

Name: Nachiketh Reddy

ID: 2117731

I have done this exam entirely on my own. I have not consulted with friends or consulted online resources that have the solution to the exam.

QUESTION 1:

Clustering Write Python code within the four files provided. Each file has instructions to write code as described in the Python triple quotes, i.e., `"""TODO:..."""`. You must write code using these files and not write your own files. To write code you will need to be familiar with numpy. The tutorial at this site <https://numpy.org/doc/stable/user/quickstart.html> will be sufficient. Once you have generated data, generated the centers, and written the code in all four files, test it with: `cat kmeans data.csv | python3 mapper-kmeans.py | sort -n | python3 reducer-kmeans.py > centers.txt` Now repeat the process and run 4 iterations of kMeans. Remember to change the centers file in mapper-kmeans.py everytime to correspond to the centers generated in the previous iteration. Note, you can maintain centers as `centers[i].txt` where `[i]` is the centers output by the `i`th iteration. The first centers file must be generated apriori and is used by mapper-kmeans.py Submit the following:

- The initial centers.txt and the kmeans data.csv. Please zip the kmeans data.csv.
- All four Python files.
- The final centers.txt after 4 iterations.
- The Hadoop command used and a screenshot of running on Hadoop cluster with time information for each iteration.

Initial centers.txt:

```
$ cat centers.txt
0.418506478974507,0.08759190798437044,0.6202098841228467,0.3340070756262328,0.7613721028137473
0.798336886574648,0.47360098667743256,0.12065307097839384,0.5759179675004569,0.4655673944438966
0.18810268101729843,0.9971773467503174,0.5126709983572116,0.9026866448432411,0.3620564522378126
0.3562429889627011,0.039461805404032546,0.10208725537669039,0.9866256772840983,0.45061509541729616
0.13902698083927612,0.9338942588983641,0.022573079076542046,0.36246402492945984,0.07085101418126205
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 16:56:07
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

First Run:

`time hadoop jar hadoop-streaming-2.6.4.jar -input /data/kmeans_data.csv -output /data/out -mapper mapper-kmeans.py -reducer reducer-kmeans.py -file mapper-kmeans.py -file reducer-kmeans.py -file centers.txt`

```

AWS Services Search [Alt+S] N. Virginia voclabs/user3013569@nparamah@depaul.edu @ 2854-

[ec2-user@ip-172-31-62-135 container_1711125434068_0003_01_000005] 2024-03-22 16:48:26
$ cd ~
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 16:49:01
$ time hadoop jar hadoop-streaming-2.6.4.jar -input /data/kmeans_data.csv -output /data/out -mapper mapper-kmeans.py -reducer reducer-kmeans.py -file mapper-kmeans.py -file reducer-kmeans.py -file centers.txt
24/03/22 16:49:28 WARN streaming.StreamJob: File option is deprecated, Please use generic option -files instead.
packageJobJar: [mapper-kmeans.py, reducer-kmeans.py, centers.txt, /tmp/hadoop-unjar228405276364153823/] [] /tmp/streamjob1207521121587065422.jar tmpDir=null
24/03/22 16:49:28 INFO client.RMProxy: Connecting to ResourceManager at /172.31.62.135:8032
24/03/22 16:49:28 INFO client.RMProxy: Connecting to ResourceManager at /172.31.62.135:8032
24/03/22 16:49:29 INFO mapred.FileInputFormat: Total input paths to process : 1
24/03/22 16:49:29 INFO net.NetworkTopology: Adding a new node: /default-rack/172.31.62.135:50010
24/03/22 16:49:29 INFO mapreduce.JobSubmitter: number of splits:2
24/03/22 16:49:29 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1711125434068_0004
24/03/22 16:49:29 INFO impl.YarnClientImpl: Submitted application application_1711125434068_0004
24/03/22 16:49:29 INFO mapreduce.Job: The url to track the job: http://ip-172-31-62-135.ec2.internal:8088/proxy/application_1711125434068_0004/
24/03/22 16:49:29 INFO mapreduce.Job: Running job: job_1711125434068_0004
24/03/22 16:49:34 INFO mapreduce.Job: Job job_1711125434068_0004 running in uber mode : false
24/03/22 16:49:34 INFO mapreduce.Job: map 0% reduce 0%
24/03/22 16:49:44 INFO mapreduce.Job: map 6% reduce 0%
24/03/22 16:49:47 INFO mapreduce.Job: map 9% reduce 0%
24/03/22 16:49:50 INFO mapreduce.Job: map 13% reduce 0%
24/03/22 16:49:53 INFO mapreduce.Job: map 16% reduce 0%
24/03/22 16:49:56 INFO mapreduce.Job: map 19% reduce 0%
24/03/22 16:49:59 INFO mapreduce.Job: map 23% reduce 0%
24/03/22 16:50:02 INFO mapreduce.Job: map 26% reduce 0%
24/03/22 16:50:05 INFO mapreduce.Job: map 29% reduce 0%
24/03/22 16:50:08 INFO mapreduce.Job: map 32% reduce 0%
24/03/22 16:50:10 INFO mapreduce.Job: map 50% reduce 0%
24/03/22 16:50:21 INFO mapreduce.Job: map 56% reduce 0%
24/03/22 16:50:24 INFO mapreduce.Job: map 59% reduce 0%
24/03/22 16:50:27 INFO mapreduce.Job: map 62% reduce 0%
24/03/22 16:50:30 INFO mapreduce.Job: map 65% reduce 0%
24/03/22 16:50:33 INFO mapreduce.Job: map 69% reduce 0%
24/03/22 16:50:36 INFO mapreduce.Job: map 72% reduce 0%
24/03/22 16:50:39 INFO mapreduce.Job: map 75% reduce 0%
24/03/22 16:50:42 INFO mapreduce.Job: map 78% reduce 0%
24/03/22 16:50:45 INFO mapreduce.Job: map 81% reduce 0%
24/03/22 16:50:48 INFO mapreduce.Job: map 100% reduce 0%
24/03/22 16:50:57 INFO mapreduce.Job: map 100% reduce 81%
24/03/22 16:51:00 INFO mapreduce.Job: map 100% reduce 89%
24/03/22 16:51:03 INFO mapreduce.Job: map 100% reduce 97%
24/03/22 16:51:05 INFO mapreduce.Job: map 100% reduce 100%
24/03/22 16:51:05 INFO mapreduce.Job: Job job_1711125434068_0004 completed successfully
24/03/22 16:51:05 INFO mapreduce.Job: Counters: 50

```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

