### IERG 4300 Fall 2019 Homework #1

Release date: Sep 20, 2019

Due date: Oct 2, 2019 (Wed) 11:59am. (i.e. noon-time)

The solution will be posted soon after the deadline. 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_	Chim	) Date:	30-9-2019	
Name	Chim Ka Long	SID	1155094482	

#### 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.

#### 1155094482 Chim ka Long

#### The overall result:

Part A small dataset:

4	74,	606	155
2	249,	610	164
2	274,	414	164
6	8,4	14	165
4	14,	610	187
4	14,	599	196
2	288,	414	202
3	880,	414	207
4	14,	448	272
4	14,	474	298

#### Part A large dataset:

8811,14463	1055
56707,59269	1061
59269,68259	1070
30723,59269	1157
27468,59269	1218
7795,59269	1356
59269,67385	1429
19635,59269	1577
7795,30723	2152
30687,31327	2330

#### Part B small dataset:

```
      582
      <123,0.097826>
      <400,0.08046>
      <526,0.056604>

      182
      <64,0.059652>
      <608,0.059498>
      <606,0.058324>

      282
      <354,0.076087>
      <166,0.071839>
      <434,0.065217>

      82
      <354,0.07013>
      <239,0.053613>
      <570,0.050992>

      382
      <220,0.055679>
      <62,0.044521>
      <348,0.042553>

      482
      <363,0.040268>
      <200,0.0358>
      <376,0.033195>
```

#### Part B large dataset:

```
      4482
      <60423,0.110345>
      <6339,0.094059>
      <6420,0.092857>

      14482
      <37803,0.078845>
      <38340,0.076394>
      <70931,0.074074>

      24482
      <11967,0.078947>
      <38982,0.065217>
      <69424,0.065217>

      34482
      <1822,0.131313>
      <16834,0.130952>
      <60127,0.121212>

      44482
      <8835,0.32906>
      <54630,0.319249>
      <22169,0.313653>

      54482
      <42697,0.0666667>
      <50783,0.066381>
      <7286,0.056338>

      64482
      <28285,0.3>
      <64119,0.294574>
      <6984,0.286624>
```

#### **Explanation**

#### Part A:

1) Map reduce 0: mapper output everything, but use movie id as key, reducer aggregate the users with same movie in array. Then process the array to emit pair of users if their ratings are same. The output like below

```
270,570 1
270,590 1
132,270 1
270,307 1
270,288 1
270,305 1
217,270 1
214,270 1
```

Map reduce 1: mapper output everything directly. Reducer sum the number and put into array. If array exceed certain number, it will do sorting and only hold the top 10 same rating. Finally, print them out

```
474,606 155
249,610 164
274,414 164
68,414 165
414,610 187
414,599 196
288,414 202
380,414 207
414,448 272
414,474 298
```

#### Part B:

1) Map reduce 0: Different from part A, reducer only emit pair which include userid same as my last 2 digits. Also, it includes the number of same movies with same rating and the number of same movies regardless to rating. The output like below

```
159,382 0 1
249,382 1 1
246,382 1 1
177,382 1 1
210,382 0 1
```

Map reduce Count: We calculate the sum of movies of every user.

175	24
114	31
534	520
412	102
504	87
480	836

Map reduce 1: mapper only emit everything, and reducer will sum up the two numbers of same movies and calculate the similarity. Reducer need to read the sum of movies of every user firstly and according to user pair, calculate similarity =  $\{no. of same movies and same rating\}/\{no. of movies of user A + no. of movies of user B - no. of same movies\}$ . The output like below

```
82,354 0.070130
166,282 0.071839
282,354 0.076087
400,582 0.080460
123,582 0.097826
```

The command like below

```
Offile sum -file relation_mapperl.py -mapper relation_mapperl.py -file relation_reducerl.py -reducer relation_reducerl.py -input homeworki/small_rela_score_outO/* -output homeworki/small_ela_score_outO/* -output homeworki/small_ela_score_outO/* -output homeworki/small_rela_score_outO/* -output homeworki/small_ela_score_outO/* -output homeworki/small_ela_score_outD/* -output homeworki/small_ela_score_output homeworki/small_e
```

