712
  Peak Reduce Physical memory (bytes)=285195328
  Peak Reduce Virtual memory (bytes)=4431630336
Shuffle Errors
  Bad Src Datanode ID=0
  CONNECTION=0
  IO ERROR=0
  WRONG LENGTH=0
  WRONG MAP=0
  WRONG REDUCE=0
File Input Format Counters
  Bytes Read=1859088
File Output Format Counters
  Bytes Written=9111
2021-11-16 02:09:24.141 INFO streaming.StreamJob: Output directory: /outputFolder5
mingruoqu123@cluster-508f:~/14848_temperature/temperature$ 

```

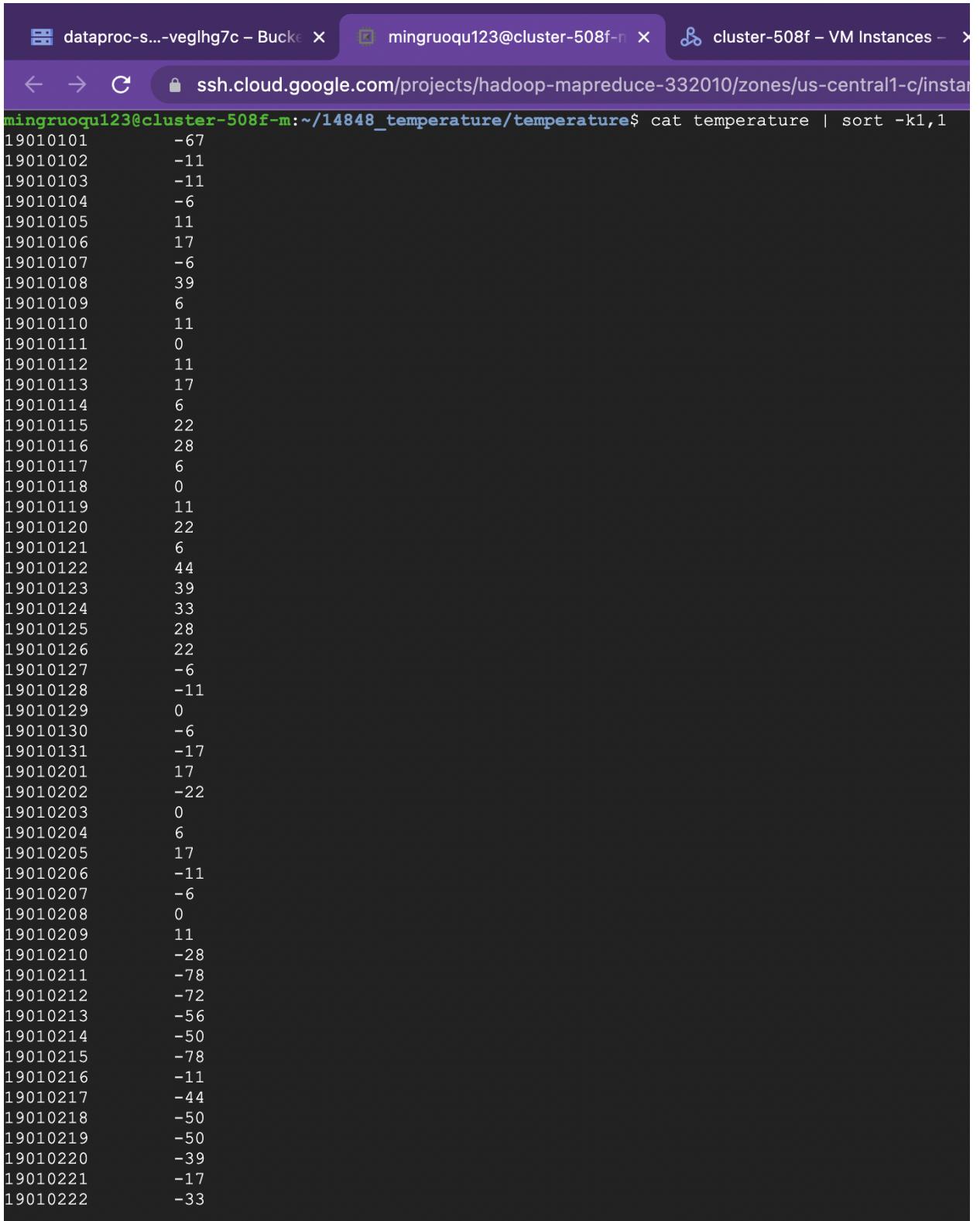
8. open the output directory, merge all the results to one file and store it on the local cluster.
You can see the results is stored as 'temperature'

