

# IERG4300 / ESTR4300 Fall 2020 Homework 1

Release date: Sept 23, 2020

Due date: Oct 16, 2020 (Friday) 11:59pm

*The solution will be posted right after the deadline, so no late homework will be accepted!*

Every Student **MUST** include the following statement, together with his/her signature in the submitted homework.

*I declare that the assignment submitted on Elearning system is original except for source material explicitly acknowledged, and that the same or related material has not been previously submitted for another course. I also acknowledge that I am aware of University policy and regulations on honesty in academic work, and of the disciplinary guidelines and procedures applicable to breaches of such policy and regulations, as contained in the website*

<http://www.cuhk.edu.hk/policy/academichonesty/>.

Signed (Student Jesse) Date: 16-10-2020

Name Chan Kai Yin SID 1155126983

## Submission notice:

- Submit your homework via the elearning system

## General homework policies:

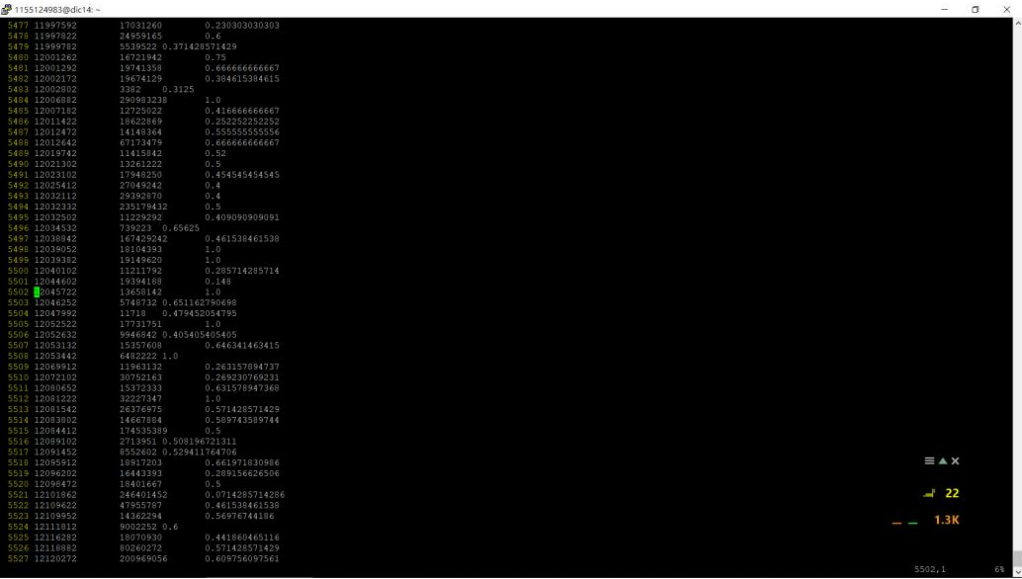
A student may discuss the problems with others. However, the work a student turns in must be created COMPLETELY by oneself ALONE. A student may not share ANY written work or pictures, nor may one copy answers from any source other than one's own brain.

Each student **MUST LIST** on the homework paper the **name of every person he/she has discussed or worked with**. If the answer includes content from any other source, the student **MUST STATE THE SOURCE**. Failure to do so is cheating and will result in sanctions. Copying answers from someone else is cheating even if one lists their name(s) on the homework.

If there is information you need to solve a problem but the information is not stated in the problem, try to find the data somewhere. If you cannot find it, state what data you need, make a reasonable estimate of its value, and justify any assumptions you make. You will be graded not only on whether your answer is correct, but also on whether you have done an intelligent analysis.

Recommend the person with the maximal number of common followees

For all users:



For the SID-version (1155124983):

```
File Input Format Counters
  Bytes Read=47199621
File Output Format Counters
  Bytes Written=34
20/10/13 14:37:28 INFO streaming.StreamJob: Output directory: md_la_sid_fin/
1155124983@dic14:~$ hadoop fs -cat md_la_sid_fin/part-* md_la_sid_fin.txt
245424983      59507002      0.489795918367
cat: 'md_la_sid_fin.txt': No such file or directory
1155124983@dic14:~$ hadoop fs -cat md_la_sid_fin/part-* > md_la_sid_fin.txt
1155124983@dic14:~$
```


For the format shown, the first column is the user. The second one is the person with the maximal number of common followees with the user. I also print the similarity in the third column as reference.