```
aws Services Search [Alt+S]

Total time spent by all map tasks (ms)=70312
Total time spent by all reduce tasks (ms)=15344
Total vcore-milliseconds taken by all map tasks=70312
Total vcore-milliseconds taken by all reduce tasks=15344
Total megabyte-milliseconds taken by all map tasks=421872000
Total megabyte-milliseconds taken by all reduce tasks=92064000
Map-Reduce Framework
  Map input records=2000000
  Map output records=2000000
  Map output bytes=200697359
  Map output materialized bytes=204697371
  Input split bytes=186
  Combine input records=0
  Combine output records=0
  Reduce input groups=2000000
  Reduce shuffle bytes=204697371
  Reduce input records=2000000
  Reduce output records=5
  Spilled Records=6000000
  Shuffled Maps =2
  Failed Shuffles=0
  Merged Map outputs=2
  GC time elapsed (ms)=199
  CPU time spent (ms)=86850
  Physical memory (bytes) snapshot=857776128
  Virtual memory (bytes) snapshot=9967841280
  Total committed heap usage (bytes)=838336512
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=192701455
File Output Format Counters
  Bytes Written=480
24/03/22 16:51:05 INFO streaming.StreamJob: Output directory: /data/out

real    1m38.293s
user    0m3.957s
sys     0m0.246s
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

centers1.txt

```
$ hadoop fs -cat /data/out/part-00000 | more
0.4466750888751059,0.3252244435386771,0.6633358814501816,0.3699391537086328,0.6190676306814223
0.7520720779788497,0.5478180588334175,0.3475748212662098,0.503967543531433,0.4704351345767741
0.3426654567922555,0.7763893876610183,0.6377502557562773,0.7037533394701059,0.45552921303595967
0.32200687462807565,0.2042678290089987,0.3010958777611171,0.809901256192348,0.427933120977061
0.25406158464323236,0.7309643097996597,0.280721980197176,0.27361726142219434,0.2685530185949394
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:01:06
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

Second Run:

```
time hadoop jar hadoop-streaming-2.6.4.jar -input /data/kmeans_data.csv -output /data/out_1 -mapper mapper-kmeans.py -reducer reducer-kmeans.py -file mapper-kmeans.py -file reducer-kmeans.py -file centers1.txt
```