Map reduce 2: mapper filter the users of pair, if its last 2 digits is my cuid last 2 digits, emit it as Key and the partner userid + similarity as Value. Reducer will only print out the top 3 similar users with their id and similarity. Output like below

```
      582
      <123,0.097826>
      <400,0.08046>
      <526,0.056604>

      182
      <64,0.059652>
      <608,0.059498>
      <606,0.058324>

      282
      <354,0.076087>
      <166,0.071839>
      <434,0.065217>

      82
      <354,0.07013>
      <239,0.053613>
      <570,0.050992>

      382
      <220,0.055679>
      <62,0.044521>
      <348,0.042553>

      482
      <363,0.040268>
      <200,0.0358>
      <376,0.033195>
```

2) Handle large data set is almost same, except we filter with last 4 digits. Output:

```
      4482
      <60423,0.110345>
      <6339,0.094059>
      <6420,0.092857>

      14482
      <37803,0.078845>
      <38340,0.076394>
      <70931,0.074074>

      24482
      <11967,0.078947>
      <38982,0.065217>
      <69424,0.065217>

      34482
      <1822,0.131313>
      <16834,0.130952>
      <60127,0.121212>

      44482
      <8835,0.32906>
      <54630,0.319249>
      <22169,0.313653>

      54482
      <42697,0.0666667>
      <50783,0.066381>
      <7286,0.056338>

      64482
      <28285,0.3>
      <64119,0.294574>
      <6984,0.286624>
```

Part C:
Our map reduce job has 2 parts
Map reduce 0:

	Max.	Min.	Average.	Max.	Min.	Average.	Total
	Mapper	Mapper	Mapper	Reducer	Reducer	Reducer	job
	time	time	time	time	time	time	
5	8s	7s	7s	1h	1h	1h	1h
mappers				29mins	1mins	14mins	41mins
5							
reducers							
10	13s	5s	11s	59mins	25mins	43mins	1h 6s
mappers				55s	9s	25s	
10							
reducers							
20	7mins	4s	6s	42mins	11mins	21mins	42mins
mappers	5s			4s	29s	21s	18s
20							
reducers							

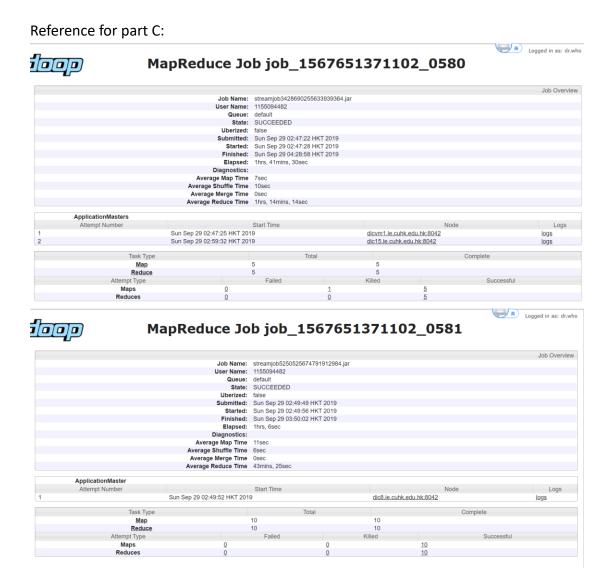
### Map reduce 1:

Although I only assign 5 mappers, the input file is too large, it is separated into many map tasks.

	Max.	Min.	Average.	Max.	Min.	Average.	Total
	Mapper	Mapper	Mapper	Reducer	Reducer	Reducer	job
	time	time	time	time	time	time	
5	19mins	13s	56s	2h	1h	1h	2h
mappers	39s			34mins	46mins	2mins	44mins
5							
reducers							
10	9mins	14s	48s	1h	57mins	25mins	1h
mappers	15sec			22mins	9s	24s	35mins
10							
reducers							
20	2mins	6s	46s	57mins	24mins	12mins	1h
mappers	2s			16s	36s	31s	6mins
20							
reducers							

Obviously, more mappers and reducers can make entire job faster. But it is not

linearly. For example, if the mappers and reducers increase become 2 times, the total job elapsed time will not become half. There are mappers or reducers which completing job slowly, because their hardware is worse than others, or their jobs is much difficult. For example, in map reducer1, if the movie was watched by many users, the number of pairs emitted will exponentially increase. The job will become more difficult.





