

Introducción a Hadoop y Spark

Procesamiento de Datos a Gran Escala

José Antonio Álvarez Ocete
Francisco Javier Sáez Maldonado

7 de octubre de 2021

Índice

1. Hadoop

1.1. Proceso

Instalamos Java

```
1 yum install java-1.8.0-openjdk-devel
```

Cambiamos la versión por defecto que se usaba en las transparencias. Incluso si instalamos la versión 1.7 se instalará la versión 1.8. Hemos de cambiar los paths a las versión 1.8 para que funcione correctamente:

```
1 export JAVA_HOME=/usr/lib/jvm/jre-1.8.0-openjdk
```

Editamos el archivo `/opt/hadoop/etc/hadoop/hadoop.env.sh`:

```
1 export JAVA_HOME=/usr/lib/jvm/jre-1.8.0-openjdk
```

Compilación

Tomamos el código que se ha proporcionado y lo compilamos utilizando la orden

```
1 bin/hadoop fs -cat /user/bigdata/compilar.bash | exec bash -s WordCount
```

Ejecución

Primero hay que iniciar el NameNode y el DataNode:

```
1 sbin/start-dfs.sh
```

Iniciar el ResourceManager y el NodeManager:

```
1 sbin/start-yarn.sh
```

Recuerda que hemos de subir el archivo utilizando:

```
1 /opt/hadoop/bin/hdfs dfs -put Quijote.txt /user/root
```

Lanzamos nuestro trabajo de MapReduce:

```
1 sudo /opt/hadoop/bin/hadoop jar WordCount.jar uam.WordCount Quijote.txt output/
```

Obtenemos en el directorio output la salida. En concreto, dos archivos:

- Un archivo `SUCCESS` indicando que la tarea ha sido exitosa.
- Un archivo `part-r-00000` que tiene la salida del programa que queríamos ejecutar.

Mostramos una parte del fichero para mostrar parte de la salida.

```
1      "Tablante",          1
2      "dichosa            1
3      "el                8
4      "y                  1
5      (Y                  1
6      (a                  1
7      (al                 1
8      (como               1
9      (creyendo           1
10     (de                  2
11     (por                  2
12     (porque              2
13     (pues                 1
14     (que                  21
```

Como podemos comprobar, las palabras quedan con ciertos símbolos de puntuación que no nos interesa que estén para realizar el conteo correcto de palabras.

Modificación del programa

Para arreglar esto, solo debemos cambiar una línea en la función `Map` del archivo `WordCount.java`. En concreto, la dejamos de la siguiente forma:

```
1 StringTokenizer itr = new StringTokenizer(value.toString()).
2     toLowerCase().replaceAll("[^a-z ]", "");
```

como vemos, hemos pasado las palabras a minúsculas usando `toLowerCase` y luego hemos eliminado todo aquello que no sean letras usando `replaceAll`.

Una vez realizada esta modificación, debemos compilar de nuevo el programa como lo hemos hecho anteriormente y, seguidamente, ejecutar de nuevo el programa java para que realice el conteo de palabras. En este caso, indicamos que la salida la realice sobre la carpeta `output2/` para tener las dos salidas por separado y poder compararlas. Obtenemos los dos mismos ficheros que anteriormente.

Hecho esto, utilizamos el archivo `script.py` que hemos creado usando Python, que toma todas las palabras que hay en los ficheros de salida y el número de veces que aparece cada una, las ordena y toma las 10 primeras para mostrarlas por pantalla.

El resultado que se obtiene es que las palabras son las mismas y prácticamente en el mismo orden, salvo una pequeña variación entre dos de ellas. Podemos verlo gráficamente en la siguiente figura:

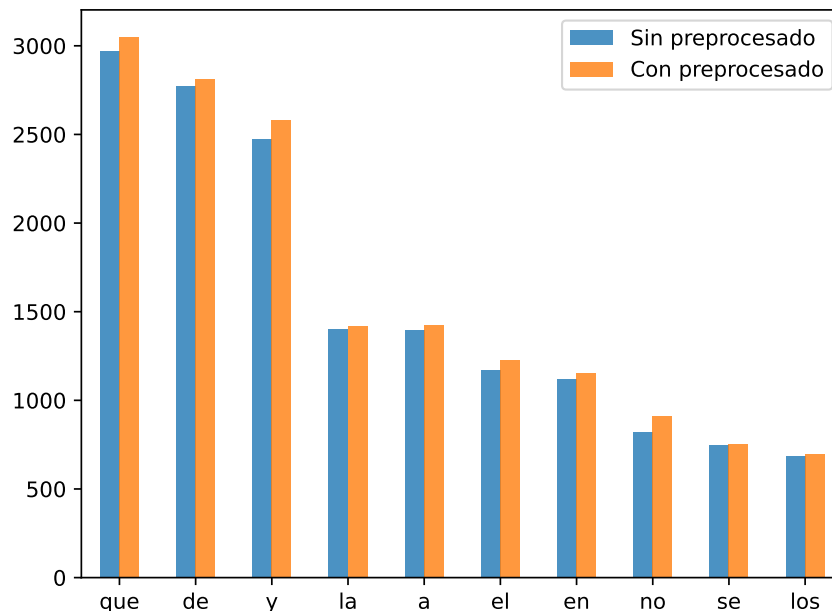


Figura 1: Comparación de número de veces que aparece cada palabra con y sin preprocesamiento.

Se puede apreciar que, al aplicar el preprocesado, la palabra *a* aparece un número de veces ligeramente superior al que lo hace cuando no tenemos el preprocesador previo del texto.

Comparación de resultados con otros ejemplos

TODO

1.2. Cuestiones planteadas

Pregunta 1.1. ¿ Qué ficheros ha modificado para activar la configuración del HDFS? ¿ Qué líneas ha sido necesario modificar?

Hemos modificado el fichero `/opt/hadoop-2.8.1/etc/hadoop/hadoop-env.sh` añadiendo la línea `export JAVA_HOME= /usr/lib/jvm/jre-1.7.0-openjdk` para especificar la instalación de Java que queremos utilizar.

Como se explica en [esta respuesta de Stack Overflow](#), el script `stop-all.sh` detiene todos los daemons de Hadoop a la vez, pero está obsoleto. En lugar de eso es recomendable parar los daemons de HDFS y YARN por separado en todas las máquinas utilizando `stop-dfs.sh` y `stop-yarn.sh`.

A continuación, para instalar Hadoop pseudo-distribuido, hemos modificado el fichero

```
1 /opt/hadoop/etc/hadoop/core-site.xml
```

del siguiente modo:

```
1 <configuration>
2   <property>
3     <name>fs.defaultFS</name>
4     <value>hdfs://localhost:9000</value>
5   </property>
6 </configuration>
```

Y añadimos al fichero /opt/hadoop/etc/hadoop/hdfs-site.xml lo siguiente:

```
1 <configuration>
2   <property>
3     <name>dfs.replication</name>
4     <value>1</value>
5   </property>
6 </configuration>
```

A continuación, hemos realizado la instalación del sistema pseudo-distribuido usando YARN, así que hemos modificado los siguientes ficheros **para configurar el uso de YARN (no de HDFS)**. etc/hadoop/mapred-site.xml:

```
1 <configuration>
2   <property>
3     <name>mapreduce.framework.name</name>
4     <value>yarn</value>
5   </property>
6 </configuration>
```

Y también etc/hadoop/yarn-site.xml:

```
1 <configuration>
2   <property>
3     <name>yarn.nodemanager.aux-services</name>
4     <value>mapreduce_shuffle</value>
5   </property>
6
7   <property>
8     <name>yarn.nodemanager.aux-services.mapreduce_shuffle.class</name>
9     <value>org.apache.hadoop.mapred.ShuffleHandler</value>
10  </property>
11 </configuration>
```

Ejercicio 1.2: Para pasar a la ejecución de Hadoop sin HDFS, ¿ es suficiente con parar el servicio con `stop-dfs.sh`? ¿ Cómo se consigue ?

TODO

Pregunta 3.1: ¿ Dónde se crea hdfs ? ¿ Cómo se puede decidir su localización ?