```
aws Services Search [Alt+S]
24/03/22 17:06:28 INFO mapreduce.Job: Counters: 50
  File System Counters
    FILE: Number of bytes read=409394754
    FILE: Number of bytes written=614424194
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=192701641
    HDFS: Number of bytes written=484
    HDFS: Number of read operations=9
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=455628
    Total time spent by all reduces in occupied slots (ms)=91866
    Total time spent by all map tasks (ms)=75938
    Total time spent by all reduce tasks (ms)=15311
    Total vcore-milliseconds taken by all map tasks=75938
    Total vcore-milliseconds taken by all reduce tasks=15311
    Total megabyte-milliseconds taken by all map tasks=455628000
    Total megabyte-milliseconds taken by all reduce tasks=91866000
  Map-Reduce Framework
    Map input records=2000000
    Map output records=2000000
    Map output bytes=200697359
    Map output materialized bytes=204697371
    Input split bytes=186
    Combine input records=0
    Combine output records=0
    Reduce input groups=2000000
    Reduce shuffle bytes=204697371
    Reduce input records=2000000
    Reduce output records=5
    Spilled Records=6000000
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=227
    CPU time spent (ms)=92940
    Physical memory (bytes) snapshot=839688192
    Virtual memory (bytes) snapshot=9370372608
    Total committed heap usage (bytes)=789577728
  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=192701455
  File Output Format Counters
    Bytes Written=484
24/03/22 17:06:28 INFO streaming.StreamJob: Output directory: /data/out_1
real    1m44.789s
user    0m3.902s
sys     0m0.366s
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:06:28
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

centers2.txt:

```
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:06:28
$ hadoop fs -cat /data/out_1/part-00000 | more
0.4734409778155092,0.31433915081627994,0.6986501515506048,0.33940488359268,0.6265724075484513
0.8031124396264666,0.5546227727457731,0.3595995575618838,0.49196922450537167,0.4901047027459025
0.3831126043898117,0.761100105539219,0.6763279032276892,0.7093220751550083,0.4929102468820437
0.34878141350228486,0.23329973429683096,0.3339663188372609,0.7623973942110422,0.4473854823347234
0.28189829426449003,0.6805021345177569,0.31218729709060067,0.25860423145094574,0.3326422763007367
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:07:29
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

Third Run:

```
time hadoop jar hadoop-streaming-2.6.4.jar -input /data/kmeans_data.csv -output /data/out_2 -mapper mapper-kmeans.py -reducer reducer-kmeans.py -file mapper-kmeans.py -file reducer-kmeans.py -file centers2.txt
```

```
aws Services Search [Alt+S]
24/03/22 17:11:02 INFO mapreduce.Job: Counters: 50
  File System Counters
    FILE: Number of bytes read=409394754
    FILE: Number of bytes written=614424194
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=192701641
    HDFS: Number of bytes written=486
    HDFS: Number of read operations=9
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=434508
    Total time spent by all reduces in occupied slots (ms)=91596
    Total time spent by all map tasks (ms)=72418
    Total time spent by all reduce tasks (ms)=15266
    Total vcore-milliseconds taken by all map tasks=72418
    Total vcore-milliseconds taken by all reduce tasks=15266
    Total megabyte-milliseconds taken by all map tasks=434508000
    Total megabyte-milliseconds taken by all reduce tasks=91596000
  Map-Reduce Framework
    Map input records=2000000
    Map output records=2000000
    Map output bytes=200697359
    Map output materialized bytes=204697371
    Input split bytes=186
    Combine input records=0
    Combine output records=0
    Reduce input groups=2000000
    Reduce shuffle bytes=204697371
    Reduce input records=2000000
    Reduce output records=5
    Spilled Records=6000000
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=211
    CPU time spent (ms)=90180
    Physical memory (bytes) snapshot=838172672
    Virtual memory (bytes) snapshot=9972330496
    Total committed heap usage (bytes)=804782080
  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=192701455
  File Output Format Counters
    Bytes Written=486
24/03/22 17:11:02 INFO streaming.StreamJob: Output directory: /data/out_2

real    1m40.371s
user    0m4.009s
sys     0m0.349s
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:11:02
$

i-062d65a3dd2e0d019 (nachiketh_server)
PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135
```

Centers3.txt:

```
$ hadoop fs -cat /data/out_2/part-00000 | more
0.4937146240474645,0.3036055110369187,0.7221833498730471,0.32118091322708636,0.6218470497306814
0.8201501479677525,0.5599202231741127,0.3580808876313664,0.48630996221515654,0.49609954018696495
0.40760238775264707,0.7521683493085846,0.6931897494272781,0.7136434667723321,0.5082152874760778
0.36340281394218504,0.24447665867134197,0.3484375930606069,0.746892746410468,0.457046003110177
0.29121142049482523,0.6605776564436759,0.3258643024647322,0.2568438608784413,0.36767398303408844
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:12:47
$

i-062d65a3dd2e0d019 (nachiketh_server)
PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135
```

Fourth run:

```
time hadoop jar hadoop-streaming-2.6.4.jar -input /data/kmeans_data.csv -output /data/out_3 -mapper mapper-kmeans.py -reducer reducer-kmeans.py -file mapper-kmeans.py -file reducer-kmeans.py -file centers3.txt
```

