

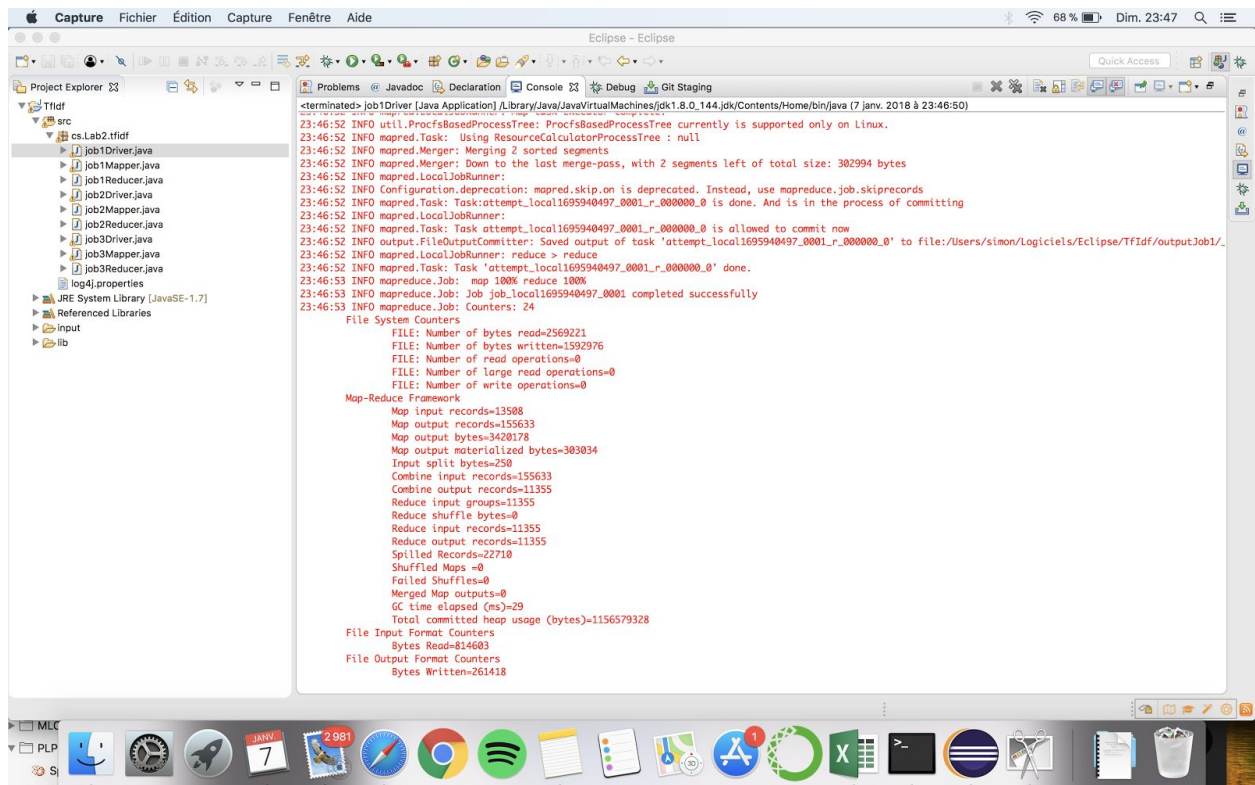
Rapport PLP

I. Hadoop

Q5.1

We did three MapReduce jobs to complete answer to the problem. Here is the screenshots of each job, launched one by one :

