

# Massive Data Processing - Lab 2

Emilie Leblanc, emilie.leblanc@student.ecp.fr

March 18, 2017

I am initially a student from ESSEC and I follow the MSc. in Data Sciences and Business Analytics (Centrale-ESSEC) as well as Centrale's Applied Mathematics (OMA) option. My git repository is here : [https://github.com/EmilieLeb/BDPA\\_Assign2\\_ELEBLAN](https://github.com/EmilieLeb/BDPA_Assign2_ELEBLAN) and there should be the commit tree inside.

## 1 Introduction - Hadoop setup

For my Hadoop setup, I have used a VirtualBox in which I have installed Ubuntu 16.04 and Java. I have also created a new user, hduser, on which I installed Hadoop 2.7.3.

## 2 Pre-processing the input

### 2.1 Do simple wordcount

I have started by doing a simple wordcount because it will be useful for us later. It also allows me to show the process I followed.

After formating the namenode and starting the dfs and the yarn, I checked that everything was up and running with jps.

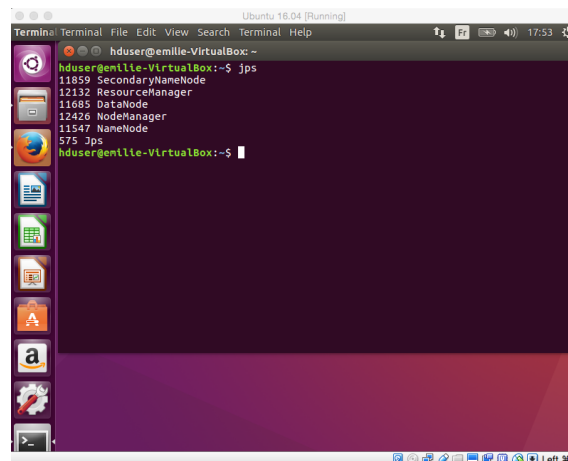


Figure 1: Jps command line.

Since everything was fine, I downloaded with curl that pg100.txt file from gutenber, as instructed in the assignment. I then put the inputdata in the HDFS.

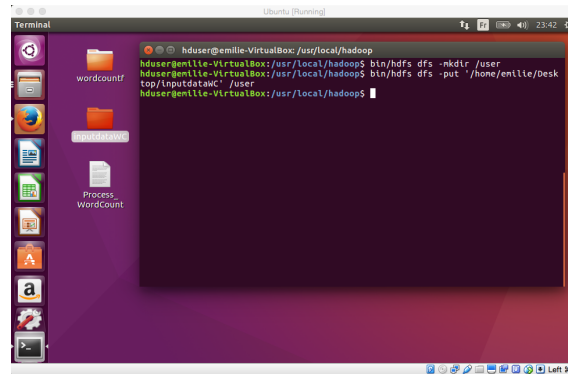


Figure 2: Inserting inputdata in HDFS.

On my Desktop, I had a file containing the WordCountpro.java code I wished to run.

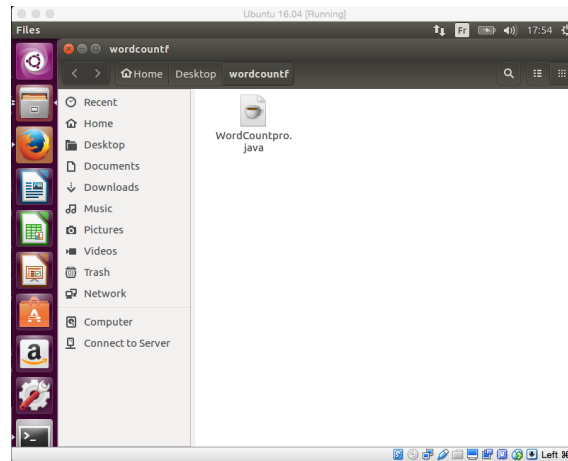


Figure 3: Wordcount initial folder containing code.