```
aws Services [Alt+S]
24/03/22 17:15:49 INFO mapreduce.Job: Counters: 50
  File System Counters
    FILE: Number of bytes read=409394754
    FILE: Number of bytes written=614424194
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=192701641
    HDFS: Number of bytes written=485
    HDFS: Number of read operations=9
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=453996
    Total time spent by all reduces in occupied slots (ms)=93816
    Total time spent by all map tasks (ms)=75666
    Total time spent by all reduce tasks (ms)=15636
    Total vcore-milliseconds taken by all map tasks=75666
    Total vcore-milliseconds taken by all reduce tasks=15636
    Total megabyte-milliseconds taken by all map tasks=453996000
    Total megabyte-milliseconds taken by all reduce tasks=93816000
  Map-Reduce Framework
    Map input records=2000000
    Map output records=2000000
    Map output bytes=200697359
    Map output materialized bytes=204697371
    Input split bytes=186
    Combine input records=0
    Combine output records=0
    Reduce input groups=2000000
    Reduce shuffle bytes=204697371
    Reduce input records=2000000
    Reduce output records=5
    Spilled Records=6000000
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=213
    CPU time spent (ms)=93230
    Physical memory (bytes) snapshot=837906432
    Virtual memory (bytes) snapshot=9970728960
    Total committed heap usage (bytes)=760217600
  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=192701455
  File Output Format Counters
    Bytes Written=485
24/03/22 17:15:49 INFO streaming.StreamJob: Output directory: /data/out_3

real    1m45.476s
user    0m4.127s
sys     0m0.303s
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:15:49
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

Centers4.txt:

```
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:15:49
$ hadoop fs -cat /data/out_3/part-000000 | more
0.5103926053310995,0.2977545901367192,0.7366440296430301,0.31258292999415127,0.6101416282578812
0.8242878188037962,0.5655850345644502,0.3483211768832058,0.48126983360050707,0.497323176535033
0.42362716273242834,0.7463217670226793,0.7030348101671381,0.715393510079285,0.5143133002628412
0.37082399732522636,0.24839848545648657,0.35587840538792614,0.7425507539694319,0.4630451944577867
0.29027488652129885,0.6523560172490893,0.331962013239862,0.2576935771924344,0.39065532819201376
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 17:16:38
$
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

QUESTION 2:

BloomFilter-based 2-way map-side Join The objective of this question is to implement a 1-pass MR that uses bloom filter for the larger table. Write Python code within the three files provided. Each file has instructions to write code as described in the Python triple quotes, i.e., `“““TODO:...”””`. You must write code using these files and not write your own files. Consider the following query:

```
select * from lineorder, dwdate where lo_orderdate = d_datekey and d_sellingseason = 'Fall'
```

Lineorder and Dwdate are to be downloaded from <http://cdmgcsarprd01.dpu.depaul.edu/CSC555/SSBM1/dwdate.tbl>
<http://cdmgcsarprd01.dpu.depaul.edu/CSC555/SSBM1/lineorder.tbl>


In the map phase, the bloom filter of the large table (lineorder) will be read. The small table (dwdate) will also be read. In the map phase, make both the where clause checks and output the join result. The reduce phase just passes the result. The bloomfilter should be setup as described in bloomfilter.py.

Python bloomfilter.py: bits set

```
aws Services [Alt+S] N. Virginia vocabs/user3013569-nparamah@depaul.edu @ 2854-6811-230

