

Q1a)

I divided this question into two map-reduce job.

The first map task is print a key-value pair similar to the input file.

The first reducer task print a key and a set, which the set will contain all the followee of that corresponding key.

The second mapper takes the first reducer output and generate every pair of possible combinations. If the length of their intersection is greater than 0, print it out.

The second reducer is just choose the greatest number of intersection of each blog id and print the information out.

1st MapReduce job:

Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file  
~/hw1/q_one_a_first_mapper.py -mapper ~/hw1/q_one_a_first_mapper.py -file  
~/hw1/q_one_a_first_reducer.py -reducer ~/hw1/q_one_a_first_reducer.py -input  
medium/medium_relation -output q_one_a_1_output
```

```
[s1155157657@dicvmc4 hadoop]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file ~/hw1/q_one_a_first_mapper.py -mapper ~/hw1/q_one_a_first_mapper.py -file ~/hw1/q_one_a_first_reducer.py -reducer ~/hw1/q_one_a_first_reducer.py -input medium/medium_relation -output q_one_a_1_output  
23/10/08 05:36:59 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.  
packageJobJar: [/home/s1155157657/hw1/q_one_a_first_mapper.py, /home/s1155157657/hw1/q_one_a_first_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob561981118697819953.jar  
tmpDir=null  
23/10/08 05:37:00 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/08 05:37:00 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/08 05:37:01 INFO mapred.FileInputFormat: Total input files to process : 1  
23/10/08 05:37:01 INFO mapreduce.JobSubmitter: number of splits:2  
23/10/08 05:37:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1036  
23/10/08 05:37:01 INFO conf.Configuration: resource-types.xml not found  
23/10/08 05:37:01 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.  
23/10/08 05:37:01 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE  
23/10/08 05:37:01 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE  
23/10/08 05:37:01 INFO impl.YarnClientImpl: Submitted application application_1694578679658_1036  
23/10/08 05:37:01 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1036/  
23/10/08 05:37:01 INFO mapreduce.Job: Running job: job_1694578679658_1036  
23/10/08 05:37:06 INFO mapreduce.Job: Job job_1694578679658_1036 running in uber mode : false  
23/10/08 05:37:06 INFO mapreduce.Job: map 0% reduce 0%  
23/10/08 05:37:12 INFO mapreduce.Job: map 100% reduce 0%  
23/10/08 05:37:19 INFO mapreduce.Job: map 100% reduce 100%  
23/10/08 05:37:19 INFO mapreduce.Job: Job job_1694578679658_1036 completed successfully  
23/10/08 05:37:19 INFO mapreduce.Job: Counters: 49  
File System Counters  
FILE: Number of bytes read=36003774  
FILE: Number of bytes written=72663317  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=32486915  
HDFS: Number of bytes written=18658436  
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  
Rack-local map tasks=2  
Total time spent by all maps in occupied slots (ms)=37756  
Total time spent by all reduces in occupied slots (ms)=36264  
Total time spent by all map tasks (ms)=9439  
Total time spent by all reduce tasks (ms)=4533  
Total vcore-milliseconds taken by all map tasks=9439  
Total vcore-milliseconds taken by all reduce tasks=4533  
Total megabyte-milliseconds taken by all map tasks=38662144  
Total megabyte-milliseconds taken by all reduce tasks=37134336  
Map-Reduce Framework  
Map input records=1768149  
Map output records=1768149  
Map output bytes=32467470  
Map output materialized bytes=36003780  
Input split bytes=252  
Combine input records=0  
Combine output records=0  
Reduce input groups=70097  
Reduce shuffle bytes=36003780
```

Code:

Mapper:

```
1 #!/usr/bin/env python3  
2 """q_one_a_first_mapper.py"""  
3  
4 import sys  
5  
6 # input comes from STDIN (standard input)  
7 for line in sys.stdin:  
8     # remove leading and trailing whitespace  
9     line = line.strip()  
10    # split the line into words  
11    words = line.split(" ")  
12    #print the pairs  
13    print('%s\t%s' % (words[0], words[1]))
```

Reducer:

```
1  #!/usr/bin/env python3
2  """q_one_a_first_reducer.py"""
3
4  import sys
5
6  current_id = None
7  current_set = set()
8
9  # input comes from STDIN
10 for line in sys.stdin:
11     # remove leading and trailing whitespace
12     line = line.strip()
13
14     # parse the input we got from mapper.py
15     follower, followee = line.split('\t')
16
17     #Add the followee to the set if the key is the same
18     if current_id == follower:
19         current_set.add(int(followee))
20     else:
21         if current_id:
22             print(f"{current_id}\t{current_set}")
23         current_id = follower
24         current_set.clear()
25         current_set.add(int(followee))
26
27 # print the last follower too
28 if current_id == follower:
29     print(f"{current_id}\t{current_set}")
```

2nd MapReduce job:

Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -D
mapred.map.tasks=20 -D mapred.reduce.tasks=1 -files
hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_a_1_output/part-00000 -file
~/hw1/q_one_a_second_mapper.py -mapper ~/hw1/q_one_a_second_mapper.py -file
~/hw1/q_one_a_second_reducer.py -reducer ~/hw1/q_one_a_second_reducer.py -input
q_one_a_1_output/part-00000 -output q_one_a_two_16_output
```

```
[s1155157657@dicvmc4 ~]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -D mapred.map.tasks=20 -D mapred.reduce.tasks=1 -files hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_a_1_0
utput/part-00000 -file ~/hw1/q_one_a_second_mapper.py -mapper ~/hw1/q_one_a_second_mapper.py -file ~/hw1/q_one_a_second_reducer.py -reducer ~/hw1/q_one_a_second_reducer.py -input q_one_a_1_output/part-
00000 -output q_one_a_two_16_output
23/10/08 21:55:30 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/s1155157657/hw1/q_one_a_second_mapper.py, /home/s1155157657/hw1/q_one_a_second_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob2998395375244070186.j
ar tmpDir=null
23/10/08 21:55:30 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/08 21:55:30 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/08 21:55:30 INFO mapred.FileInputFormat: Total input files to process : 1
23/10/08 21:55:30 INFO mapreduce.JobSubmitter: number of splits:20
23/10/08 21:55:31 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1435
23/10/08 21:55:31 INFO conf.Configuration: resource-types.xml not found
23/10/08 21:55:31 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
23/10/08 21:55:31 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
23/10/08 21:55:31 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
23/10/08 21:55:31 INFO impl.YarnClientImpl: Submitted application application_1694578679658_1435
23/10/08 21:55:31 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1435/
23/10/08 21:55:31 INFO mapreduce.Job: Running job: job_1694578679658_1435
23/10/08 21:55:37 INFO mapreduce.Job: Job job_1694578679658_1435 running in uber mode : false
23/10/08 21:55:37 INFO mapreduce.Job: map 0% reduce 0%
23/10/08 21:55:45 INFO mapreduce.Job: Task Id : attempt_1694578679658_1435_m_000000_0, Status : FAILED
[2023-10-08 21:55:44.194]Container killed on request. Exit code is 137
[2023-10-08 21:55:44.195]Container exited with a non-zero exit code 137.
[2023-10-08 21:55:44.195]Killed by external signal
23/10/08 21:55:50 INFO mapreduce.Job: map 2% reduce 0%
23/10/08 21:55:51 INFO mapreduce.Job: map 3% reduce 0%
```

Code:

Mapper:

```
hw1> q_one_a_second_mapper.py > ...
1 #!/usr/bin/env python3
2 """q_one_a_second_mapper.py"""
3
4 import sys
5 import ast
6
7 cache_lines = []
8
9 file_path = "part-00000"
10
11 # read the file from first mapreduce output
12 with open(file_path, 'r') as file:
13     for line in file:
14         # Do something with each line
15         line = line.strip()
16         temp_follower, temp_followee = line.split("\t",2)
17         # turn the value to a set
18         cache_lines.append((temp_follower, ast.literal_eval(temp_followee)))
19
20 # input comes from STDIN (standard input)
21 # for every line of input
22 for line in sys.stdin:
23     max = 0
24     # remove leading and trailing whitespace
25     line = line.strip()
26     # split the line into words
27     follower, followee_string = line.split("\t", 2)
28
29     #process the value to a set
30     followee_set = ast.literal_eval(followee_string)
31
32     #loop through the file read from first mapreduce output
33     for cache_line in cache_lines:
34         cache_follower, cache_followee_set = cache_line
35         if cache_follower != follower:
36             # find the intersection of the two sets
37             intersection = followee_set.intersection(cache_followee_set)
38             set_length = len(intersection)
39
40             # if the length of the intersection is greater than 0
41             if set_length >= max and set_length > 0:
42                 max = set_length
43                 # print the follower pairs, their intersection set and the length of the intersection
44                 print(f"{follower}\t{cache_follower}\t{intersection}\t{len(intersection)}")
```

