

Big Data**Exercises sheet n°1 : MapReduce algorithms and programming**

Important: at the end of the lab session you are required to send to dario.colazzo@polytechnique.edu the code of your Python MapReduce jobs, by indicating the related exercise for each job.

Exercise 1. Design a variant of the WordCount algorithm where the Map function performs local aggregation. The idea is to have a Map that emits couples of the kind (w, k) where $k \geq 1$ is the number of w occurrences in the input value.

Do you see any problem with this approach?

Is the combiner still needed?

Can Reduce be considered as a combiner as it is?

Observe the result of the Job, do you see any problems due to data cleaning issues?

Exercise 2. Design a MapReduce job without Combine that takes as input a file where each line contains a string indicating a Web URL and a natural number indicating the amount of time (in seconds) a user has spent on that Web page during a visit session. Of course you can have multiple lines for the same URL. The job is expected to return the list of unique URLs together with the average visit time.

Can the Reduce be considered as a Combiner?

If not, indicate the necessary modifications at the Map and Reduce level in order to have a Combine, and provide a definition of this last one.

Exercise 3. Provide a Python encoding of the following algorithms dealt with in classes: WordCount+ (Exercise 1), WordCount solving data cleaning issues, and Average calculation (Exercise2). To this end, rely on Hadoop streaming seen in the class.

Exercise 4. Encoding SQL queries in MapReduce. Consider the following simple relational schema containing informations about clients and orders they made.

Customer(cid, startDate, name)

Order(#cid, total)

Note that, for the sake of simplicity, we do not have any primary key for Order. Also assume that all the fields are mandatory.

Provide the Python MapReduce encoding of the following SQL queries.

- SELECT name FROM Customer WHERE month(startDate)=7
- SELECT DISTINCT name FROM Customer WHERE month(startDate)=7
- SELECT O.cid, SUM(total), COUNT(DISTINCT total) FROM Order O GROUP BY O.cid
- SELECT C.cid, O.total FROM Customer C, Order O WHERE C.name LIKE 'A%' and C.cid=O.cid
- SELECT C.cid, O.total FROM Customer C LEFT OUTER JOIN ON Order O ON C.cid=O.cid WHERE C.name LIKE 'A%' and C.cid=O.cid

For initial tests manually build simple text files containing a few lines for the input tables. Consider then the bigger files Customer and Order available on the Moodle.

Exercise 5. Graph analytics. Design and implement algorithms for the following two analytics problems on directed graphs. These problems are at the base of several analytics problems on graphs. Assume that a graph is a set of edge pairs (v,w) indicating that an edge exists from node v to node w .

1. **Universal sinks.** Given a directed graph, find the set of nodes having an incoming edge from all the remaining nodes, and having no outgoing edges.
2. **Triangles enumerations.** Given a directed graph, enumerate all directed triangles.

These two problems are fundamental problems in many application tasks related to graph analysis, and find applications to Internet and social network analysis, just to mention a few.

Pay attention to the fact that each of these problems may require at least two jobs in order to be solved. The second job takes as input the output of the first job.

In order to test your Python implementation first manually build very simple input graphs in the form of text files where each line models an edge (v, w) and contains a string for v followed by a `\t`, in turn followed by a string for w .

Then consider the text file graph.txt in the Moodle containing a bigger graph.

Pay attention to the fact that this file (obtained by a variant of the graph generator GTgraph) contains an edge table and the w string for (v,w) is repeated twice. Just ignore unneeded information in the Map phase.