

Lab 3 – Alternative solutions

Lab 3

Input file

```
A1008ULQSW1006,B0017OAIQI
A100EBHBG1GF5,B0013T5Y04
A1017Y05GBINVS,B0009F35AK
A101F8MBDPFDM9,B005HY2BRO,B000H7MFI
A1012H8HCJJAB,B0007A8XV6
A1012MEYm2YW2P5,B000FKGTBW
A102OP2OSXRHV,B001EQ4SGUJ,B000EHR7S
A101TGNH1D915Z,B000RHXKC6,B0002DHNXC,B0002DHNXC,B000XJK7UG,B0008DFK5,B000
SP1CWV,B0009YD7P2,B000SP1CWV,B0008DFK5,B0009YD7P2
A1012WAL0HJWH,B000W5U5H6
A1012V04GOA7RV,B002G3JY6,B001E5E3Y,B008ZRKZ5M,B002G3JW5
```

Each line contains

- a reviewerID (AXXXXXX) and
- the list of products reviewed by her/him (BXXXXXX)

Lab 3 – Ex. 1

- Your goal is to find the top 100 pairs of products most often reviewed (and so bought) together
- We consider two products as reviewed (i.e., bought) together if they appear in the same line of the input file

Lab 3 – Ex.1: Possible solutions

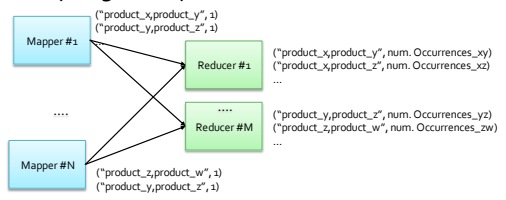
- At least three different “approaches” can be used to solve Ex. 1 of Lab 3

Lab 3 – Ex.1: Solution #1

- A chain of two MapReduce jobs is used
 - The first job computes the number of occurrences of each pair of products that occur together in at least one line of the input file
 - It is like a word count where each “word” is a pair of products
 - The second job, selects the top-k pairs of products, in terms of num. of occurrences, among the pairs emitted by the first job
 - It implements the top-k pattern

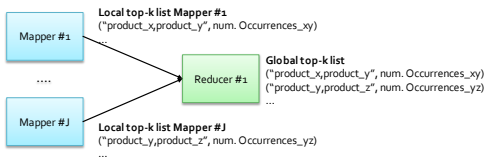
Lab 3 – Ex.1: Solution #1

- The first job computes the number of occurrences of each pair of products analyzing the input file



Lab 3 – Ex.1: Solution #1

- The second job computes the global top-k pairs of products in terms of num. of occurrences



7

Lab 3 – Ex.1: Solution #2

2. One single MapReduce jobs is used
 - The job
 - Computes the number of occurrences of each pair of products that occur together in at least one line of the input file
 - It is again like a word count where each "word" is a pair of products
 - However, the reducer does not emit all the pairs (pair of products, #of occurrences) that it computes
 - The top-k list is computed in the reducer and is emitted in its cleanup method

8

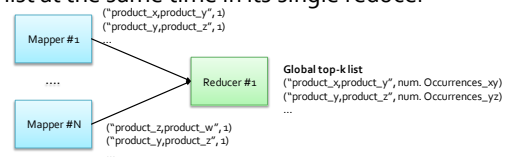
Lab 3 – Ex.1: Solution #2

- In the reducer, the job computes **also** the top-k list
 - By initializing the top-k list in the setup method of the reducer
 - By updating the top-k list in the reduce method (immediately after the computation of the frequency of the current pair of products)
 - By emitting the final top-k list in the cleanup method of the reducer
- There must be **one single reducer** in order to compute the final global top-k list

9

Lab 3 – Ex.1: Solution #2

- There is one single job that computes the number of occurrences and the global top-k list at the same time in its single reducer



10

Lab 3 – Ex.1: Solution #3

3. A chain of two MapReduce jobs is used
 - The first job is the same job used by Solution #2
 - However, in this case the number of reducers is set to a value greater than one
 - This setting allows parallelizing this intermediate step
 - Each reducer emits a local top-k list
 - The first job returns a number of local top-k lists equal to the number of reducers of the first job

11

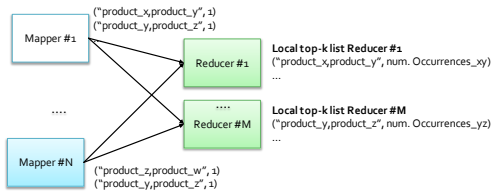
Lab 3 – Ex.1: Solution #3

- The second job computes the final top-k list merging the pairs of the local top-k lists emitted by the first job
 - It is based on the standard Top-k pattern

12

Lab 3 – Ex.1: Solution #3

- The first job computes the number of occurrences of each pair of products but each reducer emits only its local top-k pairs



33

Lab 3 – Ex.1: Solution #3

- The second job computes the global top-k pairs of products in terms of num. of occurrences merging the local list of job #1



34

Lab 3 – Ex.1: Comparison of the proposed solutions

- Solution #1
 - + Adopts two standard patterns
 - However, the output of the first job is very large
 - One pair for each pair of products occurring together at least one time in the input file

35

Lab 3 – Ex.1: Comparison of the proposed solutions

- Solution #2
 - + Only one job is instantiated and executed (there is only one job in Solution #2) and its output is already the final top-k list
 - However, only one reducer is instantiated
 - It could become a bottleneck because one single reducer must analyze the potentially large set of pairs emitted by the mappers
 - It is not a standard pattern

36

Lab 3 – Ex.1: Comparison of the proposed solutions

- Solution #3
 - + Each reducer of the first job emits only the pair contained in its local top-k lists
 - One top-k list for each reducer
 - The pairs of the top-k lists emitted by the reducers are significantly smaller than all the pairs of products occurring together at least one time
 - Since the first job instantiates many reducers, the parallelism is maintained
 - It is not a standard pattern

37