LAB DISTRIBUTED DATA ANALYTICS I TUTORIAL VI

Juan Fernando Espinosa Reinoso 303158

Group 1: Monday Tutorial

Exercise 1: Basic Map and Reduce (2 Points)

1.1 Dataset

First, for this exercise, the requirement is to use the **Gutenberg_ebook.txt** file. The goal is to get the **WordCount** in the ebook. First, a directory is created by aplying the following command: **hadoop fs -mkdir /dda/exercise1**. Next, putting the dataset in the folder is necessary.

Now it is time to dive in in the procedure.

1.2 Mapper

By taking advantage of **sys** library in python which allows to pass data between the mapper and reducer, initialization and extraction of each line of text is executed first. Splitting each word for the purpose and finally, by iterating through all the words print an output as follows: (word, 1)

```
#!/usr/bin/env python
import sys

for line in sys.stdin:
    line = line.split()

    for word in line:
        print('%s\t%s' % (word, 1))
```

The output of the mapper looks as follow (small sample):

```
Project 1
Gutenberg
                 1
Literary
Archive 1
Foundation,
                 1
how
         1
        1
to
help
produce 1
our
new
eBooks, 1
and
```

```
1
how
         1
to
subscribe
                  1
to
         1
         1
our
         1
email
newsletter
to
         1
         1
hear
         1
about
```

1.3 Reducer

The Reducer receives as input the (word, 1) printed out by the mapper. Next, an iteration through all the words is applied: First, splitting the input data between words and value. A condition is applied for storing the counting words in a dictionary. If the word already exists as a **key** in the dictionary called **word_count** add +1 to the value, otherwise add the word as a key with value of 1.

1.4 Execution

For the execution of the Map-Reduce process the HadoopStreaming comes to play. Check and copy the path of the streaming jar file saved within hadoop folders. It allows to run different codes without the default jar file. Taking the local paths from the mapper and reducer, as well as the path of the input from hadoop by applying the command: hadoop fs -ls /path/to/file. Once all the above-mentioned requirements are gotten it is time to execute.

We take the following information as execution command:

```
hadoop jar /usr/local/Cellar/hadoop/3.2.1_1/libexec/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar /
-mapper "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-1/mapper.py" /
-reducer "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-1/reducer.py" /
-input /dda/exercise1/gutenberg_ebook.txt /
-output /dda/exercise1/output
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-1/mapper.py /
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-1/reducer.py
```

he folaqefs We have to insert two times the path for the mapper and reducer: The first one points hadoop to the executables and the latter is for distributing the executables around all nodes. Finally, the output shown has the structure of (Word, count)

Note: Since there were no preprocessing steps, some punctuation signs are not identified.

```
0\x99': 1, 'weep,': 2, 'besmeared': 1, 'sadly.\xe2\x80\x99': 1, 'bought': 10, 'sure,\xe2\x80\x99': 1, 'daughters.': 2, 'daughters,': 3, 'o': 1, 'joy': 13, '\xe2\x80\x98\cord,': 2, 'job': 3, 'daughters;': 1, 'spoil': 1, 'NO': 2, 'arts,': 1, 'window-seat': 1, 'commands': 2, 'pie niversally)': 1, '\xe2\x80\x98\cord,': 2, 'job': 3, 'daughters;': 1, 'spoil': 1, 'NO': 2, 'arts,': 1, 'window-seat': 1, 'commands': 2, 'pie niversally)': 1, '\xe2\x80\x99\cord,': 1, 'grain': 3, 'this?': 2, '1.F.l.': 1, 'close.': 1, 'mouse-hole.': 1, 'close,': 1, 'mouse-hole.': 1, 'xe80\x99\cord,': 1, 'walk': 19, 'subscribe': 1, 'numbers,': 1, 'hewed': 1, 'THEY': 2, '\xe2\x80\x99\cord,': 4, 'house-door.\xe2\x80\x99': 80\x98\line\xe2\x80\x99': 1, 'painted': 3, 'cellarful': 1, 'ensuring': 1, 'visible.': 1, 'painter': 2, 'wishing-ring,': 1, 'kerchief': 8, 'nt,': 6, 'princess\xe2\x80\x99': 1, 'lived,': 1, 'will': 541, 'hovering': 1, 'intig-place.': 1, 'ask?\xe2\x80\x99': 1, 'carelessness': 16, 'ringing': 1, 'thus': 28, 'non': 1, 'ground, \xe2\x80\x99': 1, 'elbow,': 1, 'ants,': 1, 'kind,': 3, 'perhaps': 12, 'linger': 3, 'these,' ful': 1, 'wheelbarrow.': 1, 'ants,': 1, 'gets': 5, '\xe2\x80\x98\w99': 1, 'difficult': 2, 'rake': 3, 'beast': 11, 'killed!\xe2\x80\x99': 1 ent': 4, 'collar': 2, 'warming': 1, 'sons;': 1, 'noise,': 3, 'fighting': 3, 'States': 4, 'bridge.\xe2\x80\x99': 2, 'load': 4, 'cried': 103 3, 'think,\xe2\x80\x99': 3, 'cries': 2, 'provision,': 1, 'rooms,': 2, 'happiness': 2, 'price.': 1, 'wonder,': 3, 'smite': 1, 'kins:': 1, 'care,': 32, 'colours.\xe2\x80\x99': 1, 'cloak.': 4, 'cloak,': 5, 'know': 78, 'snail-shell--and': 1, 'LIMITED': 3, '99712..': 1, 'xe2\x80\x99': 1, 'did,\xe2\x80\x99': 1, 'corfined,': 1, 'mite': 1, 'raise': 2, 'throne': 3, 'herself;': 2, 'herself:' 8, 'rare': 1, 'carried': 42, '59 9, 'himself,': 3, 'himself.': 5, 'courtyard,': 4, 'cutting-board: 3, 'courtyard.': 1, 'feral': 14, 'hunger,\xe2\x80\x99': 1, 'mather': 2, 'unks\x99': 1, 'mather': 20, 'this,': 3, 'pumise': 6, 'brush\x99': 1, 'round?\xe2\x80\x99': 1, 'mad
```