Then, I created the jar file and ran it by using the following commands (all of these commands, for every process, can be found in the Process file). And the job ran correctly.

```

1 #####
2 #### CREATE WORDCOUNT ####
3 #####
4
5 #Compile
6
7 javac -classpath /usr/local/hadoop/share/hadoop/common/hadoop-common
      -2.7.3.jar:/usr/local/hadoop/share/hadoop/mapreduce/hadoop-
      mapreduce-client-core-2.7.3.jar:/usr/local/hadoop/share/hadoop/
      common/lib/commons-cli-1.2.jar -d /home/emilie/Desktop/wordcountf
      *.java
8
9 #Convert into Jar File
10
11 jar -cvf wordcountj.jar -C /home/emilie/Desktop/wordcountf/wordcountc .
12
13 #Run JAR File
14
15 bin/hadoop jar /home/emilie/Desktop/wordcountf/wordcountj.jar
      WordCountpro /user/inputdata/pg100.txt output

```

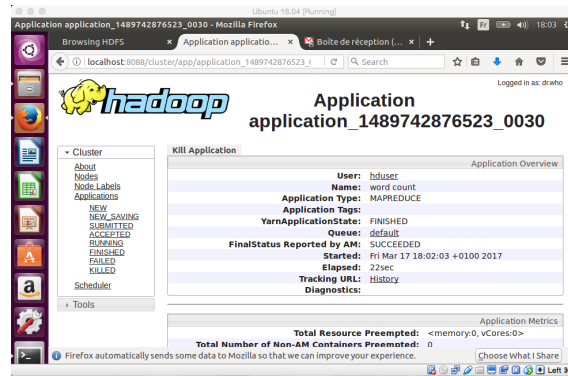


Figure 4: Success of the Wordcount job.

Extract of the obtained wordcount file:

```

1 address 19
2 addresses 2
3 addressing 1
4 address 2
5 adds 5
6 adhere 2
7 adheres 3
8 adieu 101

```

## 2.2 Create stopwords file

I proceeded the same way to create the stopwords file.

Extract of the obtained stopwords file:

```

1 a 14800
2 and 26847
3 as 5988
4 be 7146
5 but 6287
6 by 4490
7 d 8961

```

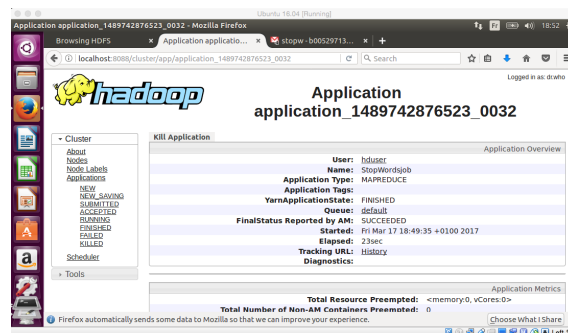


Figure 5: Success of the Stopwords job.

## 2.3 Preprocessing

In order to finalize the preprocessing, I proceeded in the following manner.

**Main function :**

```
1 public static void main(String[] args) throws Exception {
2
3     // Setting the configuration such that we can upload
4     // the documents in command line
5     Configuration conf = new Configuration();
6     conf.set("stopWords", args[2]);
7     conf.set("wordCount", args[3]);
8
9     Job job = Job.getInstance(conf, "PreProcessing");
10    job.setJarByClass(PreProcessingpro.class);
11    job.setMapperClass(Map.class);
12    job.setReducerClass(Reduce.class);
13
14    job.setOutputKeyClass(IntWritable.class);
15    job.setOutputValueClass(Text.class);
16
17    FileInputFormat.addInputPath(job, new Path(args[0]));
18    FileOutputFormat.setOutputPath(job, new Path(args[1]));
19
20    job.waitForCompletion(true);
21
22    // Problem with the number of lines: does not seem to
23    // work
24    System.out.println("Number of lines = " + totalLines);
25 }
```

