

UNIVERSIDAD DIEGO PORTALES
FACULTAD DE INGENIERÍA Y CIENCIAS
ESCUELA DE INFORMÁTICA Y TELECOMUNICACIONES



Sistemas Distribuidos
Tarea 3

Profesor:
Nicolás Hidalgo

Ayudantes:
Nicolás Nuñez
Cristian Villavivencio
Joaquin Fernandez

Estudiantes:
Nicolás Moncada
Dayana Ruminot

Índice

1. Problema y Solución	1
2. Códigos	2
2.1. Extracción de 10 documentos de texto y división de documentos en carpetas	2
2.2. WordCount para cada documento	2
2.3. Guardar datos en base de datos y buscador	3
3. Hadoop Streaming	4
4. Bibliografía	7

1. Problema y Solución

Los buscadores implementados dentro del curso no son completamente eficientes. Se han sacado datos de distintas páginas y han sido guardados directamente en la base de datos. Sin embargo, la técnica utilizada es lenta y redundante. ¿Para qué tener datos repetidos (los documentos de texto y los datos en la base de datos)? ¿Hay una forma más eficiente? Dicho esto, se propone armar un buscado que funcione de mejor manera a modo de desafío personal, utilizando MapReduce.

Para dar solución a este problema es realizar lo siguiente:

- Obtener 10 documentos de texto utilizando la API de wikipedia. Para este ítem es necesario mostrar el código empleado.
- Dividir los 10 documentos del ítem anterior en dos carpetas distintas. Desde el documento 1 al 5, irían en la carpeta 1. Desde el documento 6 al 10 irían en la carpeta 2.
- Investigar un código de MapReduce que permita contar palabras en un documento. En caso de ser necesario se solicita mostrar el código. Adicionalmente, se debe mostrar un diagrama que explique el funcionamiento del código.
- Generar un trabajo de WordCount para todos los documentos, utilizando Hadoop. Este trabajo de WordCount generará Key-values. Por ejemplo:
 - Documento1: Hola soy Nicolás y él es Joaquín.
 - Documento2: Hola soy Nicolás. Hola soy Joaquín.

El output del trabajo de WordCount generará a los siguientes valores:

Word	[(Document1, Count1), ...]
Hola	(1, 2) (2, 2)
soy	(1, 1) (2, 2)
y	(1, 1)
...

Figura 1: Output del trabajo de WordCount.

- Guardar los valores anteriores en alguna base de datos. Puede ser de cualquier tipo.
- Generar un buscador el cual deberá utilizar como regla el buscar la mayor cantidad de coincidencias posibles utilizando el índice invertido. Finalmente, deberá retornar URLs que permitan ver los documentos de manera específica.
- Como último paso se debe documentar en vídeo de no más de 10 minutos todo el funcionamiento del buscador.

2. Códigos

Todos los códigos se pueden encontrar en el repositorio <https://github.com/NicoMonc/Tarea3SD>

2.1. Extracción de 10 documentos de texto y división de documentos en carpetas

Se solicita extraer 10 documentos de Wikipedia, para posteriormente dividir estos 10 documentos en 2 grupos paralelos, para realizar esta tarea inicial se implementó el siguiente código:

```
1 import urllib.request
2 import os
3
4 for i in range(10):
5     imprimir=i+1
6     opcion = input ("Texto a buscar: ")
7     print(opcion)
8     if(i<5):
9         site = urllib.request.urlopen('https://es.wikipedia.org/w/api.php?action=query&list=search&srprop=snippet&format=json&origin=*&utf8=&srsearch='+opcion)
10        data = site.read()
11        print(data)
12
13        isFile = os.path.isfile("./1_5/search"+str(imprimir)+".txt")
14        if (isFile):
15            os.remove("./1_5/search"+str(imprimir)+".txt")
16
17        crear = open("./1_5/search"+str(imprimir)+".txt","wb") #open file in binary mode
18        crear.write(data)
19        crear.close()
20    else:
21        site = urllib.request.urlopen('https://es.wikipedia.org/w/api.php?action=query&list=search&srprop=snippet&format=json&origin=*&utf8=&srsearch='+opcion)
22        data = site.read()
23        print(data)
24
25        isFile = os.path.isfile("./6_10/search"+str(imprimir)+".txt")
26        if (isFile):
27            os.remove("./6_10/search"+str(imprimir)+".txt")
28
29        crear = open("./6_10/search"+str(imprimir)+".txt","wb") #open file in binary mode
30        crear.write(data)
31        crear.close()
```

Figura 2: API.py

En el código anterior se puede ver como mediante la api se accede al sitio web solicitado en la tarea, para posteriormente imprimir y separar en cada una de las carpetas solicitadas, integrando en la primera carpeta los primeros 5 archivos de texto y luego los 5 restantes en la segunda carpeta.

Previamente a separar los datos se recorren con mapperreduce y reducer como se verá a continuación.

2.2. WordCount para cada documento

Esta sección se realiza con mapreduce.py y reducer.py, siendo el primero una revisión de cada línea, para posteriormente en reducer realizar el recuento de palabras por cada archivo. Por lo cual, lo primero que se realiza posteriormente a el WordCount de los 10 archivos es guardarlos directamente en la base de datos, para posteriormente ser utilizado en api.py

```
1 import sys
2
3 for line in sys.stdin:
4     line = line.strip()
5     words = line.split()
6
7     for word in words:
8         print ("%s\t%s" % (word, 1))
```

Figura 3: mapreduce.py

```
1 from operator import itemgetter
2 import sys
3
4 current_word = None
5 current_count = 0
6 word = None
7
8 for line in sys.stdin:
9     line = line.strip()
10    word, count = line.split('\t', 1)
11    try:
12        count = int(count)
13    except ValueError:
14        continue
15
16    if current_word == word:
17        current_count += count
18    else:
19        if current_word:
20            print ("%s\t%s" % (current_word, current_count))
21            current_count = count
22            current_word = word
23
24 if current_word == word:
25     print ("%s\t%s" % (current_word, current_count))
```

Figura 4: reducer.py

2.3. Guardar datos en base de datos y buscador

Respecto a la base de datos, se tienen dos archivos para el funcionamiento, uno es bd.py y otra que trabaja con sql generando las tablas para guardar la información obtenida. El código de bd se ejecuta al previo funcionamiento de mapreduce.py y reducer.py con la intención de guardar en la base de datos sql los datos como; ID, palabra, número (wordcount), archivo.

Posteriormente atender los datos registrados se tiene el buscador que funciona dentro de la api, de momento no nos fue posible conectar nuestro proceso con hadoop. Sin embargo idea planteada es levantar una instancia en docker y realizar una conexión dentro de la terminal ingresar los comandos para conectar la api a hadoop como lo son HDFS y finalmente ingresar los comandos que están en el repositorio reverenciado (Hadoop Streaming Using Python).

3. Hadoop Streaming

El código de MapReduce elegido para contar las palabras de un documento es Hadoop Streaming. Hadoop Streaming es una característica que viene con Hadoop y permite a los usuarios o desarrolladores usar varios lenguajes diferentes para escribir programas de MapReduce como Python, C++, Ruby, etc. Es compatible con todos los lenguajes que pueden leer desde la entrada estándar y escribir en la salida estándar. Implementaremos Python con Hadoop Streaming y observaremos cómo funciona. Implementaremos el problema de conteo de palabras en python para comprender Hadoop Streaming.

Para realizar el mapeo y la reducción se utilizan `mapper.py` y `reducer.py`

- **mapper.py** toma un archivo de texto, lee el archivo, separa las palabras de acuerdo a los espacios que existan entre ellas, guarda la palabra en un arreglo, recorre la palabra y las imprime con un salto de línea entre cada una.

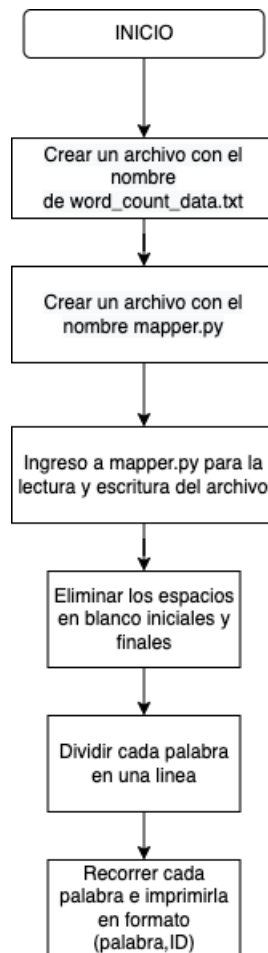


Figura 5: Diagrama del funcionamiento de `mapper.py`

- **reducer.py** Posteriormente a crear `mapper.py` se inicializan los demonios de hadoop, se crea una carpeta y dentro de este `crearword_count_data.txt`

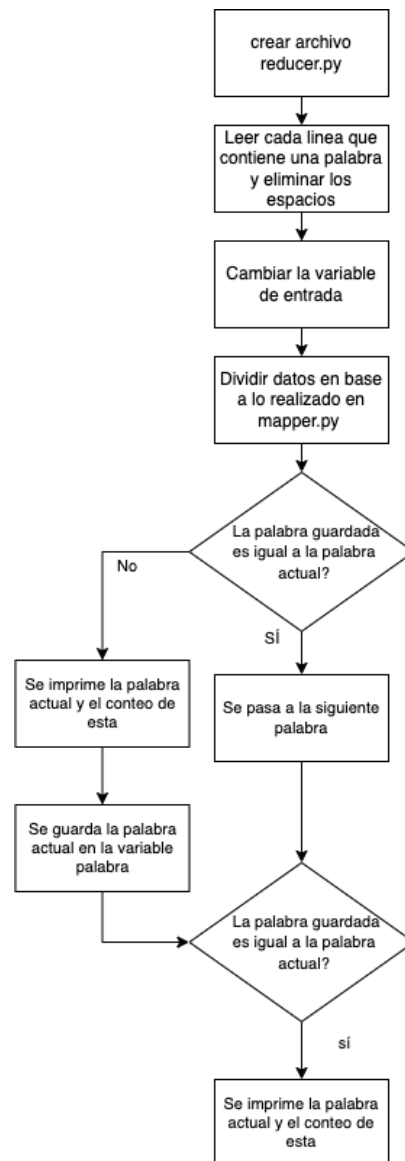


Figura 6: Diagrama del funcionamiento de `reducer.py`

Finalmente, el diagrama del funcionamiento en conjunto de ambos códigos y la configuración por terminal para autorizar lectura y escritura es el siguiente:

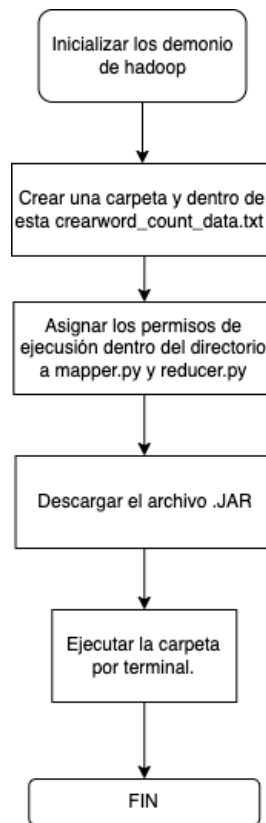


Figura 7: Diagrama de funcionamiento

4. Bibliografía

- **Vídeo del funcionamiento**
- **Repositorio**
<https://github.com/NicoMonc/Tarea3SD>
- **Hadoop Streaming Using Python**
<https://www.geeksforgeeks.org/hadoop-streaming-using-python-word-count-problem/>