JOB1



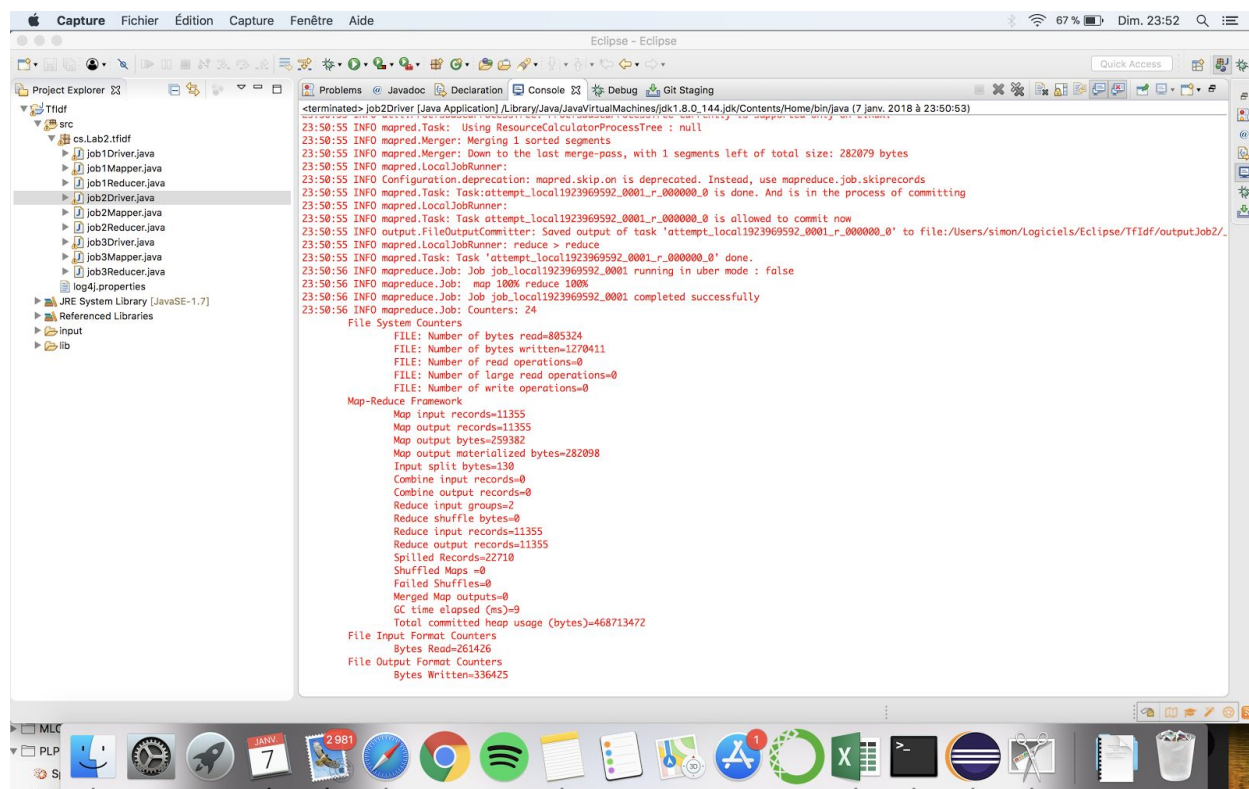
```
<terminated> job1Driver [Java Application] [Library/Java/JavaVirtualMachines/jdk1.8.0_144.jdk/Contents/Home/bin/java (7 janv. 2018 à 23:46:50)]
23:46:52 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
23:46:52 INFO mapred.Task: Using ResourceCalculatorProcessTree : null
23:46:52 INFO mapred.Merger: Merging 2 sorted segments
23:46:52 INFO mapred.Merger: Down to the last merge-pass, with 2 segments left of total size: 302994 bytes
23:46:52 INFO mapred.LocalJobRunner:
23:46:52 INFO Configuration.deprecation: mapred.skip on is deprecated. Instead, use mapreduce.job.skiprecords
23:46:52 INFO mapred.Task: Task:attempt_local1695940497_0001_r_000000_0 is done. And is in the process of committing
23:46:52 INFO mapred.LocalJobRunner:
23:46:52 INFO mapred.Task: Task attempt_local1695940497_0001_r_000000_0 is allowed to commit now
23:46:52 INFO output.FileOutputCommitter: Saved output of task 'attempt_local1695940497_0001_r_000000_0' to file:/Users/simon/Logiciels/Eclipse/Tf1df/outputJob1/.