**Mapper :**

```
1 public static class Map extends Mapper<Object, Text,
2     IntWritable, Text>{
3
4     // Initialise the stop words
5     private Text wordText = new Text();
6     private String word = new String();
7     String stopWords = "";
8
9     // Upload the stop words
10    protected void setup(Context context) throws
11        IOException, InterruptedException {
12
13        // Go get the file
14        String stopWordsPath = context.getConfiguration
15            ().get("stopWords");
16        Path openStopWords = new Path(stopWordsPath);
17        FileSystem fs = FileSystem.get(new
18            Configuration());
19        BufferedReader stopWordsFile = new
20            BufferedReader(new InputStreamReader(fs.
21                open(openStopWords)));
22        String line;
23        line= stopWordsFile.readLine();
24
25        // Upload each word one by one
26        try{
27            while (line != null){
28                stopWords += line + ",";
29            }
30        }
31    }
```

```

23         line = stopWordsFile.readLine()
24         ;
25     }
26     } finally {
27         stopWordsFile.close();
28     }
29
30     public void map(Object key, Text value, Context context
31         ) throws IOException, InterruptedException {
32         // Initialisation
33         String wordsoffline = value.toString();
34         String seenwords = "";
35         StringTokenizer token = new StringTokenizer(
36             wordsoffline, "
37             thesevalueswontshowcorrectlyinlatex");
38
39         // If word has not been seen before in the line
40         and is not a stop word, keep it
41         while (token.hasMoreTokens()) {
42             wordText.set(token.nextToken());
43             word = wordText.toString();
44             if (!stopWords.contains(word) && !
45                 seenwords.contains(word)){
46                 seenwords = seenwords + word + " ";
47             }
48         }
49         if (seenwords.length()>0) {
50             totalLines.set(totalLines.get()+1);
51             seenwords = seenwords.substring(0,
52                 seenwords.length()-1);
53             Text seenwordsText = new Text(seenwords
54             );
55             context.write(totalLines, seenwordsText
56             );
57         }
58     }
59 }

```

#### Reducer :

```

1 public static class Reduce extends Reducer<IntWritable,Text,IntWritable
2     ,Text> {
3
4     // Create a sort of dictionary to count the occurrences of
5     words to order them
6     private HashMap<String, Integer> dict = new HashMap<String,
7         Integer>();
8
9     // Import word count
10    protected void setup(Context context) throws IOException,
11        InterruptedException {
12
13        // Load file
14        String wordCountPath = context.getConfiguration().get("
15            wordCount");
16        Path openWordCount = new Path(wordCountPath);
17        FileSystem fs = FileSystem.get(new Configuration());
18        BufferedReader wordCountFile =new BufferedReader(new
19            InputStreamReader(fs.open(openWordCount)));
20        String line;

```

```

15         line = wordCountFile.readLine();
16
17         // For each line, get the word and the number of values
18         try{
19             while (line != null){
20                 String keyvalue[] = line.split("\t");
21                 String word = keyvalue[0];
22                 int nbcount = Integer.parseInt(keyvalue[1]);
23                 dict.put(word, nbcount);
24                 line = wordCountFile.readLine();
25             }
26             } finally {
27                 wordCountFile.close();
28             }
29     }
30
31     public void reduce(IntWritable key, Iterable<Text> values,
32                       Context context) throws IOException, InterruptedException {
33
34         for (Text value : values){
35
36             // Create a sort of sub dictionary per line associating
37             // word - count
38             HashMap<String, Integer> subdict = new HashMap<String,
39             Integer>();
40
41             // Take the words and counts of the associated line
42             String line = value.toString();
43             StringTokenizer token = new StringTokenizer(line);
44
45             while (token.hasMoreTokens()) {
46                 String word = token.nextToken();
47                 subdict.put(word.toLowerCase(), dict.get(word));
48             }
49
50             // Order the words of that line
51             List<String> HMKeys = new ArrayList<String>(subdict.
52             keySet());
53             List<Integer> HMValues = new ArrayList<Integer>(subdict
54             .values());
55             Collections.sort(HMValues);
56
57             // Use that ordered list
58             LinkedList<String> list = new LinkedList<String>();
59             Iterator<Integer> HMItvalue = HMValues.iterator();
60
61             while (HMItvalue.hasNext()) {
62
63                 Integer v = HMItvalue.next();
64                 Iterator<String> HMItkey = HMKeys.iterator();
65
66                 while (HMItkey.hasNext()) {
67
68                     String sortedWords = HMItkey.next();
69                     Integer compl = subdict.get(sortedWords);
70                     Integer comp2 = v;
71
72                     if (compl.equals(comp2)) {

```

```

68         HMIkey.remove();
69         list.add(sortedWords);
70     }
71 }
72 }

```

Noticeably, these settings allow us to directly set the source, the wordcount and the stopwords files in a simple command line, such as:

```

1 bin/hadoop jar /home/emilie/Desktop/preprocessingf/preprocessingj.jar
  PreProcessingpro /user/inputdata/pg100.txt output /user/inputdata/
  stopwords.txt /user/inputdata/wordcount.txt

```

As seen in the screenshot, the preprocessing job took 23sec.