6633, 236638, 236642, 236644, 236716, 236723, 236731, 236735, 236787, 236797, 236799, 236830, 236904, 236906, 236945, 236961, 236974, 236978, 236990, 237009, 237022, 237029, 237081, 237083, 237089, 237117, 237136, 237138, 237148, 237185, 237186, 237209, 237211, 237224, 237232, 237271, 237314, 237350, 237376, 237404, 237480, 237519, 237547, 237551, 237604, 237630, 237635, 237657, 237660, 237703, 237716, 237721, 237772, 237774, 237786, 237790, 237795, 237815, 237823, 237832, 237839, 237893, 237943, 237956, 237985, 238005, 238010, 238014, 238019, 238027, 238030, 238043, 238046, 238051, 238069, 238163, 238173, 238193, 238209, 238216, 238217, 238221, 238256, 238301, 238302, 238305, 238306, 238342, 238351, 238385, 238394, 238403, 238407, 238421, 238422, 238428, 238430, 238432, 238443, 238451, 238467, 238471, 238472, 238504, 238518, 238549, 238609, 238614, 238624, 238668, 238738, 238746, 238765, 238768, 238858, 238917, 238920, 238938, 238946, 238951, 238963, 239016, 239025, 239031, 239042, 239060, 239109, 239129, 239136, 239153, 239164, 239178, 239179, 239203, 239240, 239257, 239296, 239309, 239330, 239339, 239414, 239426, 239432, 239439, 239455, 239470, 239485, 239488, 239492, 239494, 239496, 239503, 239526, 239533, 239567, 239582, 239585, 239586, 239622, 239624, 239638, 239641, 239654, 239658, 239667, 239731, 239732, 239733, 239747, 239779, 239798, 239809, 239832, 239835, 239893, 239906, 239940, 239982, 240002, 240019, 240044, 240116, 240176, 240178, 240220, 240243, 240246, 240252, 240269, 240278, 240329, 240342, 240359, 240390, 240394, 240397, 240432, 240434, 240452, 240476, 240487, 240511, 240534, 240535, 240537, 240580, 240593, 240598, 240617, 240619, 240633, 240639, 240649, 240650, 240652, 240653, 240687, 240695, 240712, 240724, 240736, 240767, 240817, 240850, 240872, 240879, 240882, 240890, 240923, 240930, 240935, 240977, 241079, 241082, 241084, 241094, 241098, 241103, 241105, 241119, 241154, 241155, 241167, 241176, 241178, 241184, 241194, 241216, 241250, 241254, 241279, 241282, 241292, 241342, 241344, 241388, 241412, 241417, 241427, 241441, 241449, 241491, 241514, 241536, 241561, 241585, 241632, 241643, 241672, 241687, 241700, 241728, 241729, 241756, 241789, 241812, 241814, 241820, 241836, 241840, 241848, 241854, 241873, 241916, 241919, 241927, 241955, 241962, 241965, 241971, 241995, 242002, 242022, 242065, 242070, 242084, 242087, 242112, 242144, 242146, 242203, 242211, 242230, 242235, 242259, 242306, 242318, 242320, 242325, 242332, 242362, 242381, 242398, 242400, 242407, 242416, 242436, 242452, 242497, 242568, 242610, 242631, 242638, 242655, 242672, 242676, 242686, 242696, 242708, 242727, 242731, 242753, 242764, 242765, 242800, 242818, 242865, 242867, 242897, 242922, 242942, 242950, 242953, 243014, 243054, 243059, 243060, 243063, 243091, 243085, 243090, 243094, 243095, 243096, 243115, 243120, 243125, 243132, 243149, 243179, 243196, 243213, 243238, 243246, 243269, 243309, 243313, 243330, 243344, 243445, 243447, 243479, 243499, 243515, 243521, 243541, 243542, 243559, 243562, 243563, 243599, 243616, 243642, 243649, 243652, 243661, 243667, 243668, 243678, 243684, 243704, 243718, 243747, 243777, 243801, 243809, 243818, 243826, 243840, 243842, 243861, 243866, 243878, 243881, 243897, 243902, 243918, 243950, 243970, 244014, 244029, 244044, 244059, 244084, 244093, 244112, 244127, 244133, 244137, 244145, 244157, 244182, 244201, 244240, 244319, 244320, 244336, 244337, 244350, 244357, 244374, 244375, 244385, 244390, 244394, 244403, 244426, 244452, 244469, 244492, 244499, 244502, 244507, 244515, 244548, 244552, 244563, 244597, 244598, 244607, 244624, 244627, 244633, 244640, 244642, 244654, 244702, 244714, 244724, 244737, 244738, 244797, 244798, 244827, 244846, 244873, 244891, 244977, 244980, 245023, 245041, 245050, 245052, 245064, 245095, 245106, 245107, 245125, 245173, 245175, 245176, 245177, 245182, 245195, 245248, 245261, 245283, 245312, 245316, 245317, 245325, 245333, 245337, 245354, 245355, 245356, 245365, 245396, 245407, 245408, 245418, 245432, 245442, 245446, 245459, 245503, 245542, 245558, 245606, 245633, 245635, 245671, 245686, 245689, 245731, 245740, 245748, 245769, 245779, 245798, 245800, 245840, 245853, 245887, 245891, 245898, 245899, 245902, 245949, 245952, 245956, 245990, 246008, 246010, 246019, 246037, 246076, 246133, 246164, 246167, 246183, 246191, 246214, 246283, 246289, 246296, 246313, 246315, 246342, 246372, 246390, 246412, 246438, 246440, 246454, 246459, 246470, 246490, 246497, 246521, 246566, 246606, 246611, 246626, 246631, 246654, 246683, 246708, 246715, 246742, 246756, 246781, 246791, 246792, 246821, 246834, 246855, 246866, 246876, 246880, 246888, 246890, 246904, 246920, 246962, 246964, 246971, 246972, 246985, 246996, 247001, 247007, 247053, 247071, 247089, 247094, 247118, 247133, 247214, 247221, 247227, 247229, 247233, 247248, 247258, 247262, 247269, 247277, 247289, 247293, 247298, 247342, 247369, 247370, 247394, 247399, 247404, 247415, 247433, 247461, 247497, 247546, 247571, 247585, 247621, 247647, 247654, 247670, 247673, 247729, 247776, 247813, 247835, 247836, 247864, 247875, 247887, 247891, 247892, 247915, 247937, 247975, 248006, 248011, 248052, 248072, 248076, 248116, 248122, 248123, 248132, 248136, 248147, 248150, 248160, 248192, 248210, 248219, 248278, 248281, 248286, 248331, 248380, 248382, 248385, 248389, 248395, 248404, 248406, 248418, 248457, 248467, 248496, 248496, 248500, 248509, 248522, 248527, 248557, 248565, 248627, 248628, 248650, 248659, 248673, 248693, 248720, 248724, 248752, 248786, 248802, 248827, 248890, 248911, 248914, 248924, 248928, 248934, 248988, 248991, 249024, 249029, 249030, 249036, 249068, 249087, 249089, 249110, 249122, 249148, 249149, 249151, 249176, 249195, 249257, 249260, 249269, 249289, 249311, 249336, 249346, 249352, 249357, 249365, 249406, 249415, 249432, 249433, 249444, 249468, 249480, 249504, 249543, 249581, 249591, 249606, 249612, 249620, 249631, 249679, 249687, 249701, 249719, 249763, 249800, 249815, 249847, 249864, 249872, 249885, 249899, 