

CSE 482: Big Data Analysis (Fall 2017) Homework 4

Due date: Monday, December 12, 2017

Please make sure you submit the homework via the handin system.

1. Write the corresponding HDFS commands to perform the following tasks.

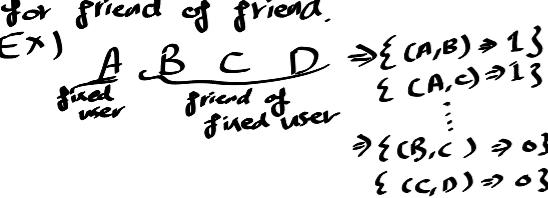
Each of these tasks must be accomplished with a single HDFS command.

Hint: type `hadoop fs -help` for the list of commands available. Note the difference between the local (Linux) file directory and HDFS directory. To double-check your answers, you're encouraged to test the commands (on your own HDFS directory) to make sure they work correctly.

- (a) Make a directory on HDFS named `/user/hadoop`. $\Rightarrow \text{hadoop fs -mkdir /user/hadoop}$
- (b) Upload a file named `data.txt` from the local filesystem to HDFS. $\text{hadoop fs -copyFromLocal data.txt}$
- (c) Rename the uploaded file on HDFS from `data.txt` to `temp.txt`. $\text{hadoop fs -mv data.txt temp.txt}$
- (d) Move the file from its current location at `/user/hadoop/temp.txt` to `/user/hduser/temp.txt`. $\text{hadoop fs -mv /user/hadoop/temp.txt /user/hduser/temp.txt}$
- (e) Copy the file from `/user/hduser/temp.txt` to `/user/hadoop/data.txt` (note that the filename is also changed). $\text{hadoop fs -cp /user/hduser/temp.txt /user/hadoop/data.txt}$
- (f) Display the content of the file `/user/hduser/temp.txt`. $\Rightarrow \text{hadoop fs -cat /user/hduser/temp.txt}$
- (g) Delete the file `/user/hadoop/data.txt`. $\Rightarrow \text{hadoop fs -rm /user/hadoop/data.txt}$
- (h) Download the file named `/user/hduser/temp.txt` to the local filesystem directory `/user/cse482`. $\Rightarrow \text{hadoop fs -copyToLocal /user/hduser/temp.txt /user/cse482}$
- (i) Merge all the files under the HDFS directory `/user/hadoop/results` into a single file `output.txt` to be stored in your current working directory. $\Rightarrow \text{hadoop fs -getmerge /user/hadoop/results output.txt}$
- (j) Changing permission of the HDFS directory `/user/hduser` so that only the user `hduser` can read and write files stored in the directory (but not other users).
 $\text{hadoop fs chmod -R 600 /user/hduser}$

b)

- **Mapper function:**
tokenize each user and set first row of first user is fixed user and last of user list are friends of fixed user. Then make key pair of friendship value 1 then make another pair with fixed friend users for friend of friend.

Ex) 

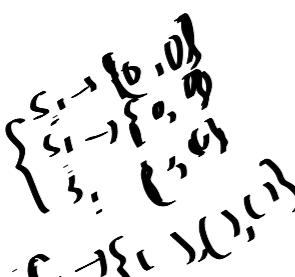
$$\begin{aligned} \Sigma(A,B) &= 1 \\ \Sigma(A,C) &= 1 \\ \Sigma(B,C) &= 0 \\ \Sigma(C,D) &= 0 \end{aligned}$$

- **Mapper output:** Key is tuple of friends relationship value is 1 and another tuple of friend of friend relationship value is 0.

- **Reducer input:** all list of tuples, value 1 and 0

- **Reducer function:** check the tuple is duplicate and if it is duplicate then drop.
Do not choose 1 value tuple because they are friend relationship. Then select only 0 value tuples.

- **Reducer output:**
tuple(fixedU, freqU) where fixedU < fixedU, key is fixedU
value: freqU
(0 value tuple was friend of friend relationship)



a)

- **Mapper function:**

tokenize each key information such as A123456074 / CSE260 / FS2013

then set key and value

- **Mapper output:** key is student ID, value: tuple type (credits, GPA)

- **Reducer input:** key is student ID, value is list of tuples (credits, GPA)

- **Reduce function:** multiply credits and GPA in each tuple and add up the multiply numbers, and add up just credits column and devide multiply add-up number.

- **Reduce output:**
Key = student ID, value = GPA Average.

- (a) **Data set:** Student transcript database. Each record consists the following information: student ID, course code, semester enrolled, number of credits, and GPA for the course. For example:

A123456074 CSE260 FS2013 3 3.5

GPA is a numeric-valued attribute ranging from 0.0 to 4.0.

Problem: Compute the average GPA for each student.