		Job Overview
Job Name:	streamjob1937540332724634808.jar	
User Name:	1155124983	
Queue:	default	
State:	SUCCEEDED	
Uberized:	false	
Submitted:	Mon Oct 12 16:39:16 HKT 2020	
Started:	Mon Oct 12 16:39:23 HKT 2020	
Finished:	Mon Oct 12 17:12:46 HKT 2020	
Elapsed:	33mins, 23sec	
Diagnostics:		
Average Map Time	10sec	
Average Shuffle Time	9mins, 12sec	
Average Merge Time	0sec	
Average Reduce Time	5mins, 22sec	

It takes about 33 minutes to run all user version.

		Job Overview
Job Name:	streamjob6326505746168795817.jar	
User Name:	1155124983	
Queue:	default	
State:	SUCCEEDED	
Uberized:	false	
Submitted:	Tue Oct 13 14:36:35 HKT 2020	
Started:	Tue Oct 13 14:36:43 HKT 2020	
Finished:	Tue Oct 13 14:37:26 HKT 2020	
Elapsed:	42sec	
Diagnostics:		
Average Map Time	8sec	
Average Shuffle Time	4sec	
Average Merge Time	0sec	
Average Reduce Time	16sec	

It takes about 5 minutes to run the SID-version.

 dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""sm_mapper.py"""
import sys
from collections import defaultdict

# input comes from STDIN (standard input)
a = defaultdict(set)
for line in sys.stdin:
    nums = line.strip().split()
    #print(nums[1],nums[0])
    print("%s\t%s"%(nums[1],nums[0]))

~
```

i Mapper

dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""md_mapper.py"""
import sys
from collections import defaultdict

# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "twitter_raw.txt"
with open(infile, 'rt') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
for line in sys.stdin:
    # for i in line.split()[1:]:
    a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():
    max_v = 0
    max_k = 0
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > max_v:
                max_v = new_v
                max_k = j
    print("%s\t%s\t%s" % (v, max_k, max_v))
~
~
~
~
```

ii Reducer for all

dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""md_mapper.py"""
import sys
from collections import defaultdict

