

# **CS550: Massive Data Mining and Learning**

## **Homework 1**

Due 11:59pm Sunday, February 25, 2018

Only one late period is allowed for this homework (11:59pm  
Monday 2/26)

# Submission Instructions

**Assignment Submission** Include a signed agreement to the Honor Code with this assignment. Assignments are due at 11:59pm. All students must submit their homework via Sakai. Students can typeset or scan their homework. Students also need to include their code in the final submission zip file. Put all the code for a single question into a single file.

**Late Day Policy** Each student will have a total of *two* free late days, and for each homework only one late day can be used. If a late day is used, the due date is 11:59pm on the next day.

**Honor Code** Students may discuss and work on homework problems in groups. This is encouraged. However, each student must write down their solutions independently to show they understand the solution well enough in order to reconstruct it by themselves. Students should clearly mention the names of all the other students who were part of their discussion group. Using code or solutions obtained from the web is considered an honor code violation. We check all the submissions for plagiarism. We take the honor code seriously and expect students to do the same.

Discussion Group (People with whom you discussed ideas used in your answers):  
Georgios Chantzialexiou

On-line or hardcopy documents used as part of your answers:

I acknowledge and accept the Honor Code.

(Signed) CM

If you are not printing this document out, please type your initials above.

## Answer to Question 1

The algorithm consists of two phases, the Map Phase and the Reduce Phase. I will go on to describe what happens in each phase.

### Map Phase

In this phase we emit  $\langle \text{key}, r, m \rangle$ . The key value will be the user id for the user that will get the recommendation, the second value (r) will be the user id of the user that will be recommended to the key user and the third value will be the id of the mutual friend. Since we do not want to recommend users that are already friends with the key user, we will set  $m=-1$  for all users that are already friends with key user. More specifically we emit the following:

```
for (user : Users)
    for (friend : user.friendsList)
        emit <user, friend, -1 >

for (user : Users)
    for (friend1 : user.friendsList)
        for (friend2 : user.friendList)
            emit <friend1, friend2, user>
```

*Note: These are all the possible combinations of friends in the user friends list.*

### Reduce Phase

We just sum up how many mutual friends they have been between the key and r users. If any of them has mutual friend -1 ( $m=-1$ ) they are already friends and we don't make that recommendation. Finally we recommend the 10 friends with the most mutual friends (or fewer if there are not that many).

## Answer to Question 2(a)

## Answer to Question 2(b)

## Answer to Question 2(c)

## Answer to Question 2(d)

## Answer to Question 2(e)



### Answer to Question 3(a)

### Answer to Question 3(b)

### Answer to Question 3(c)