### MapReduce Job job\_1567651371102\_0582

							Job Overvie
	Job Name:	streamjob8506691652159	702830.jar				
	User Name:	1155094482					
	Queue:	default					
		SUCCEEDED					
	Uberized:						
		Sun Sep 29 02:52:12 HKT					
		Sun Sep 29 02:52:25 HK1					
		Sun Sep 29 03:34:43 HK1	2019				
	Elapsed: Diagnostics: Average Map Time Average Shuffle Time						
	Average Merge Time						
	Average Reduce Time	21mins, 21sec					
ApplicationMaster							
Attempt Number		Start Time			Node		Logs
	Sun Sep 29 02:52:17 HKT 201	9		dic18.ie.cuhk.edu.l	nk:8042		<u>logs</u>
Task Type		Total			Comp	lete	
Map		20		20			
Reduce		20		20			
Attempt Type		Failed		Killed		Successful	
Maps	0		1		20		
Reduces	0		1		20		

Logged in as: dr.who

## 

### MapReduce Job job\_1567651371102\_0745

							Job Overv
	Job Name:	streamjob898037598	82397334413.jar				
		1155094482					
	Queue:	default					
	State:	SUCCEEDED					
	Uberized:	false					
	Submitted:	Sun Sep 29 14:05:04	4 HKT 2019				
	Started:	Sun Sep 29 14:05:25	5 HKT 2019				
	Finished:	Sun Sep 29 16:49:50	0 HKT 2019				
	2hrs, 44mins, 24sec						
	Diagnostics:						
	56sec						
	Average Shuffle Time						
	Average Merge Time	10sec					
	Average Reduce Time	1hrs, 2mins, 59sec					
ApplicationMaster							
Attempt Number		Start Time			Node		Logs
	Sun Sep 29 14:05:14 HKT 2019	9		dic19.ie.cuhk	c.edu.hk:8042		logs
Task Ty	/pe	Т	otal		(	Complete	
Map		915		915			
Reduc		5		5			
Attempt Type		Failed		Killed		Successful	
Maps	0		25		915		
Reduces	0		0		5		

Logged in as: dr.wh



### MapReduce Job job\_1567651371102\_0748

							Job Overv
	Job Name:	streamjob5609767518110	0085365.jar				
	User Name:	1155094482					
	Queue:	default					
	State:	SUCCEEDED					
	Uberized:	false					
	Submitted:	Sun Sep 29 14:06:59 HK	T 2019				
	Started:	Sun Sep 29 14:15:57 HK	T 2019				
	Finished:	Sun Sep 29 15:50:59 HK	T 2019				
	Elapsed:	1hrs, 35mins, 1sec					
	Diagnostics:						
	Average Map Time						
	Average Shuffle Time						
	Average Merge Time						
	Average Reduce Time	25mins, 24sec					
ApplicationMaster							
Attempt Number		Start Time			Node		Logs
	Sun Sep 29 14:15:46 HKT 2019	9		dic6.ie.cuhk.e	edu.hk:8042		<u>logs</u>
Task Type		Total			C	Complete	
Map		918		918			
Reduce		10		10			
Attempt Type		Failed	K	illed		Successful	
Maps	0		0		918		
Reduces	0		0		10		



# MapReduce Job job\_1567651371102\_0750

							Job Over
	Job Name:	streamjob35690733200	37543332.jar				
	User Name:	1155094482					
	Queue:	default					
	State:	SUCCEEDED					
	Uberized:	false					
	Submitted:	Sun Sep 29 14:08:50 H	KT 2019				
	Started:	Sun Sep 29 14:28:35 H	KT 2019				
	Finished:	Sun Sep 29 15:35:11 H	KT 2019				
	Elapsed:	1hrs, 6mins, 36sec					
	Diagnostics:						
	Average Map Time	46sec					
	Average Shuffle Time						
	Average Merge Time	3sec					
	Average Reduce Time	12mins, 31sec					
pplicationMaster							
Attempt Number		Start Time			Node		Logs
	Sun Sep 29 14:28:24 HKT 2019	)		dic10.ie.cuhk	.edu.hk:8042		<u>logs</u>
Task Type		Tota			C	omplete	
Map		923		923		•	
Reduce		20		20			
Attempt Type		Failed		Killed		Successful	
Maps	<u>0</u>		0		923		
Reduces	0		0		20		