249905, 250034, 250042, 250060, 250076, 250147, 250150, 250151, 250164, 250174, 250177, 250181, 250192, 250269, 250272, 250283, 250289, 250310, 250313, 250329, 250354, 250422, 250439, 250472, 250484, 250524, 250539, 250547, 250572, 250573, 250578, 250600, 250601, 250630, 251063, 251066, 251069, 251073, 251075, 251076, 251079, 251081, 251082, 251083, 251084, 251085, 251086, 251087, 251088, 251089, 251090, 251091, 251092, 251093, 251094, 251095, 251096, 251097, 251098, 251099, 251100, 251101, 251102, 251103, 251104, 251105, 251106, 251107, 251108, 251109, 251110, 251111, 251112, 251113, 251114, 251115, 251116, 251117, 251118, 251119, 251120, 251121, 251122, 251123, 251124, 251125, 251126, 251127, 251128, 251129, 251130, 251131, 251132, 251133, 251134, 251135, 251136, 251137, 251138, 251139, 251140, 251141, 251142, 251143, 251144, 251145, 251146, 251147, 251148, 251149, 251150, 251151, 251152, 251153, 251154, 251155, 251156, 251157, 251158, 251159, 251160, 251161, 251162, 251163, 251164, 251165, 251166, 251167, 251168, 251169, 251170, 251171, 251172, 251173, 251174, 251175, 251176, 251177, 251178, 251179, 251180, 251181, 251182, 251183, 251184, 251185, 251186, 251187, 251188, 251189, 251190, 251191, 251192, 251193, 251194, 251195, 251196, 251197, 251198, 251199, 251200, 251201, 251202, 251203, 251204, 251205, 251206, 251207, 251208, 251209, 251210, 251211, 251212, 251213, 251214, 251215, 251216, 251217, 251218, 251219, 251220, 251221, 251222, 251223, 251224, 251225, 251226, 251227, 251228, 251229, 251230, 251231, 251232, 251233, 251234, 251235, 251236, 251237, 251238, 251239, 251240, 251241, 251242, 251243, 251244, 251245, 251246, 251247, 251248, 251249, 251250, 251251, 251252, 251253, 251254, 251255, 251256, 251257, 251258, 251259, 251260, 251261, 251262, 251263, 251264, 251265, 251266, 251267, 251268, 251269, 251270, 251271, 251272, 251273, 251274, 251275, 251276, 251277, 251278, 251279, 251280, 251281, 251282, 251283, 251284, 251285, 251286, 251287, 251288, 251289, 251290, 251291, 251292, 251293, 251294, 251295, 251296, 251297, 251298, 251299, 251300, 251301, 251302, 251303, 251304, 251305, 251306, 251307, 251308, 251309, 251310, 251311, 251312, 251313, 251314, 251315, 251316, 251317, 251318, 251319, 251320, 251321, 251322, 251323, 251324, 251325, 251326, 251327, 251328, 251329, 251330, 251331, 251332, 251333, 251334, 251335, 251336, 251337, 251338, 251339, 251340, 251341, 251342, 251343, 251344, 251345, 251346, 251347, 251348, 251349, 251350, 251351, 251352, 251353, 251354, 251355, 251356, 251357, 251358, 251359, 251360, 251361, 251362, 251363, 251364, 251365, 251366, 251367, 251368, 251369, 251370, 251371, 251372, 251373, 251374, 251375, 251376, 251377, 251378, 251379, 251380, 251381, 251382, 251383, 251384, 251385, 251386, 251387, 251388, 251389, 251390, 251391, 251392, 251393, 251394, 251395, 251396, 251397, 251398, 251399, 251400, 251401, 251402, 251403, 251404, 251405, 251406, 251407, 251408, 251409, 251410, 251411, 251412, 251413, 251414, 251415, 251416, 251417, 251418, 251419, 251420, 251421, 251422, 251423, 251424, 251425, 251426, 251427, 251428, 251429, 251430, 251431, 251432, 251433, 251434, 251435, 251436, 251437, 251438, 251439, 251440, 251441, 251442, 251443, 251444, 251445, 251446, 251447, 251448, 251449, 251450, 251451, 251452, 251453, 251454, 251455, 251456, 251457, 251458, 251459, 251460, 251461, 251462, 251463, 251464, 251465, 251466, 251467, 251468, 251469, 251470, 251471, 251472, 251473, 251474, 251475, 251476, 251477, 251478, 251479, 251480, 251481, 251482, 251483, 251484, 251485, 251486, 251487, 251488, 251489, 251490, 251491, 251492, 251493, 251494, 251495, 251496, 251497, 251498, 251499, 251500, 251501, 251502, 251503, 251504, 251505, 251506, 251507, 251508, 251509, 251510, 251511, 251512, 251513, 251514, 251515, 251516, 251517, 251518, 251519, 251520, 251521, 251522, 251523, 251524, 251525, 251526, 251527, 251528, 251529, 251530, 251531, 251532, 251533, 251534, 251535, 251536, 251537, 251538, 251539, 251540, 251541, 251542, 251543, 251544, 251545, 251546, 251547, 251548, 251549, 251550, 251551, 251552, 251553, 251554, 251555, 251556, 251557, 251558, 251559, 251560, 251561, 251562, 251563, 251564, 251565, 251566, 251567, 251568, 251569, 251570, 251571, 251572, 251573, 251574, 251575, 251576, 251577, 251578, 251579, 251580, 251581, 251582, 251583, 251584, 251585, 251586, 251587, 251588, 251589, 251590, 251591, 251592, 251593, 251594, 251595, 251596, 251597, 251598, 251599, 251600, 251601, 251602, 251603, 251604, 251605, 251606, 251607, 251608, 251609, 251610, 251611, 251612, 251613, 251614, 251615, 251616, 251617, 251618, 251619, 251620, 251621, 251622, 251623, 251624, 251625, 251626, 251627, 251628, 251629, 251630, 251631, 251632, 251633, 251634, 251635, 251636, 251637, 251638, 251639, 251640, 251641, 251642, 251643, 251644, 251645, 251646, 251647, 251648, 251649, 251650, 251651, 251652, 251653, 251654, 251655, 251656, 251657, 251658, 251659, 251660, 251661, 251662, 251663, 251664, 251665, 251666, 251667, 251668, 251669, 251670, 251671, 251672, 251673, 251674, 251675, 251676, 251677, 251678, 251679, 251680, 251681, 251682, 251683, 251684, 251685, 251686, 251687, 251688, 251689, 251690, 251691, 251692, 251693, 251694, 251695, 251696, 251697, 251698, 251699, 251700, 251701, 251702, 251703, 251704, 251705, 251706, 251707, 251708, 251709, 251710, 251711, 251712, 251713, 251714, 251715, 251716, 251717, 251718, 251719, 251720, 251721, 251722, 251723, 251724, 251725, 251726, 251727, 251728, 251729, 251730, 251731, 251732, 251733, 251734, 251735, 251736, 251737, 251738, 251739, 251740, 251741, 251742, 251743, 251744, 251745, 251746, 251747, 251748, 251749, 251750, 251751, 251752, 251753, 251754, 251755, 251756, 251757, 251758, 251759, 251760, 251761, 251762,
```

time hadoop jar hadoop-streaming-2.6.4.jar -input /data/dwdate.tbl -output /data/output1 -mapper BFmapper.py -reducer BFReducer.py -file BFmapper.py -file BFReducer.py -file bloom_filter.bloom



Services

Q Search

[Alt+S]