```

Bytes Written=9111
2021-11-16 02:09:24,141 INFO streaming.StreamJob: Output directory: /outputFolder5
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /
Found 7 items
drwxr-xr-x  - mingruoqu123 hadoop          0 2021-11-16 01:46 /NewTemperatureOutputFolder
drwxr-xr-x  - mingruoqu123 hadoop          0 2021-11-16 01:40 /TemperatureOutputFolder
drwxr-xr-x  - mingruoqu123 hadoop          0 2021-11-16 00:47 /data
drwxr-xr-x  - mingruoqu123 hadoop          0 2021-11-16 02:09 /outputFolder5
drwxrwxrwt  - hdfs      hadoop          0 2021-11-15 22:33 /tmp
drwxrwxrwt  - hdfs      hadoop          0 2021-11-15 22:33 /user
drwx-wx-wx  - hive      hadoop          0 2021-11-15 22:33 /var
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /NewTemperatureOutputFolder
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -ls /outputFolder5
Found 8 items
-rw-r--r--  2 mingruoqu123 hadoop          0 2021-11-16 02:09 /outputFolder5/_SUCCESS
-rw-r--r--  2 mingruoqu123 hadoop        1297 2021-11-16 02:09 /outputFolder5/part-00000
-rw-r--r--  2 mingruoqu123 hadoop        1312 2021-11-16 02:09 /outputFolder5/part-00001
-rw-r--r--  2 mingruoqu123 hadoop        1299 2021-11-16 02:09 /outputFolder5/part-00002
-rw-r--r--  2 mingruoqu123 hadoop        1281 2021-11-16 02:09 /outputFolder5/part-00003
-rw-r--r--  2 mingruoqu123 hadoop        1324 2021-11-16 02:09 /outputFolder5/part-00004
-rw-r--r--  2 mingruoqu123 hadoop        1312 2021-11-16 02:09 /outputFolder5/part-00005
-rw-r--r--  2 mingruoqu123 hadoop        1286 2021-11-16 02:09 /outputFolder5/part-00006
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ hadoop fs -getmerge /outputFolder5 temperature
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ ls
data mapper.py reducer.py temperature
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ █

```

9. Results printed



A screenshot of a terminal window with three tabs at the top: 'dataproc-s...veglhg7c - Bucket' (purple), 'mingruoqu123@cluster-508f-m ~' (blue), and 'cluster-508f - VM Instances' (grey). The main pane shows a command-line session where the user is viewing a file named 'temperature'. The command 'cat temperature | sort -k1,1' has been run, displaying a sorted list of temperature values. The data consists of two columns: a date/time stamp and a temperature value. The first few lines of the output are:

```
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ cat temperature | sort -k1,1
19010101      -67
19010102      -11
19010103      -11
19010104      -6
19010105      11
19010106      17
19010107      -6
19010108      39
19010109      6
19010110      11
19010111      0
19010112      11
19010113      17
19010114      6
19010115      22
19010116      28
19010117      6
19010118      0
19010119      11
19010120      22
19010121      6
19010122      44
19010123      39
19010124      33
19010125      28
19010126      22
19010127      -6
19010128      -11
19010129      0
19010130      -6
19010131      -17
19010201      17
19010202      -22
19010203      0
19010204      6
19010205      17
19010206      -11
19010207      -6
19010208      0
19010209      11
19010210      -28
19010211      -78
19010212      -72
19010213      -56
19010214      -50
19010215      -78
19010216      -11
19010217      -44
19010218      -50
19010219      -50
19010220      -39
19010221      -17
19010222      -33
```

10. copy the results into your bucket. Then download the results through GUI.

```

19021231      -28
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ gsutil ls
gs://dataproc-staging-us-central1-147654780204-veglhg7c/
gs://dataproc-temp-us-central1-147654780204-fxdwdnbo/
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ gsutil cp temperature gs://dataproc-staging-us-central1-147654780204-veglhg7c/
Copying file://temperature [Content-Type=application/octet-stream]...
/ [1 files][ 8.9 KiB/ 8.9 KiB]
Operation completed over 1 objects/8.9 KiB.
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ gsutil cp temperature gs://dataproc-staging-us-central1-147654780204-veglhg7c/temperature/
Copying file://temperature [Content-Type=application/octet-stream]...
/ [1 files][ 8.9 KiB/ 8.9 KiB]
Operation completed over 1 objects/8.9 KiB.
mingruoqu123@cluster-508f-m:~/14848_temperature/temperature$ 

```

The screenshot shows the Google Cloud Platform Storage Browser interface. The left sidebar has 'Cloud Storage' selected. The main area shows 'Bucket details' for 'dataproc-staging-us-central1-147654780204-veglhg7c'. The 'OBJECTS' tab is active, displaying a list of objects:

Name	Type	Created	Storage class	Last modified
.DS_Store	application/octet-stream	Nov 15, 2018	Standard	Nov 15, 2018
data/	Folder	—	—	—
mapper.py	text/x-python-script	Nov 15, 2018	Standard	Nov 15, 2018
reducer.py	text/x-python-script	Nov 15, 2018	Standard	Nov 15, 2018
temperature	application/octet-stream	Nov 15, 2018	Standard	Nov 15, 2018

A red circle highlights the 'temperature' object in the list. The right sidebar includes sections for 'Recommended for you', 'Control access to data', 'Make data public', 'Manage object lifecycles', 'You might also like', and links to 'API & references', 'Access control', and 'Resources'.