El sistema de archivos HDFS se crea donde la variable `dfs.datanode.data.dir` indique. Esta variable se puede modificar en el archivo `hdfs-site.xml`. Su valor por defecto es

```
1 file://${hadoop.tmp.dir}/dfs/data
```

y, podemos ver que la variable `hadoop.tmp.dir` tiene el valor `/tmp/hadoop-${user.name}`.

Pregunta 3.2: ¿ Cómo se puede borrar todo el contenido del HDFS, incluido su estructura ?

Pregunta 3.3: Si estás usando hdfs, ¿ cómo puedes volver a ejecutar WordCount como si fuese single.node ?

Recordamos que hemos hecho una serie de cambios en archivos `xml` para configurar Hadoop para que funcionase de forma pseudo-distribuida. Para ejecutarlo como si fuese *single-node*, deberíamos eliminar los cambios que hemos hecho en estos ficheros `xml`: `core-site.xml` y `hdfs-site.xml`. Así, volveríamos al modo por defecto de Hadoop y conseguiríamos que funcionase como si fuese *single-node*.

Pregunta 3.4: ¿ Cuáles son las 10 palabras más utilizadas ?

Esta pregunta ya ha sido mostrada anteriormente en la Figura 1. Las palabras más utilizadas son:

```
1 que, de, y, la, a, el, en, no, se, los
```

Pregunta 3.5: Cuántas veces aparecen las siguientes palabras: el, dijo

Para esta pregunta, volvemos a usar nuestro fichero `script.py` que nos imprimirá cuántas veces aparece cada una de ellas en el texto, primero usando preprocesado y a continuación sin utilizarlo.

```
1 Numero de veces que aparece la palabra
2   el
3       - Sin preprocesado 1173
4       - Con preprocesado 1228
5   dijo
6       - Sin preprocesado 196
7       - Con preprocesado 271
```

Puede comprobarse ejecutando el script indicado.

1.3. Ejercicio 4

Creamos un archivo con 15 copias del Quijote llamado `quijote15.txt`. Lo subimos al sistema de archivos de hdfs:

```
1 bin/hdfs dfs -put quijote15.txt /user/root/quijote15-128mb.txt
```

A continuación, volvemos a subir el archivo indicando el tamaño de bloque a utilizar, 2MB:

```
1 bin/hdfs dfs -D dfs.blocksize=2097152 -put quijote15.txt /user/root/quijote15-2mb
```

Por otro lado, también podemos cambiar el valor por defecto del tamaño de bloque para no tener que especificar el tamaño 2MB cada vez que subamos un archivo. Para ello cambiamos el respectivo parámetro (`dfs.block.size`) en los archivos de configuración de hdfs editando `hdfs-site.xml`, añadiendo:

```
1 <property>
2     <name>dfs.block.size</name>
3     <value>2097152</value>
4 </property>
```

Y reiniciamos dfs y yarn para asegurarnos que la configuración activa incluye los cambios realizados. Subimos de nuevo el archivo `quijote15.txt` sin especificar el tamaño de bloque utilizando la siguiente orden:

```
1 bin/hdfs dfs -put quijote15.txt /user/bigdata/dreji/quijote15-2mb-default.txt
```

A continuación accedemos a `http://localhost:50070/explorer.html#/user/root` para comprobar los resultados.