```
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 16:38:18
$ time hadoop jar hadoop-streaming-2.6.4.jar -input /data/dwdate.tbl -output /data/output1 -mapper BFmapper.py -reducer BFReducer.py -file BFmapper.py -file BFReducer.py -file bloom_filter.bloom
24/03/22 16:38:26 WARN Streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [BFmapper.py, BFReducer.py, bloom_filter.bloom, /tmp/hadoop-unjar210374695525448556/] [] /tmp/streamjob6254127611655366535.jar tmpDir=null
24/03/22 16:38:27 INFO client.RMProxy: Connecting to ResourceManager at /172.31.62.135:8032
24/03/22 16:38:27 INFO client.RMProxy: Connecting to ResourceManager at /172.31.62.135:8032
24/03/22 16:38:28 INFO mapred.FileInputFormat: Total input paths to process : 1
24/03/22 16:38:28 INFO mapreduce.JobSubmitter: number of splits:2
24/03/22 16:38:28 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1711125434068_0002
24/03/22 16:38:28 INFO impl.YarnClientImpl: Submitted application application_1711125434068_0002
24/03/22 16:38:28 INFO mapreduce.Job: The url to track the job: http://ip-172-31-62-135.ec2.internal:8088/proxy/application_1711125434068_0002/
24/03/22 16:38:28 INFO mapreduce.Job: Running job: job_1711125434068_0002
24/03/22 16:38:34 INFO mapreduce.Job: Job job_1711125434068_0002 running in uber mode : false
24/03/22 16:38:34 INFO mapreduce.Job: map 0% reduce 0%
24/03/22 16:38:40 INFO mapreduce.Job: map 50% reduce 0%
24/03/22 16:38:44 INFO mapreduce.Job: map 100% reduce 0%
24/03/22 16:38:48 INFO mapreduce.Job: map 100% reduce 100%
24/03/22 16:38:48 INFO mapreduce.Job: Job job_1711125434068_0002 completed successfully
24/03/22 16:38:48 INFO mapreduce.Job: Counters: 49

