

Girish Rajani  
A20503736  
CSP-554 Big Data Technologies  
Homework 3

```

giris@LAPTOP-18TD7ESB MINGW64 ~
$ ssh -i d:/users/giris/downloads/emr-key-pair.pem hadoop@ec2-54-237-51-201.compute-1.amazonaws.com
Last login: Wed Sep 21 17:03:52 2022

 _ | _ | _ )
 _ | ( _ /   Amazon Linux 2 AMI
 _ | \ _ | _ |

https://aws.amazon.com/amazon-linux-2/

EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::::M          M:::::::::M R:::::::::R
EE:::::EEEEEEEE::::E M:::::::::M          M:::::::::M R:::::RRRRRR:::::R
  E::::E          EEEEE M:::::::::M          M:::::::::M RR::::R          R::::R
  E::::E          M:::::::::M M:::::::::M M:::::::::M R::::R          R::::R
  E:::::EEEEEEEEEEEE M:::::M M:::::M M:::::M M:::::M R::::RRRRRR:::::R
  E:::::EEEEEEEEEEEE M:::::M M:::::M M:::::M M:::::M R:::::RRRRRR:::::R
  E:::::EEEEEEEEEEEE M:::::M M:::::M M:::::M M:::::M R::::RRRRRR:::::R
  E::::E          M:::::::::M M:::::M M:::::M M:::::M R::::R          R::::R
  E::::E          EEEEE M:::::M          MMM M:::::M M:::::M R::::R          R::::R
EE:::::EEEEEEEE::::E M:::::M          M:::::M M:::::M R::::R          R::::R
E:::::EEEEEEEE::::E M:::::M          M:::::M RR::::R          R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRR          RRRRRR

[hadoop@ip-172-31-61-97 ~]$ sudo /usr/bin/pip3.7 install mrjob[aws]

```

Figure 1 - Showing successful connection to Master Node using SSH and installing mrjob

```
Successfully installed boto3-1.24.77 botocore-1.27.77 mrjob-0.7.4 python-dateutil-2.8.2 s3t
ransfer-0.6.0 urllib3-1.26.12
[hadoop@ip-172-31-61-97 ~]$
```

Figure 2 - Showing successful installation of mrjob

[illegible]

Figure 3 - Moved both WordCount.py and w.data files to the /home/hadoop directory using scp

```
[hadoop@ip-172-31-61-97 ~]$ hadoop fs -put /home/hadoop/w.data /user/hadoop
```

Figure 4 - Using -put command to copy w.data file from the /home/hadoop directory to HDFS /user/hadoop directory

```

36/output...
"a"      3
"all"    1
"an"     1
"and"    1
"are"    1
"as"     4
"available" 1
"be"     3
"by"     1
"cluster" 2
"combine" 1
"contained" 1
"defined" 1
"dependencies" 1
"do"     1
"either" 1
"executed" 1
"explains" 1
"file"   2
"first"  1
"following" 1
"for"    1
"hadoop" 1
"how"    2
"in"     1
"individual" 1
"is"     2
"job"    4
"machine" 1
"map"    1
"more"   2
"mrjob"  1

"must" 1
"nodes" 1
"of" 1
"on" 4
"or" 2
"oriented" 1
"our" 1
"program" 1
"python" 1
"reduce" 1
"reference" 1
"run" 1
"runners" 1
"script" 1
"second" 1
"sections" 1
"see" 1
"submitted" 1
"task" 2
"that" 1
"the" 4
"things" 1
"those" 1
"to" 3
"two" 1
"uploaded" 1
"versions" 1
"well" 1
"when" 1
"will" 1
"within" 1
"writing" 2
"your" 5
Removing HDFS temp direct

```

Figure 5 - Showing output after reading the WordCount.py file using the following command:  
python WordCount.py -r hadoop hdfs:///user/hadoop/w.data

7) (5 points) Submit a copy of this modified program and a screen shot of the results of the program's execution as the output of your assignment.

```

job output is in hdfs:///user/hadoop/tmp/mrjob/WordCount2.hadoop.20220921.201844.289504/output
Streaming final output from hdfs:///user/hadoop/tmp/mrjob/WordCount2.hadoop.20220921.201844.289504/output...
"a_to_n"      46
"other"      49
Removing HDFS temp directory hdfs:///user/hadoop/tmp/mrjob/WordCount2.hadoop.20220921.201844.289504...
Removing temp directory /tmp/WordCount2.hadoop.20220921.201844.289504...
[hadoop@ip-172-31-61-97 ~]$

