



ENSEA

Beyond Engineering

AI & BIG DATA

Big Data

Lab report : Big Textual Data Analytics

2nd year student

Erwan BROUDIN
15 février 2023

Exercise 0 : Download and Prepare your Document Corpus

In this part we have to download the three texts Romeo and Juliet, Hamlet and Richard III. To do this, we use the `addFile` (present in spark) and `textFile` functions to create an RDD that we can use later. Here is how we import the different texts :

```
1 ##### Download all the text files
2 import time
3 from pyspark import SparkFiles
4
5 RJ_url = "https://www.gutenberg.org/files/1112/1112.txt"
6 spark.sparkContext.addFile(RJ_url)
7 RomeoJuliet_rdd=sc.textFile(SparkFiles.get("1112.txt"))
8
9 Hamlet_url = "https://www.gutenberg.org/files/1524/1524-0.txt"
10 spark.sparkContext.addFile(Hamlet_url)
11 Hamlet_rdd=sc.textFile(SparkFiles.get("1524-0.txt"))
12
13 Richard_url = "https://www.gutenberg.org/cache/epub/1776/pg1776.txt"
14 spark.sparkContext.addFile(Richard_url)
15 Richard_rdd=sc.textFile(SparkFiles.get("pg1776.txt"))
```

Code 1 – Importation of the three texts

In the text file on Project Gutenberg, there are many lines at the beginning talking about the project. We decided not to delete them to avoid wasting time on a point that would not have given us anything more in the following, or we would not want to analyze the data but simply show that we know how to handle them.

Afterward, we need to clean up the text : there are some extracts that do not interest us. For example, words consisting of numbers, or even line breaks that only appear as `[""]`. We will first make the text "cleaner" with the function we call "cleaningFunc" which will remove all the function from a string, as well as put all the characters of a word in lower case. Here is how we do it :

```
1 ##### Create the cleaning function
2
3 import string
4
5 def cleaningFunc(text):
6     text = text.lower()
7     text = text.translate(text.maketrans("", "", string.punctuation))
8     return text
```

Code 2 – cleaningFunc

Finally, we want to be able to apply it to the three previously imported texts, so we need to make a `flatMap` that will allow us to define how to separate lines, as well as to apply the `cleaningFunc` function to the whole text. As we want to measure the execution time, we also set two timers that will allow us to calculate the execution time.

```
1 ##### Clean all the texts
2 start = time.time()
3
4 RJ_clean = RomeoJuliet_rdd.flatMap(lambda line : line.split(" ")).map(cleaningFunc)
5
6 Hamlet_clean = Hamlet_rdd.flatMap(lambda line : line.split(" ")).map(cleaningFunc)
7
8 Richard_clean = Richard_rdd.flatMap(lambda line : line.split(" ")).map(cleaningFunc)
9
10 end = time.time()
11
12 print(end-start)
13
14 Richard_clean.take(20)
```

Code 3 – Cleaning the three texts

The execution time of this part is 0.01 seconds.

We can check that the code has run by using the `take` function and seeing what happens then :

```
[',
'this',
'etext',
'file',
'is',
'presented',
'by',
'project',
'gutenberg',
'in',
'cooperation',
'with',
'world',
'library',
'inc',
'from',
'their',
'library',
'of',
'the']
```

IMAGE 1 – Richard III text cleaned

All the words have been separated, miniscule and the punctuation has been removed.

Exercise 1 : Count Words

In this exercise, we have two main objectives :

- to remove words that are not words ("2012" for example)
- to count the words so as not to have repetition

To do this, we will use the filter function which allows us to filter the words in the RDD according to a certain criterion : in this case, the fact that the value is only composed of a string of letters. We also do some mapping : each word is weighted by one, and when the program encounters a duplicate with "reduceByKey", we will add the two values together to see the value at the end of the process for each word.

