

Práctica 3.1

Introducción

MapReduce es un modelo de programación para dar soporte a la computación paralela sobre grandes colecciones de datos en grupos de computadoras y al commodity computing. El nombre del framework está inspirado en los nombres de dos importantes métodos, macros o funciones en programación funcional: Map y Reduce.

MapReduce ha sido adoptado mundialmente, ya que existe una implementación OpenSource denominada Hadoop.

MapReduce se emplea en la resolución práctica de algunos algoritmos susceptibles de ser paralelizados. No obstante, MapReduce no es la solución para cualquier problema, de la misma forma que cualquier problema no puede ser resuelto eficientemente por MapReduce. Por regla general se abordan problemas con datasets de gran tamaño, alcanzando los petabytes de tamaño. Es por esta razón por la que este framework suele ejecutarse en sistema de archivos distribuidos (HDFS).

Hadoop es un framework de software que soporta aplicaciones distribuidas bajo una licencia libre. Permite a las aplicaciones trabajar con miles de nodos y petabytes de datos. Hadoop se inspiró en los documentos Google para MapReduce y Google File System (GFS).

La función Map recibe como parámetros un par de (clave, valor) y devuelve una lista de pares. Esta función se encarga del mapeo y se aplica a cada elemento de la entrada de datos, por lo que se obtendrá una lista de pares por cada llamada a la función Map. Después se agrupan todos los pares con la misma clave de todas las listas, creando un grupo por cada una de las diferentes claves generadas. No hay requisito de que el tipo de datos para la entrada coincida con la salida y no es necesario que las claves de salida sean únicas.

Map (clave1, valor1) → lista (clave2, valor2)

La función Reduce se aplica en paralelo para cada grupo creado por la función Map(). La función Reduce se llama una vez para cada clave única de la salida de la función Map. Junto con esta clave, se pasa una lista de todos los valores asociados con la clave para que pueda realizar alguna fusión para producir un conjunto más pequeño de los valores.

Reduce (clave2, lista(valor2)) → lista(valor2)

Desarrollo

En esta práctica, tuve problemas al momento de correr los ejemplos de los diferentes tutoriales que intenté seguir de Hadoop, tras instalarlo. Intenté cosas diferentes y en todas me marcaba errores extraños que no supe como arreglar.

Así que lo que hice fue pedirle ayuda a mi compañero Dominic Márquez. Él sí pudo correr los ejemplos y realizar la parte práctica necesaria para esta práctica. Así que con su ayuda, para la primera demostración utilicé un archivo de texto, el cual contenía las siguientes palabras:

“Jorge Adrian Padilla Velasco
Bases de datos avanzadas”

Obteniendo el siguiente resultado:

```
Adrian 1
Bases 1
Jorge 1
Padilla 1
Velasco 1
avanzadas 1
datos 1
de 1
```

Para la segunda demostración utilice un archivo de texto, el cual contenía la definición de MapReduce, obtenida de la página de [TutorialsPoint](https://www.tutorialspoint.com/).

Obteniendo el siguiente resultado:

```
Como 1
El 1
En 1
MapReduce 3
Mapa 2
a 1
algoritmo 1
basada 1
clave/valor). 1
combina 1
como 1
computación 1
conjunto 3
contiene 1
convierte 1
datos 2
datos, 1
de 8
después 1
distribuida 1
dividen 1
dos 1
el 3
elementos 1
en 5
```

```

entrada 1
es      1
implica,      1
importantes,  1
java.    1
la       3
los      2
lugar,   1
mapa     1
mapa.    1
modelo   1
más      1
nombre   1
otro     1
pequeño  1
procesamiento 1
programa  1
que      3
realiza  1
reducción      1
reducir 1
reducir.      1
saber    1
salida   1
se       3

```

```

secuencia      1
segundo 1
siempre 1
tarea, 1
tareas 1
toma 2
tuplas 2
tuplas. 1
técnica 1
un 4
una 1
y 4

```

Conclusiones

En la primera demostración, no se encontraron palabras repetidas, ya que es un archivo pequeño y simple que yo mismo creé. Por otra parte, para la segunda demostración se utilizó un archivo de texto más extenso. En este caso, podemos ver que se encontraron varias palabras duplicadas.

Referencias

<https://es.wikipedia.org/wiki/MapReduce>

https://es.wikipedia.org/wiki/Apache_Hadoop

<http://blogs.solidq.com/es/big-data/que-es-mapreduce/>

https://www.tutorialspoint.com/es/hadoop/hadoop_mapreduce.htm