23:46:52 INFO mapred.LocalJobRunner: reduce > reduce
23:46:52 INFO mapred.Task: Task 'attempt_local1695940497_0001_r_000000_0' done.
23:46:53 INFO mapreduce.Job: map 100% reduce 100%
23:46:53 INFO mapreduce.Job: Job job_local1695940497_0001 completed successfully
23:46:53 INFO mapreduce.Job: Counters: 24

File System Counters
  FILE: Number of bytes read=2569221
  FILE: Number of bytes written=1592976
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0

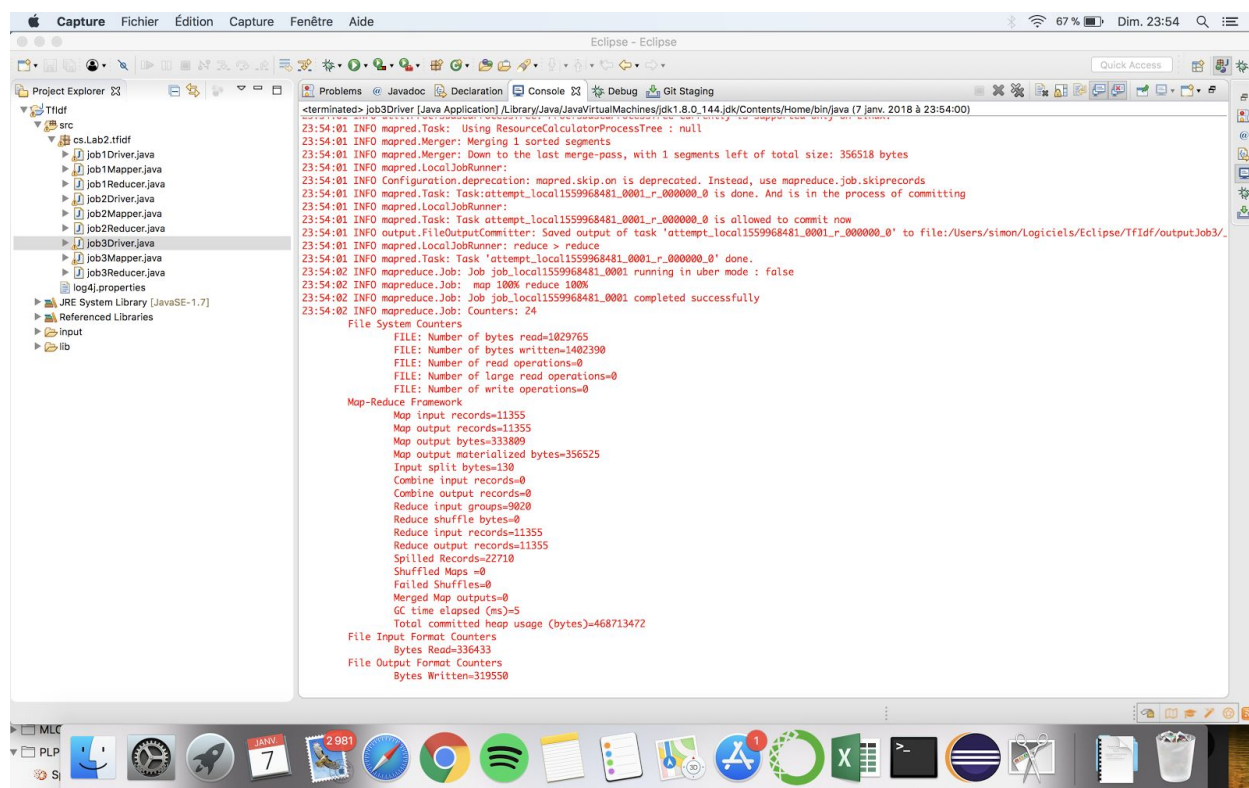
Map-Reduce Framework
  Map input records=13508
  Map output records=155633
  Map output bytes=3420178
  Map output materialized bytes=303034
  Input split bytes=250
  Combine input records=155633
  Combine output records=11355
  Reduce input groups=11355
  Reduce shuffle bytes=0
  Reduce input records=11355
  Reduce output records=11355
  Spilled Records=22710
  Shuffled Maps=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=29
  Total committed heap usage (bytes)=1156579328

File Input Format Counters
  Bytes Read=814603
File Output Format Counters
  Bytes Written=261418
```

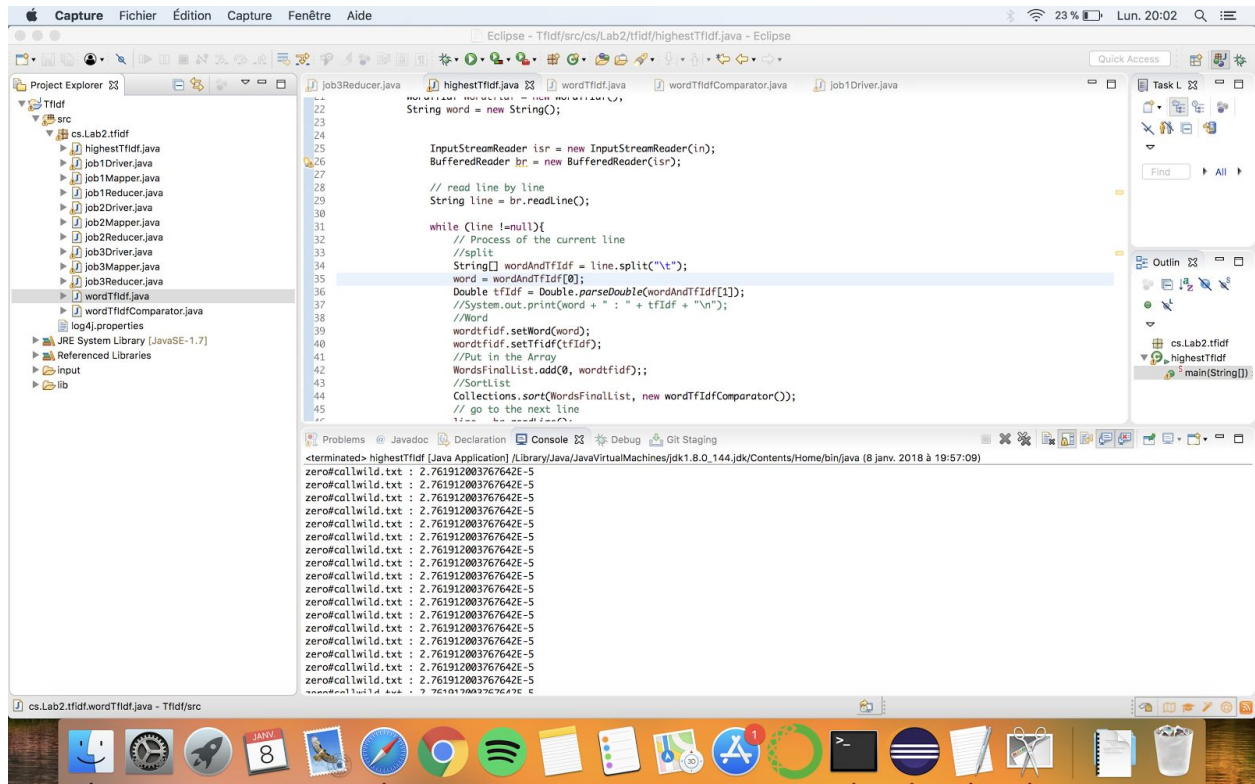
JOB2



JOB3

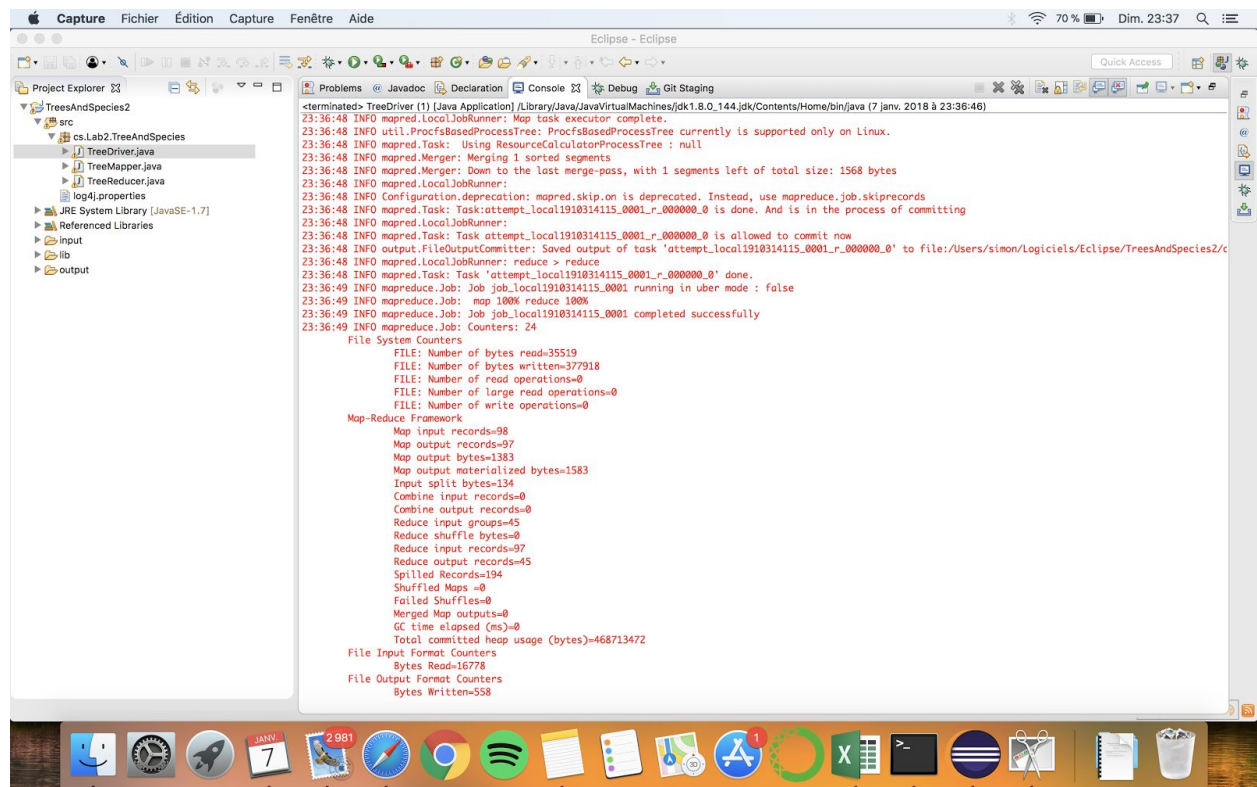


In the end I don't know why but I couldn't compute the ordered list of TFIDF. I don't know why but when adding one element to my arrayList, the array List become a List of that unique element, see the screenshot :