Reducer:

```
1  #!/usr/bin/env python3
2  """q_one_a_second_reducer.py"""
3
4  import sys
5  import ast
6
7  current_website_one = None
8  current_website_two = None
9  current_followee_set = set()
10 current_followee_number = None
11 website_one = -1
12 # input comes from STDIN
13 for line in sys.stdin:
14     # remove leading and trailing whitespace
15     line = line.strip()
16
17     # parse the input we got from mapper.py
18     website_one, website_two, common_followee_string, common_followee_number = line.split('\t')
19
20     # output the greatest length of a follower found
21     if website_one != current_website_one:
22         if current_website_one:
23             print(f"{current_website_one}:{current_website_two}, {current_followee_set}, {current_followee_number}")
24
25         current_website_one = website_one
26         current_website_two = website_two
27         current_followee_set.clear()
28         current_followee_set = ast.literal_eval(common_followee_string)
29         current_followee_number = len(current_followee_set)
30
31     else:
32         # check if the id is smaller when it is a tie in length/ length of set is longer
33         if int(common_followee_number) > int(current_followee_number):
34             current_website_one = website_one
35             current_website_two = website_two
36             current_followee_set.clear()
37             current_followee_set = ast.literal_eval(common_followee_string)
38             current_followee_number = len(current_followee_set)
39         elif int(common_followee_number) == int(current_followee_number) and int(current_website_two) > int(website_two):
40             current_website_one = website_one
41             current_website_two = website_two
42             current_followee_set.clear()
43             current_followee_set = ast.literal_eval(common_followee_string)
44             current_followee_number = len(current_followee_set)
45
46 if current_website_one == website_one:
47     print(f"{current_website_one}:{current_website_two}, {current_followee_set}, {current_followee_number}")
```

Final Output:

grep the blog ID ends with 7657 (my SID ends with 7657):

Command:

hadoop fs -cat q_one_a_two_16_output/part-00000 | grep "7657:"

```
[*1155157657@dc1vmc4 ~]$ hadoop fs -cat q_one_a_two_16_output/part-00000 | grep "7657:"
115/657:816394, {21347716, 72998223, 822388, 14296729, 14525658}, 6
158/657:18472742, {549795296, 112455655, 18744963, 25948370, 12946898, 19519258}, 6
17779/657:31289, {23838217, 11447858, 111691460, 71256938}, 4
19347/657:15693499, {14089201, 116485857, 18838275}, 3
42478/657:121569952, {248261248, 35791376, 16886231, 237428963, 284422694, 121773488, 94918923, 177345934, 153972655, 3359857, 118244437, 312199255, 18671608, 105643929, 116143618, 278109085}, 16
```

Q1b)

In this question I used the output from job1 in Q1a.

Based on the output of the 1st MapReduce job, we have a key and a set pair which the set contains all followee of the key. We only do one MapReduce job this time.

Mapper:

The mapper generates all possible pairs of keys (except with itself). If the length of their intersection of their corresponding set is greater than 0, then it output the similarity of this pair of keys and the intersection set (common followee).

```
1  #!/usr/bin/env python3
2  """q_one_b_first_mapper.py"""
3
4  import sys
5  import ast
6
7  cache_lines = []
8
9  file_path = "part-00000"
10
11  with open(file_path, 'r') as file:
12      for line in file:
13          # Do something with each line
14          line = line.strip()
15          temp_follower, temp_followee = line.split("\t",2)
16          # make the value a set
17          cache_lines.append((temp_follower, ast.literal_eval(temp_followee)))
18
19
20  # input comes from STDIN (standard input)
21  for line in sys.stdin:
22      # remove leading and trailing whitespace
23      line = line.strip()
24      # split the line into words
25      follower, followee_string = line.split("\t",2)
26      followee_set = ast.literal_eval(followee_string)
27
28      # find a pair of keys (self excluded)
29      # if their length of intersection set is not 0, return the similarity and the intersection set
30      for cache_line in cache_lines:
31          cache_follower, cache_followee_set = cache_line
32          if cache_follower != follower:
33              intersection = followee_set.intersection(cache_followee_set)
34
35              if len(intersection) != 0:
36                  union_len = len(followee_set) + len(cache_followee_set) - len(intersection)
37                  print(f"{follower}\t{cache_follower}\t{len(intersection)/union_len}\t{intersection}")
```

Reducer:

In the reduce stage I perform sorting for all the values for each key and take top K row of each key.

```
1 #!/usr/bin/env python3
2 """q_one_b_first_reducer.py"""
3
4 import sys
5 import ast
6
7 current_website_one = None
8 current_website_two = None
9 current_followee_set = set()
10 current_followee_number = None
11 top_k_list = []
12 k = 3
13 website_one = -1
14
15 # input comes from STDIN
16 for line in sys.stdin:
17     # remove leading and trailing whitespace
18     line = line.strip()
19
20     # parse the input we got from mapper.py
21     website_one, website_two, similar_score, common_followee_string = line.split('\t')
22     current_followee_set = ast.literal_eval(common_followee_string)
23
24     # print top K result if key changed
25     if current_website_one != website_one:
26         if current_website_one:
27             if len(top_k_list) > k:
28                 top_k_list = sorted(top_k_list, reverse=True, key=lambda x: (x[2]))
29                 for i in range(k):
30                     print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
31             elif len(top_k_list) == k:
32                 for i in range(k):
33                     print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
34         else:
35             for i in range(len(top_k_list)):
36                 print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
37
38         current_website_one = website_one
39         top_k_list = []
40
41     # add the information for the same key to set
42     top_k_list.append((website_one, website_two, similar_score, current_followee_set))
43
44     # print top K result for the last one if key didn't changed
45     if current_website_one == website_one:
46         if current_website_one:
47             if len(top_k_list) > k:
48                 top_k_list = sorted(top_k_list, reverse=True, key=lambda x: (x[2]))
49                 for i in range(k):
50                     print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
51             elif len(top_k_list) == k:
52                 for i in range(k):
53                     print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
54         else:
55             for i in range(len(top_k_list)):
56                 print(f"{top_k_list[i][0]}:{top_k_list[i][1]}, {top_k_list[i][3]}, {top_k_list[i][2]}")
```

MapReduce job:

Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -D
mapred.map.tasks=10 -D mapred.reduce.tasks=3 -files
hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_a_1_output/part-00000 -file
~/hw1/q_one_b_first_mapper.py -mapper ~/hw1/q_one_b_first_mapper.py -file
~/hw1/q_one_b_first_reducer.py -reducer ~/hw1/q_one_b_first_reducer.py -input
q_one_a_1_output/part-00000 -output q_one_b_5_output
```

```
[s1155157657@dicvmc2 ~]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -D mapred.map.tasks=10 -D mapred.reduce.tasks=3 -files hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_a_1_output/part-00000 -file ~/hw1/q_one_b_first_mapper.py -mapper ~/hw1/q_one_b_first_mapper.py -file ~/hw1/q_one_b_first_reducer.py -reducer ~/hw1/q_one_b_first_reducer.py -input q_one_a_1_output/part-00000 -output q_one_b_5_output
23/10/08 23:18:40 INFO streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/s1155157657/hw1/q_one_b_first_mapper.py, /home/s1155157657/hw1/q_one_b_first_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob8250016589354120990.jar tmpDir=null
23/10/08 23:18:40 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/08 23:18:40 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/08 23:18:41 INFO mapred.FileInputFormat: Total input files to process : 1
23/10/08 23:18:41 INFO mapreduce.JobSubmitter: number of splits:10
23/10/08 23:18:41 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1470
23/10/08 23:18:41 INFO conf.Configuration: resource-types.xml not found
23/10/08 23:18:41 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
23/10/08 23:18:41 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
23/10/08 23:18:41 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
23/10/08 23:18:41 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1470/
23/10/08 23:18:41 INFO mapreduce.Job: Running job: job_1694578679658_1470
23/10/08 23:18:45 INFO mapreduce.Job: Job job_1694578679658_1470 running in uber mode : false
23/10/08 23:18:45 INFO mapreduce.Job: map 0% reduce 0%
23/10/08 23:18:56 INFO mapreduce.Job: map 1% reduce 0%
23/10/08 23:18:57 INFO mapreduce.Job: map 3% reduce 0%
23/10/08 23:19:00 INFO mapreduce.Job: map 4% reduce 0%
23/10/08 23:19:09 INFO mapreduce.Job: map 5% reduce 0%
23/10/08 23:19:12 INFO mapreduce.Job: map 6% reduce 0%
```

Output:

grep the blog ID ends with 7657 (my SID ends with 7657):

Command:

`hadoop fs -cat q_one_b_3_output/part-0000x | grep "7657:"`

(which x is 0,1,2 because I used 3 reducer for my MapReduce job)