File System Counters
  FILE: Number of bytes read=9522
  FILE: Number of bytes written=351044
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=234237
  HDFS: Number of bytes written=8784
  HDFS: Number of read operations=9
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2


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)=29052
  Total time spent by all reduces in occupied slots (ms)=14628
  Total time spent by all map tasks (ms)=4842
  Total time spent by all reduce tasks (ms)=2438
  Total vcore-milliseconds taken by all map tasks=4842
  Total vcore-milliseconds taken by all reduce tasks=2438
  Total megabyte-milliseconds taken by all map tasks=29052000
  Total megabyte-milliseconds taken by all reduce tasks=14628000

Map-Reduce Framework
  Map input records=2556
  Map output records=366
  Map output bytes=8784
  Map output materialized bytes=9528
  Input split bytes=176
  Combine input records=0
  Combine output records=0
  Reduce input groups=366
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135

Run Time and Output:

 Services [Alt+S]

```
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=234061
File Output Format Counters
  Bytes Written=8784
24/03/22 16:38:48 INFO streaming.StreamJob: Output directory: /data/output1

real    0m22.342s
user    0m3.758s
sys     0m0.194s
[ec2-user@ip-172-31-62-135 ~] 2024-03-22 16:38:48
$ hdfs dfs -cat /data/output1/part-00000 | more
19920901|19920901|Fall
19920902|19920902|Fall
19920903|19920903|Fall
19920904|19920904|Fall
19920905|19920905|Fall
19920906|19920906|Fall
19920907|19920907|Fall
19920908|19920908|Fall
19920909|19920909|Fall
19920910|19920910|Fall
19920911|19920911|Fall
19920912|19920912|Fall
19920913|19920913|Fall
19920914|19920914|Fall
19920915|19920915|Fall
19920916|19920916|Fall
19920917|19920917|Fall
19920918|19920918|Fall
19920919|19920919|Fall
19920920|19920920|Fall
19920921|19920921|Fall
19920922|19920922|Fall
19920923|19920923|Fall
19920924|19920924|Fall
19920925|19920925|Fall
19920926|19920926|Fall
19920927|19920927|Fall
19920928|19920928|Fall
19920929|19920929|Fall
19920930|19920930|Fall
19921001|19921001|Fall
19921002|19921002|Fall
19921003|19921003|Fall
19921004|19921004|Fall
19921005|19921005|Fall
19921006|19921006|Fall
19921007|19921007|Fall
19921008|19921008|Fall
```

i-062d65a3dd2e0d019 (nachiketh_server)

PublicIPs: 34.207.253.48 PrivateIPs: 172.31.62.135