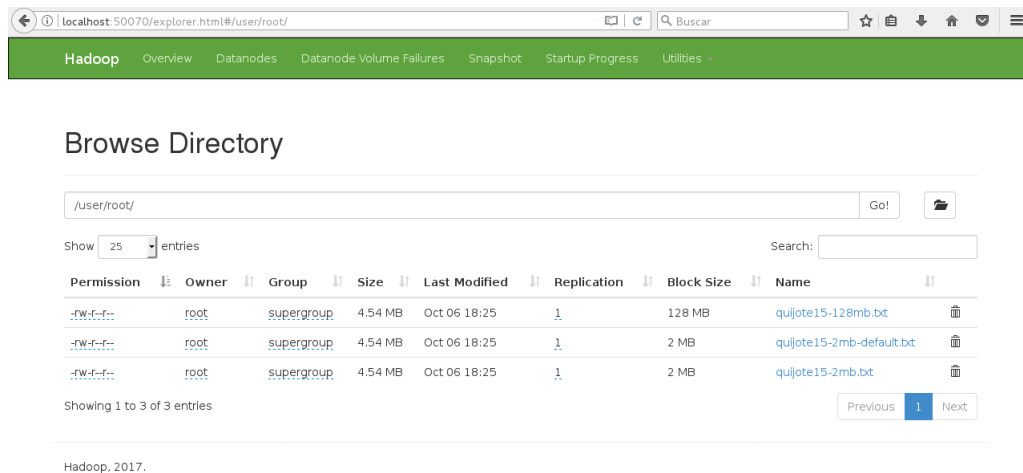


Figura 2: Captura del explorador del HDFS de hadoop

Podemos comprobar como la primera versión del archivo `quijote15.txt` subida tienen tamaño de bloque de 128MB, mientras que las dos últimas tienen tamaño de bloque 2MB. Esto responde a los apartados 4.1, 4.2 y 4.3 del guión.

Procedemos a comparar la ejecución de nuestro programa para distintos tamaños de bloque. Para ello ejecutamos el programa contador de palabras utilizando los archivos de entrada `quijote15-128mb.txt` y `quijote15-2mb.txt` subidos anteriormente. Utilizando 2MB como tamaño de bloque obtenemos la siguiente salida:

```
1 [...]
2 Shuffled Maps=3
3 Merged Map outputs=3
4 [...]
```

Indicando que se han realizado 3 operaciones *map* y 3 operaciones *reduce*. Este número de operaciones encaja con lo esperado, pues el archivo `quijote15.txt` pesa 4.54MB y se divide en tres bloques de 2MB. Al repetir la ejecución con 128MB para el tamaño de bloque obtenemos:

```
1      [...]
2      Shuffled Maps=1
3      Merged Map outputs=1
4      [...]
```

Obteniendo una única operación map y reduce, pues un único bloque de 128MB es suficiente para el archivo de 4.54MB. De esta forma hemos comprobado como afecta el tamaño de bloque en la ejecución de un programa MapReduce: Afecta al número de bloques creados y, por lo tanto, al número de operaciones map/reduce realizadas.

Las salidas de estas últimas ejecuciones pueden encontrarse en los archivos `output-quijote15-2mb.txt` y `output-quijote15-128mb.txt` respectivamente.

2. Tutorial de Spark

Seguimos el tutorial de spark y contestamos a las preguntas que se nos piden

Pregunta TS1.1 ¿Cómo hacer para obtener una lista de los elementos al cuadrado?

Bastará en este caso usar sobre el RDD la función map pasándole como parámetro la función `lambda x: x*x` que eleva el número al cuadrado. El código es este, con el resultado debajo:

```
1 numeros = sc.parallelize([1,2,3,4,5,6,7,8,9,10])
2 cuadrados = numeros.map(lambda x : x*x)
3 print(cuadrados.collect())
4
5 [1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
```

Pregunta TS1.2 ¿Cómo filtrar los impares?

El código se nos entrega resuelto en este caso. Es bastante sencillo, si usamos sobre el RDD la función `filter` dándole como parámetro la función que calcula si cada elemento es impar, esto es, si su módulo 2 es 1: `lambda x : x%2 == 1`, nos da todos los impares. El código completo es:

```
1 rddi = numeros.filter(lambda e: e%2==1)
2 print(rddi.collect())
3
4 [1, 3, 5, 7, 9]
```

Pregunta TS1.3 ¿Tiene sentido esta operación? ¿Si se repite se obtiene siempre el mismo resultado?

Nos encontramos con el siguiente código:

```
1 #Tiene sentido esta operacion?
2 numeros = sc.parallelize([1,2,3,4,5])
3
4 print(numeros.reduce(lambda elem1,elem2: elem1-elem2))
```

Esta operacion no tiene sentido, pues para hacer el reduce necesitamos que la operación sea conmutativa, es decir, que $a - b = b - a$, lo cual **NO** es cierto en todos los casos. Por eso además no producirá siempre los mismos resultados. Si ejecutamos en múltiples ocasiones veremos que los resultados no son los mismos siempre.