- (b) **Data set:** Facebook friendship graph. Each record corresponds to a Facebook user, followed by the list of his/her friends. For example, the graph data may contain the following records:

john123 mary456 tom312 lee222

mary456 john123

tom312 john123 lee222

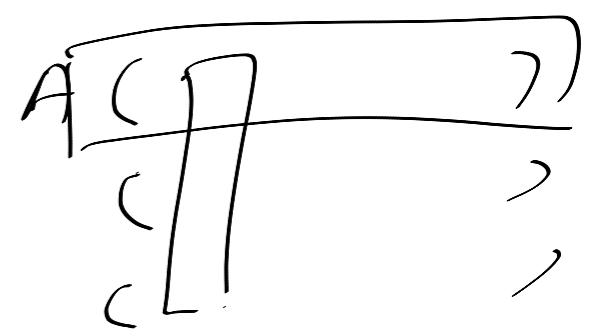
lee222 john123 tom312 → 2

Problem: Find pairs of users who are not friends with each other but share one or more common friends. This is known in network science literature as the "friend-of-a-friend" (FOF) problem. For example, mary456 and tom312 share a common friend, john123, but they are not friends with each other. In this case, the hadoop program should output the pair (mary456, tom312) only once, i.e., it should not output the duplicate pair (tom312, mary456).

2 () → 1
() → 0

c)

- mapper function: get each line and mapping the set with user_id and make tuple artist and vote
- mapper output: key is user_id and value is tuple (artist, vote)
- reducer input: key is user_id and value is list of tuple (artist, vote)
- reducer function:
select only 1 value from tuple [vote] position for each user
- reduce output:
Key: user_Id
Value: tuple (favorite artist, highest vote)



d)

- mapper function:
Get each line and mapping with each station_id and tuple with date, total_precipitation

- mapper output:
key: station_Id
value: tuple (date, total_precipitation)
- reducer Input:
key: station_Id
value: list of tuple (date, total_precipitation)
- reducer function:
Get mean value and standard deviations from each station set. Then using those values to calculate Z-score. And if Z-score is smaller than -3 or bigger than 3 means not normal. So drop the not normal value and also drop the negative total_precipitation value set.

- reducer output:
Key: station_Id
value: date (filtered normal Z-score and not negative total_precipitation value)

- reducer output

key: player1_ID

value: The ID for toughest opponents

(c) Data set: Online music streaming data. Each line in the data file has 4 columns (user_id, song_title, artist, vote), where vote is an integer value with the following values: +1 (the user liked the song), 0 (the user did not rate the song), -1 (the user disliked the song).

Problem: For each user, list his/her favorite artist, i.e., the artist that received the highest net vote. Do not output anything if the net vote is non-positive (0 or negative).

(d) Data set: Measurements of total daily precipitation for weather stations from around the world. Each line in the data files has 3 columns (station_id, date, total_precipitation).

Problem: Find the station_id and date of anomalous precipitation values. A precipitation value is considered anomalous if it is negative-valued or is 3 standard deviations or more away from its mean.

(e) Data set: Online 2-player gaming data. Each line in the data file has the following information:

gameID player1 player2 player1Score player2Score

For example, the line

123311 Wchampion lansing233 50 24

suggests that Wchampion won the game by a score of 50 to 24 (the winner is the player who achieves a score of 50 or higher). Note that a player can play against the same opponent more than once. Assume each player in the dataset has played more than 100 games.

Problem: For each player, find the toughest opponent to play against (i.e., the opponent to whom they have lowest winning percentage).

For example, if Wchampion has a winning percentage of only 0.20 when playing against msu4life and has a winning percentage higher than 0.2 when playing against other opponents, the program should output the pair (playerID, toughestOpponent):

Wchampion, msu4life

For tie-breaking cases (i.e., if a player has more than one toughest opponent with lowest winning percentage), the program may randomly select one of the toughest opponents as its output. If a player has won all of his/her games, then the program should output the pair (playerID, NoToughestOpponent).

(e) • mapper function: Get each line for games, select player1, player2 (opponents), player1 score, and player2 score, then compare player1 score and player2 score, then if player1 score is bigger than player2 score set win value "1", if is same or smaller then set win value as "0". Mapping with player1 and tuple of player2 and win value (0 or 1)

• mapper output: key player1 and values: tuple (player2, win value)

• reducer Input: key: player1 value: list of tuple (player2, win value)

• reducer function³: Count that how many game player1 done, and then add up all win value (0 or 1), it tell how many game player1 win from game. Then calculate winning percentage base on these value. Then compare each player's winning percentage and then sort order by the percentage and make pair with own player and higher winning percentage player, if it is same then randomly make pair.

~~Z-Score~~

~~El + m~~

(p1 p2 \square)
tuple (\square , 1)

p1	p2	1
p3	0	tot/100
p2	1	win 66
p4	1	0.5
⋮	⋮	⋮
		50