Here is the code that we make, in order to count the words in each of the three texts :

```
1 ##### Count the (filtered) words of each documents
2
3 import re
4 import string
5 from pyspark import SparkContext
6
7 start = time.time()
8
9 RJ_CountWords = RJ_clean.filter(lambda x: re.match('[a-z]+', x)).map(lambda word: (word,
10 1)).reduceByKey(lambda x, y : x + y)
11 Hamlet_CountWords = Hamlet_clean.filter(lambda x: re.match('[a-z]+', x)).map(lambda word:
12 (word, 1)).reduceByKey(lambda x, y : x + y)
13 Richard_CountWords = Richard_clean.filter(lambda x: re.match('[a-z]+', x)).map(lambda
14 word: (word, 1)).reduceByKey(lambda x, y : x + y)
15
16 end = time.time()
17 print(end-start)
18
19 Richard_CountWords.take(10)
```

Code 4 – Counting each words

Once executed, we then display the execution time and we can also observe that each word has actually been counted to indicate the number of occurrences :

```
0.11334037780761719
[('this', 204),
 ('is', 296),
 ('presented', 1),
 ('project', 26),
 ('gutenberg', 21),
 ('in', 306),
 ('cooperation', 1),
 ('world', 31),
 ('library', 16),
 ('inc', 11)]
```

IMAGE 2 – Richard III text, words counted

The result is as expected, and we can even notice that the execution time is 0.12 seconds.

Exercise 2 : Finding Frequent Terms and Stop Words

In this part, we need to analyze stopwords only. Some words such as "this", "a", "for" are very common words, and it is these that we will be dealing with first. In order to keep only these words, we need to import the list of stopwords into a list object, using the nltk library. Here is how we did the import :

```
1 ##### Import the list of stopwords
2
3 import nltk
4 nltk.download('stopwords')
5 from nltk.corpus import stopwords
6 stopwords = stopwords.words('english')
```

Code 5 – Importation of the stopwords

Next, we map the word - as in the previous exercise - but add a criterion for filtering : the word is only kept if it is in the list of stopwords previously imported. Here is how we code this part :

```
1 ##### Keep only the stopwords from the text files
2
3
4 RJ_StopWords = RJ_clean.filter(lambda x : x in stopwords).filter(lambda x: re.match('[a-z]
5 Hamlet_StopWords = Hamlet_clean.filter(lambda x : x in stopwords).filter(lambda x: re.
6 Richard_StopWords = Richard_clean.filter(lambda x : x in stopwords).filter(lambda x: re.
```

Code 6 – Texts without words that are not stopwords.

For the next point, we need to combine the three texts into one RDD using the "union" function. To use it, we simply call it as a method, specifying as an argument the text on which we are doing the union. Then, we do a "reduceByKey" in order to keep only one time each word iteration (thus combining the number of occurrences of "this" in Romero and Juliet, Hamlet and Richard III), and we sort by decreasing number of occurrences of the stopwords.

Here is the code we write to achieve this :

```
1 ##### Create a rdd with all the stopwords of the 3 texts
2
3 all_StopWords_V1 = RJ_StopWords.union(Hamlet_StopWords).union(Richard_StopWords).
4 reduceByKey(lambda x, y : x + y).sortBy(lambda row : row[1], ascending=False)
5 all_StopWords_V1.take(30)
```

Code 7 – Every stopwords in a single RDD

The result displayed on the screen is the sum of the stopwords occurrences in the different texts :

```

('the', 3004),
('and', 2636),
('to', 2049),
('of', 1961),
('i', 1588),
('a', 1440),
('my', 1342),
('in', 1191),

```

IMAGE 3 – Occurrence of the most used stopwords

For the next question, we are only asked to reverse the columns : the occurrences must appear before the words. We do a different mapping as follows :

```

1 all_StopWords = all_StopWords_V1.map(lambda x : (x[1],x[0]))

```

Code 8 – Every stopwords in a single RDD with occurrence first

Finally, we create a DataFrame that contains all the data we were processing, and we will convert it as a CSV file. To do this, we simply call the "createDataFrame" function which will create it from our RDD, and then use the write.csv function which will save this DataFrame in a CSV format that we can export.