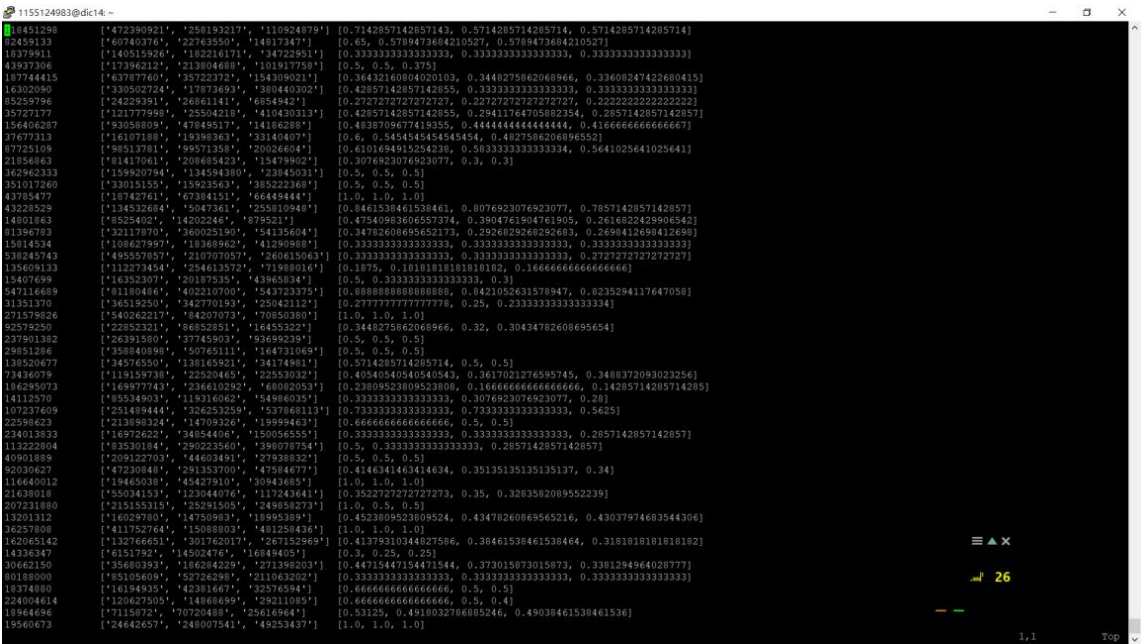
# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "twitter_raw.txt"
#set_a = [2, 4, 9, 8, 3]
#set_b = set()
with open(infile, 'rt') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
for line in sys.stdin:
    # for i in line.split()[1:]:
    a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():
    #set_b = set()
    #for j in list(str(v)):
    #    #set_b.add(int(j))
    #print(set_b)
    if not(str(v).endswith('24983')):
        continue
    max_v = 0
    max_k = 0
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > max_v:
                max_v = new_v
                max_k = j
    print("%s\t%s\t%s" % (v, max_k, max_v))
    #print(set_a , set_b)
~
~
~
~
```

iii Reducer for SID version

Most similar people of EVERY user for the medium-sized dataset

For the all user version:



1155124983@dic14~		
18451296	['472390921', '258193217', '110924879']	[0.7142857142857143, 0.5714285714285714, 0.5714285714285714]
82459133	['60740376', '22763550', '14817347']	[0.65, 0.5789473684210527, 0.5789473684210527]
18379911	['140515926', '182216171', '34722951']	[0.3333333333333333, 0.3333333333333333, 0.3333333333333333]
43937306	['17396212', '213804888', '101947758']	[0.5, 0.5, 0.378]
18774415	['63787760', '3522372', '154309021']	[0.36432160804020103, 0.346275862068966, 0.33688247422680415]
16302090	['330502724', '17873693', '380440302']	[0.42857142857142855, 0.3333333333333333, 0.3333333333333333]
85259796	['24225391', '26861141', '6854942']	[0.2727272727272727, 0.22727272727272727, 0.2222222222222222]
35727177	['12177998', '2550421', '410400233']	[0.42857142857142855, 0.29411764705882354, 0.2857142857142857]
156406287	['93058809', '47849517', '14186288']	[0.4838709677419355, 0.4444444444444444, 0.4166666666666667]
37677313	['16107188', '19398363', '33140407']	[0.6, 0.5454545454545454, 0.4827586206896552]
87725109	['98513781', '99571358', '20026604']	[0.6101694915254238, 0.5833333333333334, 0.5641025641025641]
21354963	['61417661', '208685423', '15478992']	[0.3076923076923077, 0.3, 0.3]
362962333	['15920794', '134594380', '23845031']	[0.5, 0.5, 0.5]
351017260	['33015155', '15923563', '385222368']	[0.5, 0.5, 0.5]
437825477	['18742761', '97384151', '66449444']	[1.0, 1.0, 1.0]
42226529	['13452684', '5047321', '255811048']	[0.4461384615384615, 0.8076923076923077, 0.7857142857142857]
14801863	['9525402', '14202246', '879521']	[0.47540983606557374, 0.3904761904761905, 0.2616822429906542]
81396783	['32117870', '360025190', '54135604']	[0.34782608695652173, 0.2926829268292683, 0.2698412698412698]
15814534	['108627997', '18368962', '41280988']	[0.3333333333333333, 0.3333333333333333, 0.3333333333333333]
338245743	['486557897', '21070705', '246815064']	[0.3333333333333333, 0.3333333333333333, 0.2727272727272727]
135609133	['112273454', '254613572', '71988016']	[0.1875, 0.18181818181818182, 0.16666666666666666]
15407699	['16325307', '20187335', '43965834']	[0.5, 0.3333333333333333, 0.3]
347116489	['81180488', '402210700', '54372375']	[0.8888888888888888, 0.8421052631578947, 0.8235294117647058]
31311376	['36159505', '542770193', '25042112']	[0.2777777777777778, 0.25, 0.23333333333333334]
27157926	['540262217', '84207073', '70850380']	[1.0, 1.0, 1.0]
92579250	['22852321', '86852851', '16455322']	[0.3448275862068966, 0.32, 0.30434782608695654]
237801382	['26391580', '37745993', '93689239']	[0.5, 0.5, 0.5]
28931256	['35840888', '50745111', '164731693']	[0.5, 0.5, 0.5]
138520677	['34576550', '138165921', '34174981']	[0.5714285714285714, 0.5, 0.5]
73436079	['119159738', '22520465', '22553032']	[0.40540540540540543, 0.3617021276595745, 0.3488372093023256]
186295073	['169977743', '236610292', '68082053']	[0.23809523809523808, 0.16666666666666666, 0.14285714285714285]
41112570	['85534903', '119316062', '54986035']	[0.3333333333333333, 0.2076923076923077, 0.28]
107237609	['251489444', '326253258', '537868113']	[0.7333333333333333, 0.7333333333333333, 0.5625]
22598623	['213898324', '14709326', '19999463']	[0.6666666666666666, 0.5, 0.5]
234013833	['16972622', '34854406', '150056555']	[0.3333333333333333, 0.3333333333333333, 0.2857142857142857]
113222504	['83630184', '290223460', '398097954']	[0.5, 0.3333333333333333, 0.2857142857142857]
40901889	['209122703', '44603491', '27938832']	[0.5, 0.5, 0.5]
92030627	['47230848', '291353700', '47584677']	[0.4146341463414634, 0.35135135135135137, 0.34]
116640012	['18465038', '45427910', '30943685']	[1.0, 1.0, 1.0]
21630018	['55034153', '123844076', '117243641']	[0.38272727272727273, 0.35, 0.3283582009552239]
207231880	['21515531', '25291505', '249858273']	[1.0, 0.5, 0.5]
13201312	['16029780', '14759983', '18995389']	[0.4523809523809524, 0.43478260869565216, 0.43037974683544306]
36257808	['411752764', '15088803', '481258436']	[1.0, 1.0, 1.0]
162065142	['132766651', '301762017', '267152969']	[0.41378310344897586, 0.38461538461538464, 0.3181818181818182]
14336347	['61517892', '14502476', '16849405']	[0.3, 0.25, 0.25]
30662150	['35860393', '186284229', '27138203']	[0.44715447154471544, 0.373015873015873, 0.3381294964028777]
80188090	['85105609', '52726296', '211063202']	[0.3333333333333333, 0.3333333333333333, 0.3333333333333333]
18374880	['16194935', '42381667', '32376594']	[0.6666666666666666, 0.5, 0.5]
224004014	['126627505', '14868699', '28211085']	[0.6666666666666666, 0.5, 0.4]
18964696	['7115872', '70720488', '25616964']	[0.53125, 0.4918032786885246, 0.49038461538461536]
19560673	['24642657', '248007541', '45253437']	[1.0, 1.0, 1.0]

Job Overview	
Job Name:	streamjob926582730120096016.jar
User Name:	1155124983
Queue:	default
State:	SUCCEEDED
Uberized:	false
Submitted:	Mon Oct 12 12:47:11 HKT 2020
Started:	Mon Oct 12 12:47:16 HKT 2020
Finished:	Mon Oct 12 13:22:53 HKT 2020
Elapsed:	35mins, 36sec
Diagnostics:	
Average Map Time	19sec
Average Shuffle Time	5mins, 24sec
Average Merge Time	0sec
Average Reduce Time	19mins, 30sec

It takes about 35 mins.


For the SID-version (1155124983) :

