

# Data Mining

## Homework 4

**Due:** 25/5/2014, 23:59.

You must hand in the homeworks electronically and before the due date and time. Check the web page for instructions about collaboration, about being late, and about handing in the homework.

### Hand in:

- The .py files with the source code.
- Instruction to execute them.
- The output files generated, and output at the screen.

**Problem 1.** In this assignment we will implement the  $k$ -means algorithm in Hadoop. The specifications of the program are to accept a file where each line corresponds to a document and a value  $k$  of the number of clusters, then select  $k$  random initial centers, run the  $k$ -means algorithm, and return a file where each line corresponds to a cluster.

In more detail, the format of the input file is:

```
fileIDi \t termi1:tfidfi1,termi2:tfidfi2...
```

where  $\text{fileID}_i$  is an ID (e.g., the filename) of document  $i$ ,  $\text{term}_{ij}$  the  $j$ th term of the  $i$ th document, and  $\text{tfidf}_{ij}$  the TFIDF value of the  $j$ th term in the  $i$ th document, after  $\ell_2$  normalization so that for every  $i$  we have that

$$\sum_{j=1}^n (\text{tfidf}_{ij})^2 = 1.$$

The main idea of the  $k$ -means algorithm in Hadoop should be almost obvious, and it is as follows:

1. First we select the  $k$  random centers using MapReduce. To do that we can assign a random number to each document and we select the top  $k$ .
2. The main part of the algorithm consists of several map–reduce rounds. In the map step, every document must select its closest center. In the reduce rounds, for every cluster we have a list of the documents assigned to it and we recompute the new center.
3. After convergence, for each cluster we report the documents assigned to it. We can also do this with MapReduce.

The various components of the algorithm are not too hard but you will need to combine them carefully. In particular, you will need three different MapReduce classes (one for choosing the random centers, one to perform an iteration of  $k$ -means, and one to report the final clusters) and a main program that controls them.

After you write and test your program, try to cluster some of the books of the Gutenberg project (<http://www.gutenberg.org>). Download the August 2003 CD:

[http://www.gutenberg.org/wiki/Gutenberg:The\\_CD\\_and\\_DVD\\_Project](http://www.gutenberg.org/wiki/Gutenberg:The_CD_and_DVD_Project)

The goal is to cluster all the .txt files in the CD. First preprocess the document. Create a single file and for each .txt file in the CD:

1. Remove the project Gutenberg header.

2. Remove stopwords, convert to lowercase and stem.
3. Store them in a single file in the format described above.

After you prepare the file, cluster it with  $k = 10, 20,$  and  $50$ . Try multiple times. Does it end up with the same cluster? Can you suggest a way to combine the different clusterings?