```
[s1155157657@dc1cvmc4 hw115]$ hadoop fs -cat q_one_b_5_output/part-00000 | grep "7657:"
1157657:70977926, {14296729, 14525658, 822388, 72998223}, 0.2857142857142857
1157657:15822197, {14296729, 14525658}, 0.2222222222222222
1157657:48137812, {21347716, 72998223, 822388, 14296729, 14525658}, 0.19238769238769232
19347657:171158575, {14889281, 18838275}, 0.4
19347657:17581519, {14889281, 18838275}, 0.4
19347657:509138385, {115486057}, 0.3333333333333333
[s1155157657@dc1cvmc4 hw115]$ hadoop fs -cat q_one_b_5_output/part-00001 | grep "7657:"
177707657:15910847, {11447858, 71256938}, 0.4
177707657:45883825, {71256938}, 0.25
177707657:156724888, {11447858}, 0.2
[s1155157657@dc1cvmc4 hw115]$ hadoop fs -cat q_one_b_5_output/part-00002 | grep "7657:"
1507657:72988479, {14245665, 25940379, 12946898}, 0.375
1507657:21156529, {14245665, 25940379, 12946898, 19519258}, 0.36363636363636365
1507657:18453798, {14245665, 12946898}, 0.25
424787657:118244437, {35781376, 284422694, 121773488, 153972655, 312109255, 18671688, 116143610, 278189885}, 0.32
424787657:182794595, {153972655, 3359857, 312109255, 18671688, 185643929, 278189885}, 0.25
424787657:322923214, {2480261248, 237420963, 94918923, 312109255, 185643929, 116143610, 278189885}, 0.2413793103448276
```

Q1c)

I used two MapReduce job in this question.

The first Mapper task is to take the medium_relation input and output the followee as key and the follower as value.

The first Reducer task is to output the blog id that has two or more follower.

Then, the second Mapper task is to take the output file from the first task and make it a list.

Also take medium_label as input file. Read medium_label line by line and if the blog id is in the list, output a <community, 1> intermediate key if yes.

The second reducer simply sum it up for each community key and output the final value.

1st MapReduce job:

Code:

Mapper:

```
1  #!/usr/bin/env python3
2  """q_one_c_first_mapper.py"""
3
4  import sys
5
6  # input comes from STDIN (standard input)
7  for line in sys.stdin:
8      # remove leading and trailing whitespace
9      line = line.strip()
10     # split the line into words
11     words = line.split(" ")
12     # print the followee as key and follower as value this time
13     print('%s\t%s' % (words[1], words[0]))
```

Reducer:

```
1  #!/usr/bin/env python3
2  """q_one_c_first_reducer.py"""
3
4  import sys
5
6  current_id = None
7  current_set = set()
8
9  # input comes from STDIN
10 for line in sys.stdin:
11     # remove leading and trailing whitespace
12     line = line.strip()
13
14     # parse the input we got from mapper.py
15     follower, followee = line.split('\t')
16
17     # add the follower to the set of the corresponding followee key
18     if current_id == followee:
19         current_set.add(int(follower))
20     else:
21         if current_id and len(current_set) >= 2:
22             print(f"{current_id}")
23             current_id = followee
24             current_set.clear()
25             current_set.add(int(follower))
26
27 # if the set contains two or more follower, output the followee key
28 if current_id == followee and len(current_set) >= 2:
29     print(f"{followee}")
```


Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file  
~/hw1/q_one_c_first_mapper.py -mapper ~/hw1/q_one_c_first_mapper.py -file  
~/hw1/q_one_c_first_reducer.py -reducer ~/hw1/q_one_c_first_reducer.py -input  
medium/medium_relation -output q_one_c_3_output
```

```
[s1155157457@dicvmc4 hw1]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file ~/hw1/q_one_c_first_mapper.py -mapper ~/hw1/q_one_c_first_mapper.py -file ~/hw1/q_one_c_first_reducer.p  
y -reducer ~/hw1/q_one_c_first_reducer.py -input medium/medium_relation -output q_one_c_3_output  
23/10/09 01:01:30 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.  
packageJobJar: [/home/s1155157457/hw1/q_one_c_first_mapper.py, /home/s1155157457/hw1/q_one_c_first_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob278534886454403552.jar  
tmpDir=null  
23/10/09 01:01:31 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/09 01:01:31 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/09 01:01:32 INFO mapred.FileInputFormat: Total input files to process : 1  
23/10/09 01:01:32 INFO mapreduce.JobSubmitter: number of splits:2  
23/10/09 01:01:32 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1507  
23/10/09 01:01:32 INFO conf.Configuration: resource-types.xml not found  
23/10/09 01:01:32 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.  
23/10/09 01:01:32 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE  
23/10/09 01:01:32 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE  
23/10/09 01:01:32 INFO Impl.YarnClientImpl: Submitted application application_1694578679658_1507  
23/10/09 01:01:32 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1507/  
23/10/09 01:01:32 INFO mapreduce.Job: Running job: job_1694578679658_1507  
23/10/09 01:01:39 INFO mapreduce.Job: Job job_1694578679658_1507 running in uber mode : false  
23/10/09 01:01:39 INFO mapreduce.Job: map 0% reduce 0%  
23/10/09 01:01:49 INFO mapreduce.Job: map 100% reduce 0%  
23/10/09 01:01:58 INFO mapreduce.Job: map 100% reduce 100%  
23/10/09 01:01:58 INFO mapreduce.Job: Job job_1694578679658_1507 completed successfully  
23/10/09 01:01:58 INFO mapreduce.Job: Counters: 49  
File System Counters  
FILE: Number of bytes read=36863774  
FILE: Number of bytes written=72663320  
FILE: Number of read operations=0  
FILE: Number of large read operations=0  
FILE: Number of write operations=0  
HDFS: Number of bytes read=32598794  
HDFS: Number of bytes written=730203  
HDFS: Number of read operations=9  
HDFS: Number of large read operations=0  
HDFS: Number of write operations=2
```

2nd MapReduce job:

Code:

Mapper:

```
1  #!/usr/bin/env python3
2  """q_one_c_second_mapper.py"""
3
4  import sys
5
6  cache_lines = []
7
8  file_path = "part-00000"
9
10 with open(file_path, 'r') as file:
11     for line in file:
12         # Do something with each line
13         line = line.strip()
14         cache_lines.append(int(line))
15
16 # input comes from STDIN (standard input)
17 for line in sys.stdin:
18     # remove leading and trailing whitespace
19     line = line.strip()
20     # split the line into words
21     follower, community = line.split(" ", 2)
22     # print(follower_set)
23
24     # if the follower is in the first output file, the emit <community, 1>
25     if int(follower) in cache_lines:
26         print(f"{int(community)}\t1")
```

Reducer:

```
1  #!/usr/bin/env python3
2  """q_one_c_second_reducer.py"""
3
4  import sys
5
6  current_community = None
7  current_count = 0
8  community = -1
9
10 # input comes from STDIN
11 for line in sys.stdin:
12     # remove leading and trailing whitespace
13     line = line.strip()
14
15     # parse the input we got from mapper.py
16     community, number = line.split('\t')
17     number = int(number)
18
19     # sum the numbers up for each community
20     if community != current_community:
21         if current_community:
22             print(f"Community {current_community}: {current_count}")
23
24         current_community = community
25         current_count = number
26
27     else:
28         current_count += number
29
30 # output if the last community is not outputted
31 if current_community == community:
32     print(f"Community {current_community}: {current_count}")
```

Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -files
hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_c_3_output/part-00000 -file
~/hw1/q_one_c_second_mapper.py -mapper ~/hw1/q_one_c_second_mapper.py -file
~/hw1/q_one_c_second_reducer.py -reducer ~/hw1/q_one_c_second_reducer.py -input
medium/medium_label -output q_one_c_second_2_output
```

```
[s1155157657@dicvmc4 hw1]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -files hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_c_3_output/part-00000 -file ~/hw1/q_one_c_second_mapper.py -mapper ~/hw1/q_one_c_second_mapper.py -file ~/hw1/q_one_c_second_reducer.py -reducer ~/hw1/q_one_c_second_reducer.py -input medium/medium_label -output q_one_c_second_2_output
23/10/09 01:03:57 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/s1155157657/hw1/q_one_c_second_mapper.py, /home/s1155157657/hw1/q_one_c_second_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob7544825735914073835.jar tmpDir=null
23/10/09 01:03:57 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/09 01:03:57 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200
23/10/09 01:04:01 INFO mapred.FileInputFormat: Total input files to process : 1
23/10/09 01:04:01 INFO mapreduce.JobSubmitter: number of splits:2
23/10/09 01:04:01 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1509
23/10/09 01:04:02 INFO conf.Configuration: resource-types.xml not found
23/10/09 01:04:02 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
23/10/09 01:04:02 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE
23/10/09 01:04:02 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
23/10/09 01:04:02 INFO impl.YarnClientImpl: Submitted application application_1694578679658_1509
23/10/09 01:04:02 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1509/
23/10/09 01:04:02 INFO mapreduce.Job: Running job: job_1694578679658_1509
23/10/09 01:04:11 INFO mapreduce.Job: Job job_1694578679658_1509 running in uber mode : false
23/10/09 01:04:11 INFO mapreduce.Job: map 0% reduce 0%
23/10/09 01:04:21 INFO mapreduce.Job: map 19% reduce 0%
23/10/09 01:04:22 INFO mapreduce.Job: map 30% reduce 0%
23/10/09 01:04:24 INFO mapreduce.Job: map 48% reduce 0%
23/10/09 01:04:25 INFO mapreduce.Job: map 57% reduce 0%
23/10/09 01:04:27 INFO mapreduce.Job: map 62% reduce 0%
23/10/09 01:04:31 INFO mapreduce.Job: map 67% reduce 0%
23/10/09 01:04:33 INFO mapreduce.Job: map 83% reduce 0%
23/10/09 01:04:34 INFO mapreduce.Job: map 100% reduce 0%
23/10/09 01:04:38 INFO mapreduce.Job: map 100% reduce 100%
23/10/09 01:04:38 INFO mapreduce.Job: Job job_1694578679658_1509 completed successfully
```

Output:

```
[s1155157657@dicvmc4 hw1]$ hadoop fs -cat q_one_c_second_2_output/part-00000
Community 0: 23633
Community 1: 23751
Community 2: 23705
```

Q1d)

#Job	Mapper num	Reducer num	Max mapper time	Min mapper time	Avg mapper time	Max reducer time	Min reducer time	Avg reducer time	Total time
1 (1036)	2	1	4s	4s	4s	2s	2s	2s	13s
2 (1041)	2	1	31m52s	29m46s	30m49s	9s	9s	9s	32m04s

#Job	Mapper num	Reducer num	Max mapper time	Min mapper time	Avg mapper time	Max reducer time	Min reducer time	Avg reducer time	Total time
1 (1562)	5	1	18s	3s	15s	3s	3s	3s	13s
2 (1584)	20	1	3m42s	2m58s	3m15s	17s	17s	17s	25m48s

#Job	Mapper num	Reducer num	Max mapper time	Min mapper time	Avg mapper time	Max reducer time	Min reducer time	Avg reducer time	Total time
1 (1596)	3	1	4s	4s	4s	2s	2s	2s	14s
2 (1598)	10	1	6m58s	5m56s	6m18s	16s	16s	16s	22m40s

#Job	Mapper num	Reducer num	Max mapper time	Min mapper time	Avg mapper time	Max reducer time	Min reducer time	Avg reducer time	Total time
1 (1614)	10	1	3s	3s	3s	2s	2s	2s	18s
2 (1615)	10	3	7m13s	6m1s	6m31s	4s	4s	4s	27m43s

Discovery:

For job1, since it is a very lightweight job, increasing amount of mapper does not help to boost the mapping stage but to slow it down on the whole progress.

For job2, since it requires a lot of system resources on mapping stage, so more mapper meaning more system resources can be utilized for mapper. It is generally faster if we use more mapper on mapping stage. However, it does not necessarily guarantee the whole job is

faster since it requires more time to shuffle and sort the result from more mappers. It slows down a bit when the mapper is too much.

Increasing reducer does not help in this case too because it is not reducer heavy job. No big performance boost is observed.

Bonus

I tried to run my code but it cannot work on the IE server. The first MapReduce job run successfully and it generates a key value pair where the key is the follower and the value is a set of followee of that follower.

The first Job runs fine but the second one is messed up. The second job takes up too much memory

Command:

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file  
~/hw1/q_one_a_first_mapper.py -mapper ~/hw1/q_one_a_first_mapper.py -file  
~/hw1/q_one_a_first_reducer.py -reducer ~/hw1/q_one_a_first_reducer.py -input  
large/large_relation -output q_one_e_1_output
```

```
[s1155157657@dicvmc4 ~]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -file ~/hw1/q_one_a_first_mapper.py -mapper ~/hw1/q_one_a_first_mapper.py -file ~/hw1/q_one_a_first_reducer.py -reducer ~/hw1/q_one_a_first_reducer.py -input large/large_relation -output q_one_e_1_output  
23/10/09 06:44:54 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.  
packageJobJar: [/home/s1155157657/hw1/q_one_a_first_mapper.py, /home/s1155157657/hw1/q_one_a_first_reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar] /tmp/streamjob6394153205728634333.jar  
tmpDir=null  
23/10/09 06:44:55 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/09 06:44:55 INFO client.AHSProxy: Connecting to Application History server at dicvmc1.ie.cuhk.edu.hk/172.16.5.161:10200  
23/10/09 06:44:55 INFO mapred.FileInputFormat: Total input files to process : 1  
23/10/09 06:44:55 INFO mapreduce.JobSubmitter: number of splits:3  
23/10/09 06:44:55 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1694578679658_1624  
23/10/09 06:44:55 INFO conf.Configuration: resource-types.xml not found  
23/10/09 06:44:55 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.  
23/10/09 06:44:55 INFO resource.ResourceUtils: Adding resource type - name = memory-mb, units = Mi, type = COUNTABLE  
23/10/09 06:44:55 INFO resource.ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE  
23/10/09 06:44:56 INFO mapreduce.Job: The url to track the job: http://dicvmc3.ie.cuhk.edu.hk:8088/proxy/application_1694578679658_1624/  
23/10/09 06:44:56 INFO mapreduce.Job: Running job: job_1694578679658_1624  
23/10/09 06:45:00 INFO mapreduce.Job: Job job_1694578679658_1624 running in uber mode : false  
23/10/09 06:45:00 INFO mapreduce.Job: map 0% reduce 0%  
23/10/09 06:45:07 INFO mapreduce.Job: map 33% reduce 0%  
23/10/09 06:45:10 INFO mapreduce.Job: map 78% reduce 0%  
23/10/09 06:45:15 INFO mapreduce.Job: map 100% reduce 0%  
23/10/09 06:45:17 INFO mapreduce.Job: map 100% reduce 54%  
23/10/09 06:45:20 INFO mapreduce.Job: map 100% reduce 75%  
23/10/09 06:45:23 INFO mapreduce.Job: map 100% reduce 87%  
23/10/09 06:45:26 INFO mapreduce.Job: map 100% reduce 100%  
23/10/09 06:45:26 INFO mapreduce.Job: Job job_1694578679658_1624 completed successfully
```

```
hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.10.1.jar -D  
mapred.map.tasks=15 -D mapreduce.map.output.compress=true -files  
hdfs://dicvmc2.ie.cuhk.edu.hk/user/s1155157657/q_one_e_1_output/part-00000 -file  
~/hw1/q_one_e_first_mapper.py -mapper ~/hw1/q_one_e_first_mapper.py -file  
~/hw1/q_one_e_first_reducer.py -reducer ~/hw1/q_one_e_first_reducer.py -input  
q_one_e_1_output/part-00000 -output q_one_e_second_12_output
```