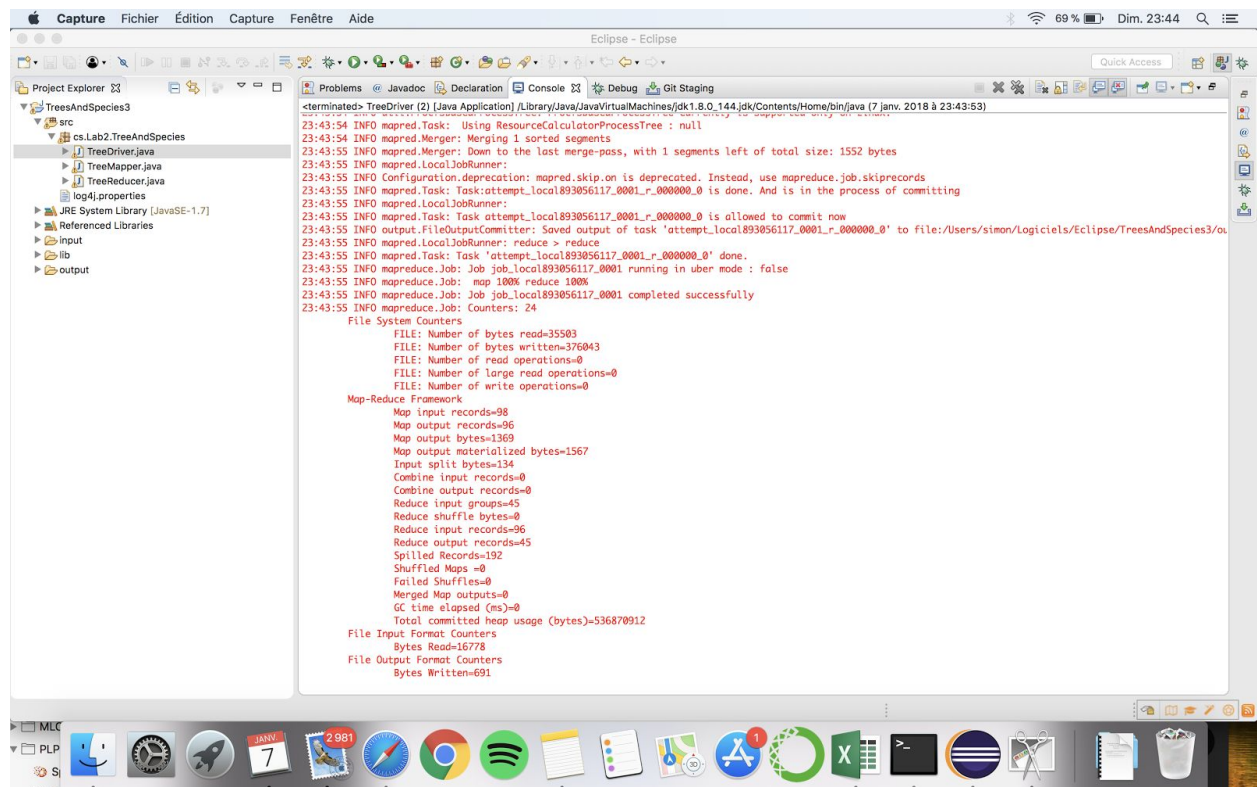


5.3 :

First Part :



Second Part :



II. Spark

A. Part 1 :Starter

1. *Print the top 10 most frequent words with their probability of appearance*

Script : *wordCount.py*

We use the `.split()` methods to split the lines into words. However that does not take into account special characters. It means : “the.” or “the:” is not taken into account to count the appearance of the.

2. *Get rid of special characters (.,:!?')*

Script : *wordCountSpeChar.py*

To get rid of the special characters, we call the `replace` method as follow :

```
line.replace('.', ' ')\n      .replace(',', ' ')\n      .replace(':', ' ')\n      .replace('!', ' ')\n      .replace('?', ' ')\n      .replace('"', ' ')\n      .split()
```

3. *Identify the transformations and the actions in our script*

Transformations : `flatMap` : to split lines into words and associate each word to a key value.

Action : `reduceByKey` : gather the key together and sum the key values.

Count: return the number of element in the dataset

4. *How many times are the transformations evaluated*

The transformations are evaluated for each action. In this case twice.

5. *Can you reduce this number ?*

We can persist the datasets to cache the results of transformations and reuse them in other actions on these datasets.

B. Part 1 :Intermediate

1. *Print the top 10 words from the Iliad that have “most disappeared” in The Odyssey*
2. *Do the same by swapping the Iliad and The Odyssey*
3. *Improve your script by getting rid of stopwords*

Script : *wordDisappeared.py*

4. *Use the Spark UI to make your script faster*

C. Part 2 : warmup

1. Create a cluster of 4-6 Spark nodes

Command on windows :

- Masternode : `bin\spark-class org.apache.spark.deploy.master.Master`
- Slavenodes : `spark-class org.apache.spark.deploy.worker.Worker spark://ip:port`
 2. If possible, launch a couple HDFS datanodes