Pregunta TS1.4 ¿Cómo lo ordenarías para que primero aparezcan los impares y luego los pares?

Lo haremos para el caso general. Nos damos cuenta primero de que si el vector tiene n posiciones y tratamos de tomar $m > n$ posiciones, Spark se quedará con todas las posibles. Por tanto, lo primero que hacemos es obtener un entero con el máximo de elementos

```
1 numeros = sc.parallelize([3,2,1,4,5])
2
3 n = numeros.count()
```


Por último, basta ahora tomar elementos ordenados especificando cuál es el criterio de orden. En este caso, queremos que los elementos pares sean los más grandes, por lo que debemos pasarle como argumento la función `lambda x: x%2 == 0`. De esta manera, escribiendo la línea siguiente, obtenemos el resultado.

```
1 print(numeros.takeOrdered(n, lambda elem: elem%2== 0))
2
3 [3, 1, 5, 2, 4]
```

Pregunta TS1.5 ¿Cuántos elementos tiene cada rdd? ¿Cuál tiene más?

Nos encontramos con el siguiente código

```
1 lineas = sc.parallelize(['', 'a', 'a b', 'a b c'])
2
3 palabras_flat = lineas.flatMap(lambda elemento: elemento.split())
4 palabras_map = lineas.map(lambda elemento: elemento.split())
5
6 print (palabras_flat.collect())
7 print (palabras_map.collect())
```

Cuya salida es:

```
1 ['a', 'a', 'b', 'a', 'b', 'c']
2 [[], ['a'], ['a', 'b'], ['a', 'b', 'c']]
```

Para ver cuántos elementos tiene cada uno, podemos hacer

```
1 print(palabras_flat.count())
2 print(palabras_map.count())
3
4 6
5 4
```

Claramente, `palabras_flat` tiene más. Esto ocurre porque en `flatMap` cada elemento de entrada al que se le aplica la función `lambda` puede tener como salida 0 o más (un vector) de elementos. Por tanto:

- Con `flatMap` para cada elemento de `lineas` obtenemos un vector que luego separaremos por elementos e ignoraremos los elementos que estén vacíos. Por ello, de 4 elementos iniciales pasamos a 6 que es el total de letras que tenemos en el RDD.
- Con `map` a cada uno de los vectores iniciales se le aplica la función indicada y se introduce **tal cual** en el RDD resultante, no se separan los elementos de los vectores obtenidos como ocurría en el caso anterior.

Pregunta TS1.6 ¿De qué tipo son los elementos del rdd `palabras_map`? ¿Por qué `palabras_map` tiene el primer elemento vacío?

Siguiendo la respuesta de la pregunta anterior, en `palabras_map` los elementos obtenidos tras aplicar a cada elemento de `lineas` la función, son vectores, por lo que los elementos de `palabras_map` son vectores. Lo podemos ver en la salida del código

```

1 print (palabras_map.collect())
2
3 [[], ['a'], ['a', 'b'], ['a', 'b', 'c']]

```

Además, tiene el primer **elemento vacío** porque la función split aplicada sobre el elemento '' devuelve un vector vacío pues no hay nada que separar.

Pregunta TS1.7. Prueba la transformación distinct si lo aplicamos a cadenas.

Para probarlo, realizamos el siguiente código en el que ponemos en un RDD ejemplos de prueba con las mismas palabras en mayúscula y minúscula, con espacios o signos de puntuación:

```

1 test = sc.parallelize(["abcd", "abcd", "dcba", "abCd", "a bcd", "hola",
2                        "HOLA", "hola!", "HoLa", "ho la"])
3
4 dis = test.distinct()
5
6 print(dis.collect())

```

La salida que obtenemos es:

```

1 ['abcd', 'abCd', 'a bcd', 'hola', 'hola!', 'dcba', 'HOLA', 'HoLa', 'