

Tarea 3

Sistemas Distribuídos

Alejandro Díaz

Docente: Nicolas Hidalgo

Índice

1. Introducción	3
2. Desarrollo	3

1. Introducción

El objetivo de este trabajo es introducir a los estudiantes a los sistemas de procesamiento para grandes volúmenes de datos. Para ello, los alumnos deberán trabajar con Hadoop, un proyecto de software de código abierto que se puede utilizar para procesar de forma eficaz conjuntos de datos de gran tamaño. Los alumnos deberán aplicar la tecnología, aplicando un código enfocado a map-reduce y reconocer las ventajas de esta. Para lograr esto, se desea recopilar 10 artículos de wikipedia y hacer un mapa de las palabras encontradas con el respectivo documento donde aparecen y la frecuencia de aparición.

2. Desarrollo

En un principio, se crea el código para recopilar los textos a través de la api de wikipedia, este código fue creado en python 3

```

1  #!/usr/bin/env python
2  # -*-coding:utf-8 -*
3
4  import wikipedia as wiki
5  import regex as re
6  wiki.set_lang("es")
7
8  pages = []
9  pages.append('Inquisicion')
10 pages.append('Guerra santa')
11 pages.append('Brujeria')
12 pages.append('Herejia en el catolicismo')
13 pages.append('Color')
14 pages.append('Idioma griego')
15 pages.append('Caballeros Templarios')
16 pages.append('Orden militar')
17 pages.append('Orden de malta')
18 pages.append('Palacio magistral de la orden de malta')
19
20 counter = 0
21
22 for page in pages:
23     a = wiki.page(page)
24
25     page = page.strip().replace(" ", "_")
26     array = a.content.split(" ")
27     path = "./"
28
29     if counter < 5:
30         path = path + "docs1/" + page + ".txt"
31     else:
32         path = path + "docs2/" + page + ".txt"
33
34     f = open(path, "w")
35     f.write("Title: " + page + "\n")
36
37     for i in array:
38         f.write(i + " ")
39     counter += 1

```

En este código se recorre un arreglo con 10 títulos de páginas, en el cual por cada página se creará un documento .txt con el contenido de las páginas conseguidos a través de la api de wikipedia, considerando un contador para separar los documentos en las carpetas docs1 y docs2.

Luego, se investiga sobre la estructura de código de map reduce, el cual consta de 5 puntos:

- Asignar los procesadores que mapearán los input a través de index's que se asociarán a cada procesador, y se pasará la data de los archivos a cada procesador dependiendo de la key que le corresponda.
- Se utiliza el código de Map que se le indica a Hadoop, en este caso examples/mapper.py y este clasifica Título, palabra y 1, de esta manera se escribe el output en un almacenamiento temporal que será leído por los reducers.
- Se asignan los procesadores que reducirán el output anterior, a través de keys que asignarán la data que deberá trabajar cada procesador, y se provee a estos con la data necesaria.
- Se utiliza el código de reducer que se le indica a Hadoop, en este caso examples/reducer.py y este cuenta la frecuencia de las palabras a través de un diccionario, para luego utilizar un diccionario de índice invertido con key: palabra y value: array de tuplas de [página donde se encuentra la palabra, frecuencia (cantidad) de aparición].
- Finalmente se junta todo el output de los reducers y se escribe un output final.

El código de mapper.py utilizado fue el siguiente:

```

1 #!/usr/bin/env python
2 # -*-coding:utf-8 -*
3
4 import sys
5 import re
6
7 title = sys.stdin.readline().split(":")[1].strip()
8
9 for line in sys.stdin:
10     line = line.strip()
11
12     if line.startswith("Title:"):
13         title = line.split(":")[1].strip()
14
15     words = re.split('\W+', line)
16
17     words = [word for word in words if not word.isdigit() and len(word) > 1]
18     words = [word.lower() for word in words]
19
20     for word in words:
21
22         if word != '':
23             print('%s\t%s\t%s' % (title, word, 1))

```

El código de reducer utilizado fue el siguiente :

```

1 #!/usr/bin/env python
2 # -*-coding:utf-8 -*
3
4 import sys
5 import re
6
7 def SaveWords (hashmap, lastTitle, wordMap):
8     for word in hashmap.keys():
9         value = hashmap[word]
10         if word in wordMap:
11             array = wordMap[word]
12             array.append([lastTitle,value])
13             wordMap[word] = array
14         else:
15             wordMap[word] = ([[lastTitle,value]])
16
17 wordMap = {}
18 hashmap = {}
19 title = ""
20 lastTitle = ""

```

```

21 counter = 0
22
23 for line in sys.stdin:
24     title, word, count = line.split('\t')
25
26     if title != lastTitle and counter != 0:
27         SaveWords(hashmap, lastTitle, wordMap)
28         hashmap = {}
29
30     lastTitle = title
31     counter += 1
32     if word in hashmap:
33         hashmap[word] += 1
34     else:
35         hashmap[word] = 1
36
37 SaveWords(hashmap, lastTitle, wordMap)
38
39 for key, value in wordMap.items():
40     print('%s\t%s' % (key,value))

```

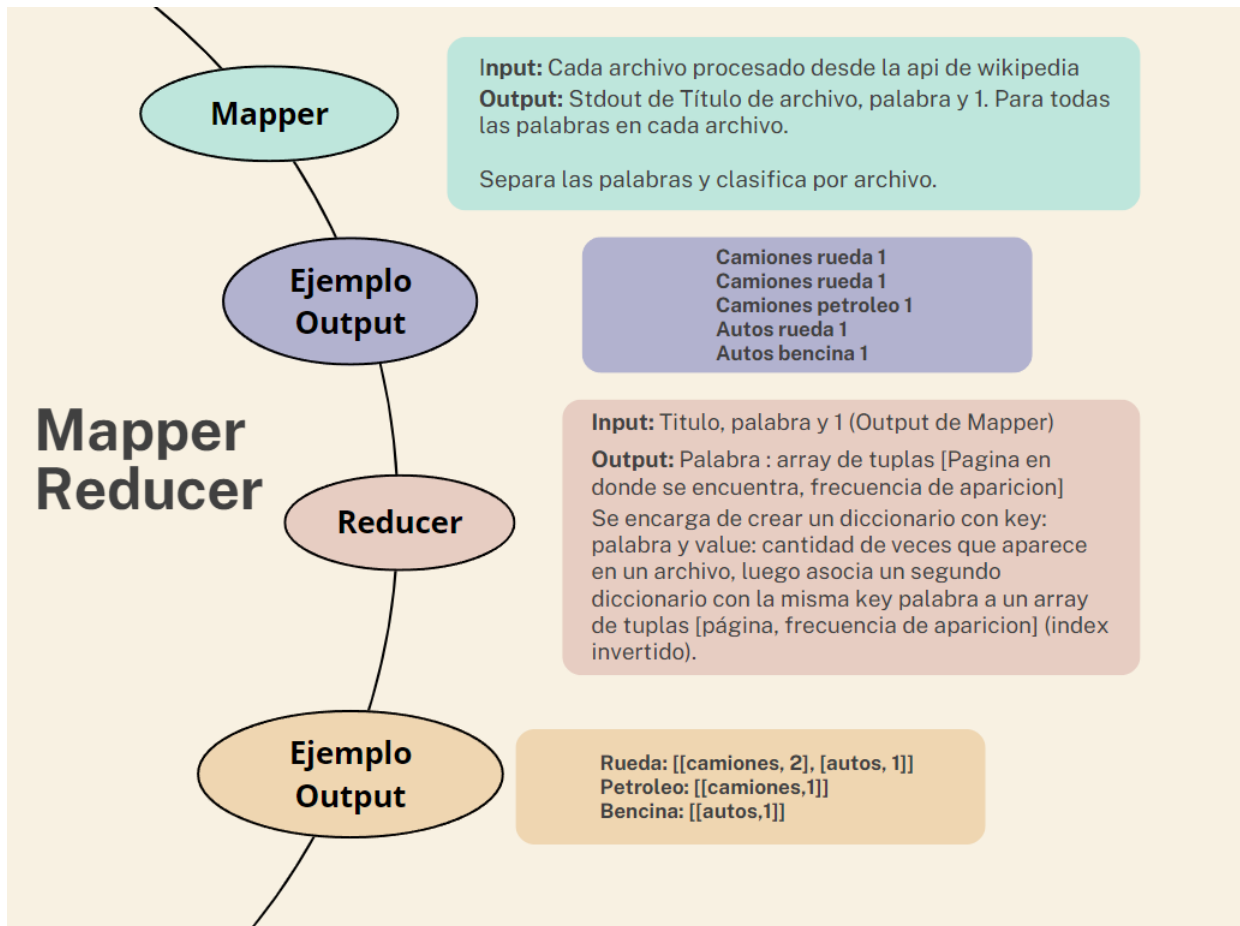


Figura 1: Diagrama de mapper reducer.

Finalmente, para el buscador, se crea un código que busca al ingresar una palabra buscará en el output de los reducer si existe la palabra y entregará los documentos en los cuales fue encontrada, con la cantidad de apariciones y la url de wikipedia del documento.

```

1 #!/usr/bin/env python
2 # -*-coding:utf-8 -*

```

```

3
4 import wikipedia as wiki
5 import regex as re
6 import sys
7 from unidecode import unidecode
8 wiki.set_lang("es")
9
10 def get_url(title):
11     return wiki.page(title).url
12
13 f = open('data.txt', 'r')
14 lines = f.readlines()
15
16 print('Ingrese palabras a buscar separadas por espacio')
17
18 # read line from sys stdin and split by space
19 words = sys.stdin.readline().split(' ')
20 # remove \n from last word
21 words[-1] = words[-1].strip()
22
23
24 for line in lines:
25     line = line.strip()
26     # split by tab
27     line = line.split('\t')
28     palabra = line[0]
29     # palabra to unidecode
30     palabra = unidecode(palabra)
31     # if palabra in words:
32     if palabra in words:
33         data = line[1]
34         # data remove first and last character
35         data = data[1:-1]
36         # split by comma ignoring commas inside brackets
37         data = re.split(r',\s*(?![^\]]*\))', data)
38         # remove special characters from every string on data
39         data = [re.sub(r'[^\w\s]', '', i) for i in data]
40         counter = 0
41
42         titles = []
43         q = []
44         for i in data:
45             if counter % 2 == 0:
46                 titles.append(i)
47             else:
48                 q.append(i)
49                 counter += 1
50
51         print('Palabra: ' + palabra)
52         for i in range(len(titles)):
53             url = get_url(titles[i])
54             toPrint = '|Titulo: ' + titles[i] + '| Cantidad: ' + q[i] + '| URL: ' + url + '|\n'
55             for i in range(len(toPrint)):
56                 print('-', end='')
57             print('\n', end='')
58             print(toPrint, end='')
59             for i in range(len(toPrint)):
60                 print('-', end='')
61             print('\n')

```

```
hduser@sd: ~/examples
hduser@sd:~/examples$ sudo python3 browser.py
Ingresa palabras a buscar separadas por espacio
santa maria imperio perdida
Palabra: maria
-----
|Titulo: Caballeros_Templarios| Cantidad: 1| URL: https://es.wikipedia.org/wiki/Caballeros_templarios|
-----
|Titulo: Orden_militar| Cantidad: 2| URL: https://es.wikipedia.org/wiki/Orden_militar|
-----
Palabra: santa
-----
|Titulo: Caballeros_Templarios| Cantidad: 18| URL: https://es.wikipedia.org/wiki/Caballeros_templarios|
-----
|Titulo: Guerra_santa| Cantidad: 6| URL: https://es.wikipedia.org/wiki/Guerra_santa|
-----
|Titulo: Herejia_en_el_catolicismo| Cantidad: 1| URL: https://es.wikipedia.org/wiki/Herej%C3%ADa_en_el_catolicismo|
-----
|Titulo: Inquisicion| Cantidad: 4| URL: https://es.wikipedia.org/wiki/Inquisici%C3%B3n|
-----
|Titulo: Orden_de_malta| Cantidad: 10| URL: https://es.wikipedia.org/wiki/Orden_de_Malta|
-----
|Titulo: Orden_militar| Cantidad: 10| URL: https://es.wikipedia.org/wiki/Orden_militar|
-----
|Titulo: Palacio_magistral_de_la_orden_de_malta| Cantidad: 2| URL: https://es.wikipedia.org/wiki/Palacio_Magistral_de_la_Orden_de_Malta|
-----
Palabra: imperio
-----
|Titulo: Caballeros_Templarios| Cantidad: 1| URL: https://es.wikipedia.org/wiki/Caballeros_templarios|
-----
|Titulo: Idioma_griego| Cantidad: 6| URL: https://es.wikipedia.org/wiki/Idioma_griego|
-----
|Titulo: Inquisicion| Cantidad: 2| URL: https://es.wikipedia.org/wiki/Inquisici%C3%B3n|
-----
```

Figura 2: Funcionamiento del buscador browser.py