```
WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=47932987
File Output Format Counters
  Bytes Written=110
20/10/13 14:44:33 INFO streaming.StreamJob: Output directory: md_1b_sid_fin/
1155124983@dic14:~$ hadoop fs -cat md_1b_sid_fin/part-*
245424983      ['59507002', '37564255', '82963512']      [0.4897959183673469, 0.46099290780141844, 0.443609022
95639095]
1155124983@dic14:~$
```

Job Name: streamjob1343057855736599733.jar		Job Overview
User Name: 1155124983		
Queue: default		
State: SUCCEEDED		
Uberized: false		
Submitted: Tue Oct 13 14:43:54 HKT 2020		
Started: Tue Oct 13 14:43:59 HKT 2020		
Finished: Tue Oct 13 14:44:31 HKT 2020		
Elapsed: 32sec		
Diagnostics:		
Average Map Time 7sec		
Average Shuffle Time 5sec		
Average Merge Time 0sec		
Average Reduce Time 12sec		

It takes about 32 secs.

For the format shown, the first column Is the user. The second one is the persons with the Top3 maximal number of common followees with the user. I also print the similaritys for the each followees in the third column as reference.

 dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""sm_mapper.py"""
import sys
from collections import defaultdict

# input comes from STDIN (standard input)
a = defaultdict(set)
for line in sys.stdin:
    nums = line.strip().split()
    #print(nums[1],nums[0])
    print("%s\t%s"%(nums[1],nums[0]))

~
```

iv mapper

```

dic14.ie.cuhk.edu.hk - PuTTY
~/usr/bin/env python
"""md_mapper.py"""
import sys
from collections import defaultdict
from collections import Counter

# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "twitter_raw.txt"
with open(infile, 'rt') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
#set_b = {"0":0,"1":0,"2":0}
set_b_k = []
for line in sys.stdin:
    # for i in line.split()[1:]:
    a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():
    max_v = 0
    max_k = 0
    set_b = {"0":0,"1":0,"2":0}
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > min(set_b.values()):
                set_b[str(j)] = new_v
                #c = Counter(set_b)
                #set_b = c.most_common(3)
                #set_b = sorted([(x, i) for (i, x) in enumerate(set_b)], reverse=True)[:3]
                #if new_v in set_b:
                #    set_b_k.append(v)
    c = Counter(set_b)
    mc = c.most_common(3)
    print("%s\t%s\t%s" % (v, [key for key, val in mc], [val for key, val in mc]))
~
~
~
~
~

```

v reducer for all

dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""md_mapper.py"""
import sys
from collections import defaultdict
from collections import Counter

# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "twitter_raw.txt"
#set_a = {2, 4, 9, 8, 3}
#set_b = set()
with open(infile, 'rt') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
#set_b = {"0":0,"1":0,"2":0}
set_b_k = []
for line in sys.stdin:
    # for i in line.split()[1:]:
        a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():

    # set_c = set()
    #for j in list(str(v)):
        #set_c.add(int(j))
    #print(set_b)
    #if len(set_a.union(set_c)) != 5:
        #continue
    if not(str(v).endswith('24983')):
        continue
    max_v = 0
    max_k = 0
    set_b = {"0":0,"1":0,"2":0}
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > min(set_b.values()):
                set_b[str(j)] = new_v
                #c = Counter(set_b)
                #set_b = c.most_common(3)
                #set_b = sorted([(x, i) for (i, x) in enumerate(set_b)], reverse=True)[:3]
                #if new_v in set_b:
                    #set_b_k.append(v)

    c = Counter(set_b)
    mc = c.most_common(3)
    print("%s\t%s\t%s" % (v, [key for key, val in mc], [val for key, val in mc]))

~
"md reducer k sid.py" 50L, 1453C
```

vi reducer for sid



Common followees shared between A and its similar users

All user version:

[illegible]

Job Overview	
Job Name:	streamjob4588403187632686413.jar
User Name:	1155124983
Queue:	default
State:	SUCCEEDED
Uberized:	false
Submitted:	Mon Oct 12 15:54:43 HKT 2020
Started:	Mon Oct 12 15:55:19 HKT 2020
Finished:	Mon Oct 12 16:25:44 HKT 2020
Elapsed:	30mins, 24sec
Diagnostics:	
Average Map Time	9sec
Average Shuffle Time	4mins, 15sec
Average Merge Time	0sec
Average Reduce Time	17mins, 57sec

It takes about 30mins to run.

SID-version (1155124983) :

```

1155124983dd1c14-
Gc time elapsed (ms)=24444
CPU time spent (ms)=789050
Physical memory (bytes) snapshot=76550922240
Virtual memory (bytes) snapshot=43954045376
Total committed heap usage (bytes)=7209083936

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=46361139
File Output Format Counters
Bytes Written=2338

20/10/13 16:08:21 INFO streaming.StreamJob: Output directory: md_lc_output_1/
1155124983dd1c14:5 hadoop fs -cat md_lc_output_1/part-0
245424983158507002 ['115114748', '185397767', '79011908', '16405372', '2944
7358', '137064188', '16042649', '22798591', '1335662', '1910408', '19253485',
'19164725', '25440543', '7458942', '116324294', '15941338', '234903540', '207331
730', '123565032', '159409151', '59010637', '21575175', '16679681', '282280747',
'134234059', '36750478', '241203719', '59238436', '82563512', '121537720', '1414
06457', '109126247', '12360922', '112428240', '20437704', '47966829', '38500966',
'33235204', '271410124', '35415477', '15485844', '14642816', '35777372', '66489
082', '26299078', '227032210', '37564255', '133875756', '59166951', '15335662',
'11591138', '29725730', '241371677', '19780462', '15265186', '15989191', '62417
275', '16028750', '12360922', '11594312', '47248844', '22702507', '115303077',
'148368403', '22060843', '259286437', '14416109', '59110408', '28676476', '16012
222', '56767953', '101569095'] check sum = 5802276903
245424983137564255 ['115114748', '185397767', '79011908', '16405372', '2944
7358', '137064188', '16042649', '22798591', '1335662', '1910408', '19253485', '2646403',
'19164725', '59110488', '336644583', '7458942', '116324294', '15941338', '23490
3540', '123565032', '159409151', '59010637', '21575175', '16679681', '282280747',
'2947358', '36750478', '241203719', '59238436', '82563512', '121537720', '141488457', '2091
3636', '173428240', '20437704', '47966829', '38500966', '271410124', '35415477',
'15485844', '14642816', '22798591', '11342502', '58486467', '26299078', '133875
756', '119235720', '59166951', '16405372', '115911638', '29725730', '19780462',
'15265186', '15989191', '62417275', '16028750', '12360922', '11594312', '47348844', '2270825
07', '25440543', '59507002', '22060843', '151186578', '14416109', '28676476', '1
6012222', '56767953'] check sum = 4806343331
245424983182963512 ['115114748', '185397767', '79011908', '16405372', '1342
4058', '137064188', '22798591', '11535662', '1910408', '19253485', '2646403',
'19164725', '59110488', '7458942', '116324294', '159409151', '30913636', '1235650
32', '21575175', '21685295', '282280747', '2947358', '36750478', '241203719',
'121537720', '59166951', '141488457', '15594313', '47966829', '15485844', '385009
66', '33235204', '271410124', '224903540', '25777272', '64849082', '26299078',
'11186578', '133875756', '19265186', '5807988', '115911638', '29725730', '197804
62', '59238436', '15989191', '62417275', '16028750', '95406467', '25440543', '15
941338', '115303077', '59507002', '271410124', '259286437', '14416109', '2168033
66', '46409236', '101569095'] check sum = 496116248