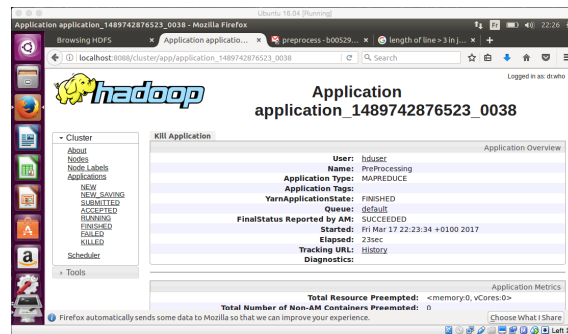


Figure 6: Preprocessing job success.

#### Extract of the preprocessing output file :

```

1 1      ebook complete works shakespeare william gutenber project the
2 2      shakespeare william
3 3      anyone anywhere ebook cost use this
4 4      restrictions whatsoever copy almost away give may you
5 5      included license re gutenber project terms under use
6 6      online gutenber www ebook org
7 7      copyrighted details below ebook gutenber project this
8 8      guidelines copyright file please follow
9 9      title complete works shakespeare william the
10 10     author shakespeare william

```

Ordering the words in the individual lines was a quite laborious task... Moreover, as you can see, I did not manage to extract the total number of lines.

### 3 Set-similarity joins

In this part, I only had time to do the pairwise comparisons with Jaccard similarity, so I did not implement the Inverted index. There also was not many space left on my VM so I limited the number of lines taken for the pairwise similarities to 50 and then to 1 000.

#### 3.1 PairWise comparisons

I proceeded as following.

**Main function :**

```
1      public static void main(String[] args) throws Exception {
2          // As usual..
3          Configuration conf = new Configuration();
4          Job job = Job.getInstance(conf, "PairWise");
5          job.setJarByClass(PairWisepro.class);
6          job.setMapperClass(Map.class);
7          job.setReducerClass(Reduce.class);
8          job.setOutputKeyClass(Text.class);
9          job.setOutputValueClass(Text.class);
10         FileInputFormat.addInputPath(job, new Path(args[0]));
11         FileOutputFormat.setOutputPath(job, new Path(args[1]));
12
13         job.waitForCompletion(true);
14     }
```

**Mapper :**

```
1     public static class Map extends Mapper<Object, Text, Text, Text>{
2
3         // Since we did not manage to find the exact line
4         // number, we set it arbitrarily (50, then 1000 and
5         // more if we can)
6         private int fullSize = 1000;
7
8         // Create another kind of counter to be sure we don't
9         // exceed the fullSize
10        private IntWritable lineKey = new IntWritable(0);
11
12        public void map(Object key, Text value, Context context
13            ) throws IOException, InterruptedException {
14            lineKey.set(lineKey.get()+1);
15
16            // For each line as long as we are under
17            // fullSize limit
18            if (lineKey.get() <= fullSize){
19
20                // We split the line between the line
21                // number and the string
22                String [] lineValueIt = value.toString
23                    ().split("\t");
24                String lineValue = lineValueIt[0];
25                int lineKey = Integer.parseInt(
26                    lineValueIt[0]);
27                String line = lineValueIt[1];
28                Text lineText = new Text(line);
29
30                // We then create all pairs of lines
31                // that we will use for the similarity
32                for (int compLineKey = 1; compLineKey <=
33                    fullSize; compLineKey++){
```



```

24         String newKey = new String();
25         String compLineValue = String.
           valueOf(compLineKey);
26         // According to their positions
           we put the first one first
           in the comparison couple
27         if (compLineKey != lineKey) {
28             if (lineKey <
               compLineKey){
29                 newKey =
                   lineValue +
                   "," +
                   compLineValue
                   ;
30             } else {
31                 newKey =
                   compLineValue
                   + "," +
                   lineValue;
32             }
33             Text newKeyText = new
               Text(newKey);
34             context.write(
               newKeyText,
               lineText);
35         }
36     }
37 }
38 }
39 }

```

#### Reducer :

```

1 public static class Reduce extends Reducer<Text,Text,Text,Text> {
2
3     public void reduce(Text key, Iterable<Text> values,
           Context context) throws IOException,
           InterruptedException {
4
5         // We then create an array with the lines of
           the same key
6         List<String> linesOfKey = new ArrayList<String>
           >(2);
7         for (Text value : values){
8             linesOfKey.add(value.toString());
9         }
10
11         // After checking that we compute a couple only
           once . .
12         if (linesOfKey.size() == 2){
13             String d1 = linesOfKey.get(0);
14             String d2 = linesOfKey.get(1);
15             String [] wordsOf1 = d1.split(" ");
16             String [] wordsOf2 = d2.split(" ");
17
18             // .. We compute the Jacquard
               similarity as enunciated in the
               assignment
19             int intersection = 0;
20             int union = 0;

```

```

21
22 // inspiration from : https://github.
    com/tdebatty/java-string-similarity
    /blob/master/src/main/java/info/
    debatty/java/stringssimilarity/
    Jaccard.java
23 for (String word : wordsOf1){
24
25     // If a word is in both, add
        one to intersection. Always
        add one to union.
26     if (d2.contains(word)){
27         intersection ++;
28     }
29     union++;
30     // Now check the leftover part:
        only add to union if we
        have not seen the word
        before, i.e. if it is only
        in wordsOf2
31     for (String word : wordsOf2){
32         if (!(d1.contains(word)
33             )) {
34             union++;
35         }
36
37 float Jaccsim = (float) intersection /
    union;
38
39 // We now only return the couples that
    have a Jacquard similarity over 0.8
40 if (Jaccsim > 0.8){
    String out = d1 + " - " + d2 +
        " - Jaccsim = " + String.
        valueOf(Jaccsim);
41    Text value = new Text(out);
42    context.write(key, value);
43 }
44 // Increment our counter !!
45 context.getCounter(counting.
    NUMBER_OF_COMPARISONS).increment(1)
    ;
46     }
47 }
48 }

```

Of course, I am here taking the preprocessed pg100.txt as input, not simply pg100.txt (see Process file if necessary).

As you can see in the two following screenshots, running the job for 50 lines took 25sec, and for lines it took 30sec.

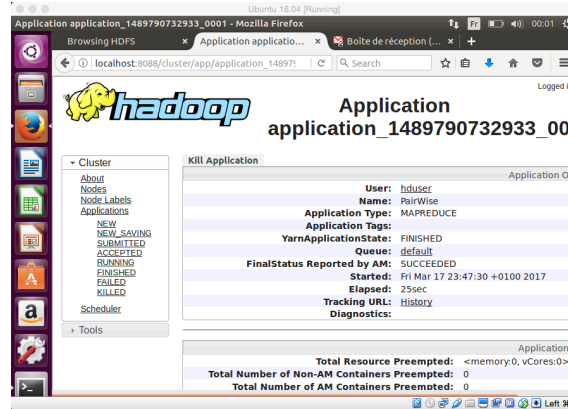


Figure 7: Pairwise job success with 50 lines.

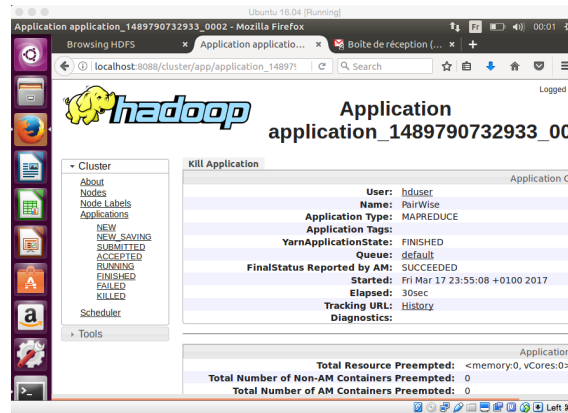


Figure 8: Pairwise job success with 1 000 lines.

The number of comparisons vary: 1 225 for 50 lines and 499 500 for 1 000 lines (as seen in the two following screenshots). This is coherent with the fact that there should be  $\frac{n*(n-1)}{2}$  computations possible (so to do), with n the number of lines.

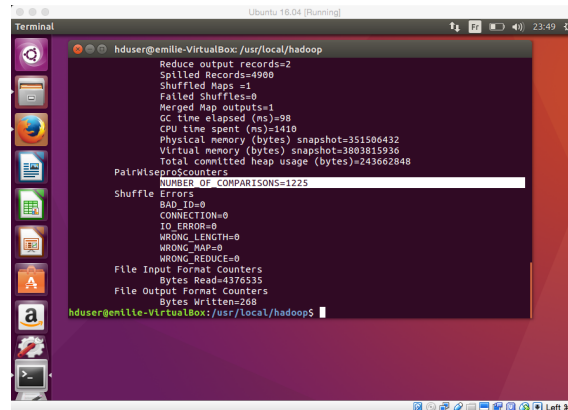


Figure 9: Number of comparisons for 50 lines.

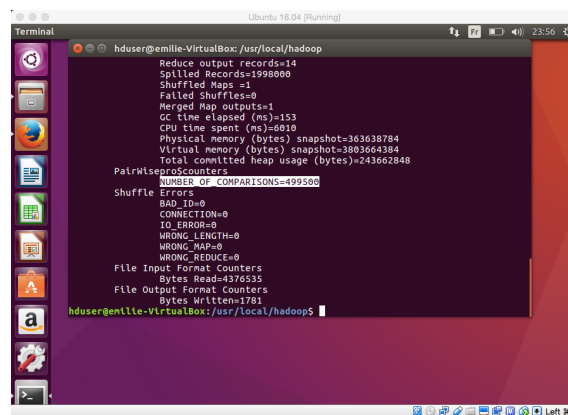


Figure 10: Number of comparisons for 1 000 lines.

I chose only to output the lines with a Jaccard similarity over 0,8 to skip comparisons between non-similar lines.

#### Extract of the pairwise comparison for 50 lines :

- 1 22,38 **version** works complete this william the electronic of – **version**  
works complete this william the electronic of – Jaccsim = 1.0
- 2 23,39 copyright world 1993 1990 inc library shakespeare is by –  
copyright world 1993 1990 inc library shakespeare is by and –  
Jaccsim = 0.9

#### Extract of the pairwise comparison for 1 000 lines :

- 1 2,132 shakespeare william – shakespeare william – Jaccsim = 1.0
- 2 22,122 **version** works complete this william the electronic of – **version**  
works complete this william the electronic of – Jaccsim = 1.0
- 3 22,38 **version** works complete this william the electronic of – **version**  
works complete this william the electronic of – Jaccsim = 1.0
- 4 23,123 copyright world 1993 1990 inc library shakespeare is by and –  
copyright world 1993 1990 inc library shakespeare is by and –  
Jaccsim = 1.0
- 5 23,39 copyright world 1993 1990 inc library shakespeare is by –  
copyright world 1993 1990 inc library shakespeare is by and –  
Jaccsim = 0.9
- 6 24,124 college illinois benedictine provided etext project gutenber  
by of – college illinois benedictine provided etext project  
gutenberg by of – Jaccsim = 1.0
- 7 25,125 readable machine permission with be may copies electronic and –  
readable machine permission with be may copies electronic and –

```
Jacsim = 1.0
s 26,126 your long such as so others distributed copies are for - your
    long such as so others distributed copies are for - Jacsim = 1.0
```

### 3.2 Inverted Index and prefix filtering

As said previously, I did not have the time to implement it.

### 3.3 Explain the difference in the number of comparisons and the computation time

I am guessing that the prefix filtering will generate a drop in execution time and reduce the number of comparisons. Indeed, as I understood it, it allows us to compute only the couples likely to be similar (thus removing many unnecessary comparisons). It is based on comparing the least common words only.