The code works satisfactory.

Exercise 2: Data cleaning and text tokenization (6 points)

2.1 Dataset

The goal of the exercise is to overcome the issues experienced in the output of the first exercise. As a consequence, 5 text datasets were given. Similarly as the first exercise, a directory called **exercise2** were created and all the 5 files were stored in it. It is important to mention that all files have the same initial name (ie exercise2-1.txt, exercise2-2.txt, and on) because it is easier to execute all files using a trick that will be explained below.

Permission	↓↑ Owner	↓↑ Group	↓ ↑	Size	Į↑	Last Modified	Ţţ	Replication	Ţţ	Block Size	J1	Name	Į†	
-rw-rr	juanfer	supergroup	<u> </u>	546.97 KE	3	Jun 21 17:10		1		128 MB		exercise3-1.txt		
-rw-rr	juanfer	supergroup	<u> </u>	1.01 MB		Jun 21 17:10		1		128 MB		exercise3-2.txt		
-rw-rr	juanfer	supergroup	<u> </u>	228.6 KB		Jun 21 17:11		1		128 MB		exercise3-3.txt		m
-rw-rr	juanfer	supergroup	2	1.51 MB		Jun 21 17:11		1		128 MB		exercise3-4.txt		
-rw-rr	juanfer	supergroup	<u> </u>	905.71 KE	3	Jun 21 17:11		1		128 MB		exercise3-5.txt		â

2.2 Mapper

As from its name it is easy to infer, in this step the goal is to map which words are stop words and also if they have punctuation. Get the lines in all documents for next take every single word in it and make a binary classification: if the word is stop-word = 0, and 1 otherwise.

2.3 Reducer

Take the output from the **Mapper** and start iterating through each (word, classification). Next step is to create a dictionary with the classification binary value as key and all the words as values. This make it possible to drop all the misclassified values (i.e. 0) and return only the words that are not stop words nor have they punctuation.

```
#!/usr/bin/env python
import numpy as np
import os
import sys

words_cleaned = {}
for i, line in enumerate(sys.stdin):
    line = line.split()
    (word, count) = line
    if count not in words_cleaned:
        words_cleaned[count] = []
        words_cleaned[count].append(word)
    else:
        words_cleaned[count].append(word)
result = words_cleaned.pop('0')
print(words_cleaned)
```

2.4 Execution

The same requirements as explained in step 1.4 are gathered here, the execution command looks as follow:

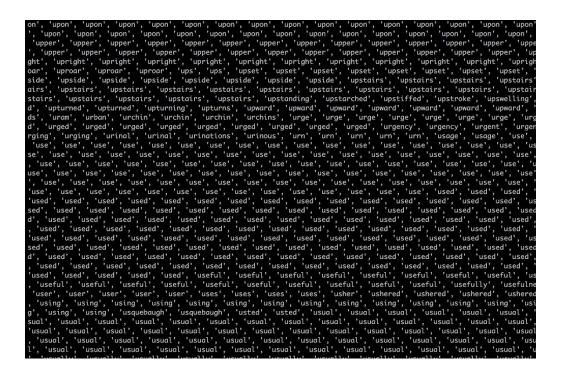
```
hadoop jar /usr/local/Cellar/hadoop/3.2.1_1/libexec/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar /
-mapper "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-2/ex2-mapper.py" /
-reducer "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-2/ex2-reducer.py" /
-input /dda/exercise2/exercise3*.txt /
-output /dda/exercise2/output /
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-2/ex2-mapper.py /
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-2/ex2-reducer.py
```