This is how it is done :

```

1 ##### Create the dataframe
2
3 start = time.time()
4
5 columns = ["count","stopword"]
6 allStopWordsDF = spark.createDataFrame(all_StopWords, columns)
7
8 allStopWordsDF.show()
9
10 ##### Create the csv file
11
12 allStopWordsFormatCSV = allStopWordsDF.write.csv("StopWord")
13
14 end = time.time()
15
16 print(end-start)

```

Code 9 – Create a DataFrame and a CSV File

This is the DataFrame we obtained :

```

+-----+-----+
|count|stopword|
+-----+-----+
| 3004|the|
| 2636|and|
| 2049|to|
| 1961|of|
| 1588|i|
| 1440|a|
| 1342|my|
| 1191|in|
| 1180|you|
| 1067|that|
| 1032|is|
| 877|with|
| 848|it|
| 841|not|
| 824|this|
| 764|for|
| 675|his|
| 672|me|
| 662|be|
| 616|but|
+-----+-----+

```

IMAGE 4 – DataFrame

We can also see the CSV file being saved in the default directory :

3004	the
2636	and
2049	to
1961	of
1588	i
1440	a
1342	my
1191	in
1180	you
1067	that
1032	is

IMAGE 5 – CSV File

The execution time of this question is about two seconds.

Exercise 3 : Simple Inverted Index

In this question, we will work on words that are not stopwords. We also do not want to know the number of occurrences of each word, but only in which text each word is represented. We then kept only the stopwords before joining the three texts :

```
1 ##### Delete the stopwords from the text files
2
3 start = time.time()
4
5 RJ_OtherWords = RJ_clean.filter(lambda x : x not in stopwords).filter(lambda x: re.match(
6     '[a-z]+', x)).map(lambda word: (word, "RJ.txt")).reduceByKey(lambda x, y : x )
7 Hamlet_OtherWords = Hamlet_clean.filter(lambda x : x not in stopwords).filter(lambda x:
8     re.match('[a-z]+', x)).map(lambda word: (word, "Hamlet.txt")).reduceByKey(lambda x, y
9     : x)
10 Richard_OtherWords = Richard_clean.filter(lambda x : x not in stopwords).filter(lambda x:
11     re.match('[a-z]+', x)).map(lambda word: (word, "Richard.txt")).reduceByKey(lambda x,
12     y : x)
13
14 ##### Create a rdd with all the other words of the 3 texts
15
16 all_OtherWords = RJ_OtherWords.union(Hamlet_OtherWords).union(Richard_OtherWords).
17     reduceByKey(lambda x, y : x + ", " + y).sortByKey().zipWithIndex().map(lambda x : (x
18     [1],x[0]))
19
20 end = time.time()
21
22 print(end-start)
23
24 all_OtherWords.take(10)
```

Code 10 – Sorting and indexing "Uncommon" words

This execution allows us to see that we have an identifier, followed by the list of common words arranged in alphabetical order and the names of the texts in which they are used.

```
[(0, ('abate', 'RJ.txt, Hamlet.txt')),
 (1, ('abatements', 'Hamlet.txt')),
 (2, ('abbey', 'RJ.txt')),
 (3, ('abbot', 'Richard.txt')),
 (4, ('abed', 'RJ.txt')),
 (5, ('abels', 'Richard.txt')),
 (6, ('abet', 'Richard.txt')),
 (7, ('abhorred', 'RJ.txt, Hamlet.txt')),
 (8, ('abhors', 'RJ.txt')),
 (9, ('abide', 'RJ.txt, Hamlet.txt, Richard.txt'))]
```

