2018 Introduction to Massive Data Analysis Assignment 2

Deadline: 2018/11/16(Fri) 23:59

Questions: Finding Similar Articles

Given a set of BBCSports articles, please write a MapReduce program in Hadoop that can find out the most similar top 10 articles for the topic "athletics".

Your program should output 2 things: top 10 indices of candidate pairs and their similarities. The output format should look like this:

(index1, index2): similarity

According to Ch.3, 3 essential steps are needed for similar docs. Thus, this assignment should contain these 3 parts.

1. **Shingling**: Convert documents to sets.

The first part of this assignment is to implement the k-shingle, but this time we are going to use "words" instead of "characters".

You should read in all the articles of the topic, and then translate them into shingles. After this part, your output may look like this example, which depends on how you implement your own Shingling.

Note that you should use 5-shingles.

Documents

| Saliki di sa | 1 | 1 | 1 | 0 |
|--------------|---|---|---|---|
| | 1 | 1 | 0 | 1 |
| | 0 | 1 | 0 | 1 |
| | 0 | 0 | 0 | 1 |
| | 1 | 0 | 0 | 1 |
| | 1 | 1 | 1 | 0 |
| | 1 | 0 | 1 | 0 |

2. **Minhashing**: Covert large sets to short signatures, while preserving similarity.

The goal in this part is that you should hash your shingles into smaller size to get "signatures".

As we know, permuting rows even once is prohibitive because the cost is too high. However, there is a implementation trick.

As described in Ch.3 P.40, instead of permuting rows, you can simply generate different hash functions to simulate the permutation step according to the algorithm.

Note that you need 100 different hash functions.

3. Locality-sensitive hashing: Focus on pairs of signatures likely to be from similar documents.

Finally, we can apply LSH to the signatures. You should first partition your signature matrix M into b bands and r rows, where b = 20 and r = 5.

Then, for each band, you should hash its portion of each column to a hash table with k buckets.

The details are all in Ch.3 P48~55.

Note that candidate pairs are those that hash to the same bucket for ≥ 1 band.

After getting all the candidate pairs, you can calculate the Jaccard Simularity and then output the top 10.

Hint. MapReduce actually is a kind of hashing, so if you can design your mapper smartly, you can avoid some difficulties of the hashing function in this part.

Report Requirements:

a. Final output. (We require top 10 candidate pairs and their similarities)

NOTE: Please sorted results decreasingly.

- b. Explain how you design every mapper and reducer.
- c. Please make sure that your .java file has the same name as your class name, which must be LSH.
- d. Please hand in a .rar or .zip file, which includes your .pdf report and your java code.