Here there are few keys that worth explain: To tell hadoop to run all the files at once and without writing all paths it is necessary to make all files have same initial name and then by adding an * the program will take all files that accomplish the condition. The output folder is created after running the code.

1=	Permission	11 Owner	↓↑ Group ↓↑	Size 1	Last Modified			11 Name	11
	-rw-rr	juanfer	supergroup	546.97 KB	Jun 21 17:10	1	128 MB	exercise3-1.txt	â
	-rw-rr	juanfer	supergroup	1.01 MB	Jun 21 17:10	1	128 MB	exercise3-2.txt	
	-rw-rr	juanfer	supergroup	228.6 KB	Jun 21 17:11	1	128 MB	exercise3-3.txt	â
	-rw-rr	juanfer	supergroup	1.51 MB	Jun 21 17:11	1	128 MB	exercise3-4.txt	
	-rw-rr	juanfer	supergroup	905.71 KB	Jun 21 17:11	1	128 MB	exercise3-5.txt	
	drwxr-xr-x	juanfer	supergroup	0 B	Jun 21 17:16	0	0 B	output	â

Finally, to run and explore the results the following command is inserted: hadoop fs -cat /dda/exercise2/output/filename.

Note: Hadoop always sorts the inputs, hence it is possible to appreciate several words next to one another.



Exercise 3: Calculate TFIDF scores of words/tokens (12 points)

3.1 Dataset

For this exercise, we use the same dataset applied in exercise 2. Hence no need to create a directory.

3.2 Mapper

As this exercise requires, we need to calculate the TF and IDF in order to get the values. The mapper is doing the following process:

- 1. preprocessing the data: erasing all stop-words and punctuaction.
- 2. It is selecting the name of the files for the IDF calculation. We need to know how many documents are there.
- 3. As the TFIDF happens to be a unique value per document it is necessary to keep the order of words, therefore, at the end the name of the file and the words are merged together, and to the right the count number 1.

```
#!/usr/bin/env python
import sys
import os
from string import punctuation
stop_words = ['a','able','about','across','after','all','almost',
         'by','can','cannot','could','dear','did','do','does','eith
        'hers','him','his','how','however','i','if','in','into','
        'my','neither','no','nor','not','of','off','often','on','o
'so','some','than','that','the','their','them','then','the
        'what', 'when', 'where', 'which', 'while', 'who', 'whom', 'why',
for line in sys.stdin:
    filename = os.environ["map_input_file"]
    line = line.translate(None, punctuation)
    line = line.split()
    for word in line:
        word = word.lower()
        if word.isalpha() == True and word not in stop_words:
             file_word = filename + "," + word
             print('%s\t%s' % (file_word, 1))
```

3.3 Reducer

The reducer task is to compute the TF and IDF simultaneously. The steps for doing so are stated next.

- 1. Call the outputs of the mapper and split it into the count and the (filename, word)
- 2. Next, split the filename and word as well by a comma-separated split.
- 3. Finally, merge into a list all the values for document and words
- 4. Iterate through the newly created list and create a dictionary with the **keys** as the filename and values as the words. This dictionary makes it pretty simple to extract the value of number of files as well as a dictionary of number of unique words with the respective number of appearances in a document.
- 5. As being said in the previous step a dictionary with the wordCount per document is created. Also by using Counter I got the number of unique words in the document. Then I proceed to calculate TF.
- 6. Followed it is necessary to get the number of documents each token appears on. Hence a small Counter iterating in the list of dictionaries created in TF makes it possible to extract the partial IDF.
- 7. Finally a dictionary applying the formula dof IDF is created with the words as keys and the values as outputs.
- A final for loop is useful for iterating through all the words in each document and multiply the respective IDF with the same word. Hence, finally the output is a document specific TFIDF.

```
import sys
import numpy as np
from collections import Counter
lista = []
TFIDF = []
documents = {}
for line in sys.stdin:
    file_word, count = line.split('\t')
    document, words = file_word.split(',')
    lista.append((document,words))
for i,e in lista:
    if i in documents.keys():
        documents[i].append(e)
        documents[i] = []
        documents[i].append(e)
number_documents = int(len(documents.keys()))
for i in documents.values():
    frequency = Counter(i)
    unique_words = int(len(list(set(i))))
    partial_tf = {k: (float(v) / float(unique_words)) for k,v in frequency.items()}
    TF.append(partial_tf)
```