some errors which shows cannot allocate that much memory:

```
[2023-10-09 17:01:16.795]Container exited with a non-zero exit code 1. Error file: prelaunch.err.  
Last 4096 bytes of prelaunch.err :  
Last 4096 bytes of stderr :  
Java HotSpot(TM) 64-Bit Server VM warning: INFO: os::commit_memory(0x00000000c2600000, 1879572480, 0) failed; error='Cannot allocate memory' (errno=12)  
  
[2023-10-09 17:01:16.795]Container exited with a non-zero exit code 1. Error file: prelaunch.err.  
Last 4096 bytes of prelaunch.err :  
Last 4096 bytes of stderr :  
Java HotSpot(TM) 64-Bit Server VM warning: INFO: os::commit_memory(0x00000000c2600000, 1879572480, 0) failed; error='Cannot allocate memory' (errno=12)
```

(Please give me some points for at least trying^^)

Q2a)

1)

Port 9000 is responsible for metadata services. It provides communication between datanode and resource manager. It is where the hdfs are connected to.

Port 8088 is yarn resource manager. It facilitates the execution of tasks in a job using existing resources available.

2)

I am using private IP. Using private IP is a saver than using public IP as public IP is exposed to security threads.

3)

No, as stated in the assignment we kept the internal communication port in gcp and it will allow communication of all machine in the same network. No extra firewall rules need to be setup.

Q2b)

In general, take up 150GB of total space in Hadoop using 4 VMs with 100GB storage each is feasible since $150\text{GB} < 400\text{GB}$. However, we should consider other factors.

If the storage in the VMs cannot be fully utilized, such as they already have some big files that takes up bunch of storage. In that case, the total storage of hdfs may be less than 150GB. So, it does not have resources to store the 150GB in total in hdfs.

Also if you are refereeing to 150GB before storing it to hdfs, since by default it will have 3 replication and be stored to different VMs. So In this case it does not have enough storage for $150\text{GB} * 3 = 450\text{GB}$ to store it.