1155124983dd1c14-5

```





Time consumption for each MapReduce job and its tasks.

Mapper num	Reducer num	Max mapper time	Min mapper time	Avg mapper time	Max reducer time	Min reducer time	Avg reducer time	Total job
20	10	1mins, 24sec	5sec	13sec	1hrs, 13mins, 34sec	56mins, 53sec	63 mins, 19sec	1hrs, 13mins, 45sec
10	20	13sec	4sec	10sec	36mins, 28sec	24mins, 59sec	31mins, 13sec	36mins, 45sec
30	50	18sec	4sec	8sec	16mins, 0sec	9mins, 59sec	14mins, 01sec	16mins, 12sec
50	30	16sec	2sec	7sec	28mins, 27sec	17mins, 41sec	21mins, 36sec	31mins, 39sec
60	60	1mins, 1sec	3sec	6sec	14mins, 2sec	8mins, 10sec	11mins, 25sec	41mins, 8sec
60	100	15sec	3sec	5sec	13mins, 34sec	4mins, 40sec	9mins, 14sec	40mins, 14sec

From the first 4 cases, the observation is that when # of reduce task > # of M reduce task will perform better and faster than vice versa version. Like 30-50 case only take about 16min, while 50-30 case almost take about twice of the time. Because the mapper output is spread across files in reducer, it is better for have more reducer than mapper.

And from the 1-6 cases, observation is the more mappers and reducers will make the avg. time of map and reduce task shorter.

But in the total time needed, more reducer and mapper task doesn't mean

better performance and speed. Like the fastest case in here is 30-50 case which only take about 16 mins, while 60-100 case takes about 100 mins. But perhaps it due to the heavy-workload and congestion of the system.

Find the TOP 3 ( $=K$ ) most similar people and the list of common followees for each user in the large dataset

For all users:

```
dic1414m3kxhukuh - PUTTY
11707275158747899459 ['10266545720184613689', '102029481403199559431', '103012084794102991018',
[0.25, 0.2, 0.2]
11171725363401423301 ['109312807119102647620', '11047540483309496739', '114744484073167880096',
[0.19230769230769232, 0.19117647058823529, 0.183900453970151]
1086689483444444245 ['1071994943219236230', '111045500143120127028', '1078467978731550642',
[0.4703703703703704, 0.43566464730290457, 0.4227941176470588]
11568309124486078076 ['112693170075201231415', '106207308906395800138', '106627210821625707000',
[0.37142857142857144, 0.2601626016260163, 0.25742574257425743]
1057159824042437759 ['107506870214464849', '112831359102447219603', '10226525655638010669',
[0.5777777777777777, 0.5306122448799592, 0.5185185185185185]
104361087486032569763 ['10168882213693965938', '104776183384905805748', '103482300621829758413',
[0.625, 0.5866666666666667, 0.541764705882353]
10490839721319130413 ['101138847601603602324', '103789101597902553217', '10036133655008061595',
[0.7818181818181819, 0.6610169491525424, 0.625]
1063513132285893699 ['11683296893573081132', '117923491423548471728', '11307167044043132974',
[0.25, 0.2, 0.18474619047619047]
1004661831319542459 ['10035251482668707250', '117278008765287520310', '100647247127009584190',
[0.633711340206185, 0.6477777777777778, 0.6470588235294118]
107786754120082323579 ['10167235788129482806', '112533250046302544371', '108019940100292441619',
[0.6306306306306307, 0.3908205252525252, 0.44436776048327668]
108175807481716626959 ['11839745767777317230', '10962942721319587159', '10343408925929861530',
[0.5714285714285714, 0.48275820686552, 0.42857142857142855]
104633353729476682346 ['101391281097346539710', '104786643585334861389', '114824780965740949647',
[0.6876666666666667, 0.597014925371343, 0.59659313004927]
11287162612521888143 ['111291864939851009207', '10894398664068380361', '1148041009934342743139',
[0.9411764705882353, 0.9411764705882353, 0.9411764705882353]
11422117575255368133 ['10679530474650056180', '107863686874402607929', '102294632541337027483',
[0.72329411764705882, 0.72329411764705882, 0.724268131594203]
112403137998583609131 ['11582129711457452448', '109649264094810211637', '114307254741024170258',
[0.23076923076923078, 0.23076923076923078, 0.18181818181818182]
11189558641602495569 ['1117212310517282832', '10120159611555213766', '11361910244480539592',
[0.19494949494949494, 0.10355555555555556, 0.08090909090909091]
10441882145591517246 ['11149687376196355326', '101499648759020119246', '1079397270578115018',
[0.1, 0.1, 0.1]
11242327545395344567 ['1153195577929363036', '105973253613548304461', '10618925012129550267',
[0.4, 0.3333333333333333, 0.3333333333333333]
10103433481075262228 ['110566643012319645135', '117930040055788289648', '11336501638198031804',
[0.40625, 0.4, 0.3593939393939393]
101801628277673976815 ['1115614547531522624015', '11557292606666198391', '117038500119265159736',
[0.1, 0.1, 0.1]
112246484101605991560 ['1132326615955555564242', '115394891386172731981', '10751720452410447474',
[0.1, 0.1, 0.1]
118001151658936861 ['10357406431963313088', '103952720945151610571', '101595419102735936432',
[0.1, 0.1, 0.1]
109040812548292539623 ['109913979662705246064', '108368853950551763047', '117271725376878070806',
[0.925, 0.9, 0.9]
1167722729151928447 ['117067539494881616527', '11669473210680381214', '1065580611437905645',
[0.6, 0.5806451612903226, 0.5714285714285714]
108104819156520935235 ['101524867499560135651', '10286439866231596579', '107063216606152401508',
[0.3333333333333333, 0.25, 0.2]
1174646788914968237 ['106240422862574287', '10440343546673630932', '114471618006333618271',
[0.9230769230769231, 0.6923076923076923, 0.6923076923076923]
```

For SID-version:

```

1041839405624039000001110835757750485679006 ['100163594096411699235', '1151526207793921795
101', '1178656314596727319', '100137120246450322354', '116455136667523112865', '1023381728449
69643494', '11014899471730462081', '100142984949036019136', '10164673244383912044', '1084004
1017099', '10452471413371515434', '10916145899293538208', '11653763444778032347', '1113901
11598152106969719269', '101083761367841810938', '1034394984303503915', '10436737829574616
101', '1050176758224682287', '107037649543687039893', '10553553214800234996', '111323916055941
8001244', '10209566989414388681', '11424281514142428932', '1163539434584510955', '10243705
11288351937073782140', '112864723042274554367', '1032467000244078824', '101
451723888075438708', '104010226630424675198', '108458224027000515910', '11355925179536825456
101', '10841775552475147385', '118276409840495259555', '104377821524346930092', '10729631454914
944984', '10112533366540715832', '105976191766074046677', '11785742785936544082', '1181910508
202494', '1056566940562935987', '10172340595902293917', '10020513240573822783', '101
867919101028045823', '101061733221561379250', '10332300833869011845', '1113070573630183757',
10597960948707087011', '117063357728895318461', '117870805029808655286', '101635681228985
21258', '101626577406833986387', '10576004874226038522', '10247567850174188605', '1116289702
171791833', '1013653274607048781', '11523812121007153117', '10234764602848097941', '100
3857646806028928', '10432756694647954178', '105998911565909100866', '102502383230554191628',
1099882023118687373071' check_sum = 65417836477618494257901
10418394056240390000111334798257177049200 ['100163594096411699235', '10164673244383912044
101', '1178656314596727319', '100644752379382846280', '109761659909253596200', '1023381728449
69643494', '11014899471730462081', '10112533366540715832', '111785742785936544082', '1020228
76706453290480', '1084004482007489874709', '11226545340461469145', '105442373415337158364', '11
11835925179536825456', '1163763444778032347', '101706255042693790681', '10013172004645032235
101', '101083761367841818362', '104717569501741189805', '1034673782957461614', '103858087687946
1587976', '1050107658224682287', '107037649543687039893', '1079691818994812988', '11203929
775340861640', '115152620779392179501', '10453498940305039915', '11583673584987391638', '102
10204425458', '1024703678545814735379', '11645136667523112865', '102338172844969643494',
'102346700023478932', '11355925179536825456', '104377821524346930092', '115228458623032
241233', '107229631454914944984', '105947063994736596850', '105976191766074046677', '105691729
617188160933', '1147490373171729301', '11424281514142428932', '1131356278416318857913', '10
1070676763640', '1172340595902293917', '11746870101028045823', '11522273240573822783', '101
1193230083386918845', '11317057376730183757', '1097960948707087011', '10597960948707087011',
55286', '100238746026920057941', '101626577406833986387', '100576009470226039522', '1013655277
46070426475', '11642487698141688437', '115232121007153117', '10506580408405235005', '1115
915219846915269', '105998911565909100866', '1025023832305541916261' check_sum = 60404573783912
72888832
104183940562403900000111334798257177049200 ['100163594096411699235', '10164673244383912044
101', '1178616599253596200', '100137120246450322354', '11645136667523112865', '10452240270
901585195', '10338172844969643494', '100644752379382846280', '115152620779392179501', '1113901
7263711696935', '1026467009470226039522', '11657363446778032347', '11131343895524406647',
10170625504269379681', '100644752379382846280', '1034394984303503915', '10436737829574616
101', '1050176758224682287', '107037649543687039893', '10553553214800234996', '11723391605
295771', '111333616058818001244', '10209566989414388681', '1020228767854350400', '11178574
278583654082', '102437056957454537375', '11288351937073782140', '113564728642274554367', '10
3426700023478932', '10723888075438708', '10401022663042675198', '106655673108390131524
10555179536825456', '10841775552475147385', '113204103382100530984', '10729631454914
944984', '105976191766074046677', '10369372678787562886', '1056912967188160535', '11160173
2219613797250', '11478490373171729301', '118191050832081858394', '1121356278416318857913', '116
292962086969736383', '1165363844545810955', '1030994069580497144', '10020513240573822783',
10170999', '11392273240573822783', '11932273240573822783', '110230633869011845', '1
```

Can't find the one that match my Sid last digit, so I use "00001" as the last 5 digit to match. The format is the same as 1b.

dic14.ie.cuhk.edu.hk - PuTTY