3.4 Execution

Similar as the previous execution, the command is as follows:

```
hadoop jar /usr/local/Cellar/hadoop/3.2.1_1/libexec/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar /
-mapper "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-3/ex3-mapper.py" /
-reducer "python /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-3/ex3-reducer.py" /
-input /dda/exercise3/exercise3*.txt /
-output /dda/exercise3/output /
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-3/ex3-mapper.py /
-file /Users/juanfer/Documents/Maestria/SemesterIII/Lab-DDA/tutorial-6/exercise-3/ex3-reducer.py
```

The output of the algorithm is a list which contains dictionaries per document with the resultant TFIDF.

```
'0.000',
                                                                     'roams': '0.000', 'moonlight': '0.000',
                          '0.000',
                                                                                                                                              '0.001'.
                                       'stripes':
       '0.001', 'lamentations': '0.000',
                                                           'tanheap': '0.001', 'craftsmen':
                                                                                                             '0.000',
                                                                                                                          'once': '0.018',
                                                                                                                                                   'oral':
                               'forget': '0.001',
                  '0.000',
                                                           'caused': '0.001', 'beware': '0.000',
                                                                                                                                                  'neigh':
thoughtful': '0.000', 'concept': '0.000', 'barons': '0.000'
hankful': '0.000', 'dreadful': '0.001', 'griffin': '0.001',
                                                                              '0.000', 'dish': '0.003', 0.001', 'ravens': '0.001'
                                                                                                                     'follow':
                                                                                                                                                   'faster':
                                                                                                                                    '0.001
                                                                                                                       'hunting'
                                                                                                                                        '0.001'.
                                                                                                                    'solicit':
                           'spirit': '0.001', 'tail': '0.003',
                                                                                 'dinnertime': '0.000',
                                                                                                                                    '0.000',
nities': '0.000',
                                                       'belonging': '0.000', 'woman': '03', 'roe': '0.000', 'medicine':
                               'case': '0.001'
                                                                                         'woman': '0.013'
                                                                                                                    'returned': '0.002
.000', 'robes': '0.000', 'fat': '0.003',
                                                                                                        '0.000'.
                                                                                                                      'bucketful':
                                                                                                                                         '0.000'
ful': '0.000',
                      'livelihood': '0.000', 'tiresome': '0.000', 'condition': '0.000',
                                                                                                                          'washingday': '0.000', 'heaven':
                                                                                                                         'was'.
'0.002', 'sank':
                                  '0.000',
                                               'grandfather': '0.000', 'large': '0.005',
                                                                                                              'sang':
er': '0.000', 'alms': '0.000', 'small': '0.004', 'study': '0.000', 'indemnity': '0.000', 'sank': '0.000', '001', 'willowwrens': '0.000', 'honourably': '0.000', 'grasshopper': '0.000', 'healed': '0.000', 'past': 'e': '0.000', 'pass': '0.003', 'further': '0.001', 'httppglaforgfundraising': '0.000', 'situated': '0.000' heeses': '0.000', 'stood': '0.01', 'darkness': '0.001', 'sweetly': '0.000', 'clock': '0.001', 'deeply':
                                                                                                                                                  '0.000',
                                                                                                                                                             '0.000'
                    'stoop': '0.000', 'public': '0.001', 'version': '0.000', 'sunrise': '0.000', 'undecided':
300', 'full': '0.009', 'ranged': '0.000', 'godchildren': '0.000', 'compilation': '0.000', '
we': '0.000'.
                                                                                                                                                             '0.000',
                                        'threads': '0.000',
                                                                      'compressed': '0.000', 'trunk': '0.001', 'hast':
```

As we can see, since the words are repeated seldomly and the number of unique words per document are really high and also the number of words repeated across files makes the output to be low with some values as zero because it requires more decimals.

The output were saved in the directory and the executable output command is: hadoop fs -cat /dda/exercise3/output/part-00000



References:

- [1] https://towardsdatascience.com/installing-hadoop-on-a-mac-ec01c67b003c
- [2] https://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/
- [3] https://www.roseindia.net/bigdata/hadoop-shell-commands.shtml

- [4] http://rare-chiller-615.appspot.com/mr1.html
- [5] https://maelfabien.github.io/bigdata/MRJobP/#hadoop-streaming