```

Figure 6 - Showing screenshot of the results after reading the WordCount2.py file using the following command:  
python WordCount2.py -r hadoop hdfs:///user/hadoop/w.data

The program counted 46 words beginning with small letters a to n and 49 words beginning with anything else.

```
[hadoop@ip-172-31-61-97 ~]$ python Salaries.py -r hadoop hdfs:///user/hadoop/Salaries.tsv
```

Figure 7 - Uploaded Salaries.py and Salaries.tsv to HDFS using scp and then read the python file using the following command:

```
python Salaries.py -r hadoop hdfs:///user/hadoop/Salaries.tsv
```

```
"Waste Water Opns Tech II Sanit"      81
"Waste Water Tech Supv I Pump"      6
"Waste Water Tech Supv II Pump"      1
"Waste Water Tech Supv II Sanit"     10
"Waste Water Techn Supv I Sanit"     19
"Water Systems Pumping Supv"         1
"Water Systems Treatment Manage"     1
"Water Systems Treatment Supv"       2
"YOUTH DEVELOPMENT TECH"             3
"ZONING ADMINISTRATOR"               1
"ZONING APPEALS ADVISOR BMZA"        1
"ZONING APPEALS OFFICER"             1
"ZONING ENFORCEMENT OFFICER"        1
"ZONING EXAMINER I"                  2
"ZONING EXAMINER II"                 1
Removing HDFS temp directory hdfs:///user/hadoop/tmp/mrjob/Salaries.hadoop.20220921.211610.724953...
Removing temp directory /tmp/Salaries.hadoop.20220921.211610.724953...
[hadoop@ip-172-31-61-97 ~]$
```

Figure 8 - Illustrating the output, which shows the number of workers who share the same job

11) (5 points) Submit a copy of this modified program and a screen shot of the results of the program's execution as the output of your assignment.

```
job output is in hdfs:///user/hadoop/tmp/mrjob/Salaries2.hadoop.20220921.222046.731251/output
Streaming final output from hdfs:///user/hadoop/tmp/mrjob/Salaries2.hadoop.20220921.222046.731251/output...
"High" 442
"Low" 7064
"Medium" 6312
Removing HDFS temp directory hdfs:///user/hadoop/tmp/mrjob/Salaries2.hadoop.20220921.222046.731251...
Removing temp directory /tmp/Salaries2.hadoop.20220921.222046.731251...
[hadoop@ip-172-31-61-97 ~]$ client_loop: send disconnect: Connection reset by peer
```

Figure 9 - Showing screenshot of the results after reading the Salaries2.py file using the following command:

```
python Salaries2.py -r hadoop hdfs:///user/hadoop/Salaries.tsv
```

The program counted 442 workers who have 'High' annual salaries, 7064 workers who have 'Low' annual salaries, and 6312 employees who have 'Medium' annual salaries.

**Note: From this question onwards, the cluster had been terminated and then cloned at a later date. That is why the Hadoop IP address will be different in the new screenshots below**

```
[hadoop@ip-172-31-63-51 ~]$ hadoop fs -put /home/hadoop/u.data /user/hadoop
```

Figure 10 - Using -put command to copy u.data file from the /home/hadoop directory to HDFS /user/hadoop directory

13) (5 points) Write a program to perform the task of outputting a count of the number of movies each user (identified via their user id) reviewed.

```
"80" 37
"81" 160
"82" 39
"83" 161
"84" 116
"85" 107
"86" 190
"87" 31
"88" 255
"89" 66
"9" 45
"90" 50
"91" 150
"92" 123
"93" 159
"94" 196
"95" 299
"96" 76
"97" 128
"98" 71
"99" 188
Removing HDFS temp directory hdfs:///user/hadoop/tmp/mrjob/Movies.hadoop.20220922.015826.509108...
Removing temp directory /tmp/Movies.hadoop.20220922.015826.509108...
[hadoop@ip-172-31-63-51 ~]$ 20
-bash: 20: command not found
[hadoop@ip-172-31-63-51 ~]$
```

Figure 11 - Showing screenshot of the results after reading the Movies.py file using the following command:

```
python Movies.py -r hadoop hdfs:///user/hadoop/u.data
```

The program shows the number of movies each user (identified via their user id) reviewed. So for example, from figure 11 above, we can see that user 9 reviewed 45 movies.

In total, this movie dataset contains over 10,000 review samples, and to be able to count the number of movies each user reviewed (identified via their user id) within a few lines of code is astounding.