 3. Launch `wordcount.py` on `iliad.mb.txt`

`bin\spark-submit --master spark://192.168.0.17:7077 wordCount.py iliad100.txt`

4. Launch `wordcount.py` on `raw.txt`
5. Launch two jobs at the same time. Make them run at the same time

The master process only one at a time, and the other one is "waiting".

6. What happens when a Spark node is brutally shutdown

The master notice the a slavenote shut down : "Lost executor 0 on 192.168.0.17: worker lost"
Then the master redistribute the work and still manage to complete the job thanks to the other workers.

D. Part 2 : Intermediate

1. What is Word2Vec ?

It is a technique that convert words into numeric vectors which can then be "fed into" various machine learning models to perform natural language processing.

2. Create a Word2Vec model of the Iliad
3. Who is Achilles + (Priam - Hector)

Script : `word2Vec.py`

E. Part 3 : Velib

1. *Create a Spark Streaming app that reads data from `velib.behmo.com`*

Any station has its "contract_name" + "name" that identify the station.

We have several times the same station for each RDDs. We use the ".transform(lambda rdd : rdd.distinct())" to solve this issue.

2. *Every 5 sec : print the empty Velib stations*

Script: `velibEmpty.py`

3. Every 5s: print the Velib stations that have become empty

Script: `velibGetEmpty.py`

We use the ".updateStateByKey" method that stores the following information:

- 1 : Station just got empty
- 0 : Station is empty but haven't just turn empty during the last batch
- 1 : Station is not empty

Then we can filter by value to get the stations that just got empty.

4. *Every 1 min: print the stations that were most active during the last 5 min
(activity = number of bikes borrowed and return)*

Script: `velibMostActive.py`

We use “.reduceByKeyAndWindow” method to take into account a 5 min window and to print the results every minute.