```
#!/usr/bin/env python
"""sm_mapper.py"""
import sys
from collections import defaultdict

# input comes from STDIN (standard input)
a = defaultdict(set)
for line in sys.stdin:
    nums = line.strip().split()
    #print(nums[1],nums[0])
    print("%s\t%s"%(nums[1],nums[0]))

~
~
~
~
```

x mapper

```
dic14.ie.cuhk.edu.hk - PuTTY
#!/usr/bin/env python
"""xi_mapper.py"""
import sys
from collections import defaultdict
from collections import Counter

# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "qpius_raw.txt"
with open(infile, 'r') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
set_b = ["*", "*", "0", "2":0]
set_b_k = []
for line in sys.stdin:
    # for i in line.split()[1:]:
    a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():
    max_v = 0
    max_k = 0
    set_b = ["*", "*", "0", "2":0]
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > min(set_b.values()):
                set_b[str(j)] = new_v
                #c = Counter(set_b)
                #set_b = c.most_common(3)
                #set_b = sorted([(k, v) for (k, v) in enumerate(set_b)], reverse=True)[0]
                #if new_v in set_b:
                #    set_b_k.append(v)

c = Counter(set_b)
mc = c.most_common(3)
print("%s\t%s\t%s" % (v, [key for key, val in mc], [val for key, val in mc]))

"qp_reducer_k.py" 39L, 1203C
```

xi reducer for all

```
dic14.ie.cuhk.edu.hk - PuTTY
#!/usr/bin/env python
"""xi_mapper.py"""
import sys
from collections import defaultdict
from collections import Counter

# input comes from STDIN (standard input)
b = defaultdict(set)
infile = "qpius_raw.txt"
set_a = [{"*", "*", "0", "2":0}]
set_b = set()
with open(infile, 'r') as f_i:
    for line in f_i.readlines():
        for i in line.strip().split()[1:]:
            b[str(line.strip().split(' ')[0])].add(i)
a = defaultdict(set)
set_b = ["*", "*", "0", "2":0]
set_b_k = []
for line in sys.stdin:
    # for i in line.split()[1:]:
    a[str(line.strip().split('\t')[0])].add(line.strip().split('\t')[1])

for v,k in a.items():
    if not(str(v).endswith('00001')):
        continue
    max_v = 0
    max_k = 0
    set_b = ["*", "*", "0", "2":0]
    for j,y in b.items():
        if len(k.union(y)) == 0 or v == j:
            pass
        else:
            new_v = float(len(k.intersection(y))/len(k.union(y)))
            if new_v > min(set_b.values()):
                set_b[str(j)] = new_v
                #c = Counter(set_b)
                #set_b = c.most_common(3)
                #set_b = sorted([(k, v) for (k, v) in enumerate(set_b)], reverse=True)[0]
                #if new_v in set_b:
                #    set_b_k.append(v)

c = Counter(set_b)
mc = c.most_common(3)

for key,val in mc:
    set_c = b[str(key)].intersection(k)
    set_sum = sum(map(int,list(set_c)))
    print("%s\t%s\t%s\t%s\t%s\t%s" % (v,key ,list(set_c), set_sum))
    print("[" + ",".join(str(i) for i in set_c) + "]")

"qp_reducer_k_sid.py" 49L, 1407C
```

xii reducer for sid