IMAGE 6 – RDD : Common words

Exercise 4 : Extended Inverted Index

For this last exercise, we need to extend the interest of the previous exercise. We want to see not only in which text the different words appear, but also how often they are used. To do this, we use a counter which we will increment. This counter is a little different from the ones we are used to seeing because it is put into a string in each iteration. Here is how we do this :

```
1 ##### Delete the stopwords from the text files
2
3 start = time.time()
4
5 RJ_OtherWords_Extended = RJ_clean.filter(lambda x : x not in stopwords).filter(lambda x:
    re.match('[a-z]+', x)).map(lambda word: (word, "RJ.txt #1")).reduceByKey(lambda x, y
    : "RJ.txt #" + str(int(x[x.find("#")+1:])+int(y[y.find("#")+1:])) )
6 Hamlet_OtherWords_Extended = Hamlet_clean.filter(lambda x : x not in stopwords).filter(
    lambda x: re.match('[a-z]+', x)).map(lambda word: (word, "Hamlet.txt #1")).
    reduceByKey(lambda x, y : "Hamlet.txt #" + str(int(x[x.find("#")+1:])+int(y[y.find("#"
    +1:])) )
7 Richard_OtherWords_Extended = Richard_clean.filter(lambda x : x not in stopwords).filter(
    lambda x: re.match('[a-z]+', x)).map(lambda word: (word, "Richard.txt #1")).
    reduceByKey(lambda x, y : "Richard.txt #" + str(int(x[x.find("#")+1:])+int(y[y.find("#"
    +1:])) )
8
9 end = time.time()
10 print(end-start)
11
12 RJ_OtherWords_Extended.take(10)
```

Code 11 – Counting words in texts

We can see with the text of Romeo and Juliet that this code works well as shown by the display of the first 10 elements of the RDD :

```
[('project', 'RJ.txt #90'),
 ('gutenberg', 'RJ.txt #32'),
 ('ebook', 'RJ.txt #13'),
 ('juliet', 'RJ.txt #65'),
 ('shakespeare', 'RJ.txt #8'),
 ('files', 'RJ.txt #3'),
 ('produced', 'RJ.txt #2'),
 ('proofing', 'RJ.txt #1'),
 ('tools', 'RJ.txt #1'),
 ('developed', 'RJ.txt #1')]
```

IMAGE 7 – Romeo and Juliet : words counted

Then we simply need to combine these three texts into a single RDD and index them. Before doing this, we have arranged them in alphabetical order for easier reading.

```
1 ##### Create a rdd with all the other words of the 3 texts
2
3 start = time.time()
4
5 all_OtherWords_Extended = RJ_OtherWords_Extended.union(Hamlet_OtherWords_Extended).union(
    Richard_OtherWords_Extended).reduceByKey(lambda x, y : x + ", " + y).sortByKey().
    zipWithIndex().map(lambda x : (x[1],x[0]))
6
7 end = time.time()
8
9 print(end-start)
10
11 all_OtherWords_Extended.take(10)
```

Code 12 – Indexing sorting and make the union of all "Uncommon words"

We can see from the following display that we have an indexed rdd, sorted alphabetically by word, which allows us to retrieve both the text from which it is taken, as well as the number of occurrences in each.

```
[(0, ('abate', 'RJ.txt #1, Hamlet.txt #1')),  
 (1, ('abatements', 'Hamlet.txt #1')),  
 (2, ('abbey', 'RJ.txt #1')),  
 (3, ('abbot', 'Richard.txt #7')),  
 (4, ('abed', 'RJ.txt #1')),  
 (5, ('abels', 'Richard.txt #1')),  
 (6, ('abet', 'Richard.txt #1')),  
 (7, ('abhorred', 'RJ.txt #1, Hamlet.txt #1')),  
 (8, ('abhors', 'RJ.txt #1')),  
 (9, ('abide', 'RJ.txt #1, Hamlet.txt #1, Richard.txt #1'))]
```

IMAGE 8 – RDD : Indexed, sorted and counted common words

This exercise took 1.28 second to complete.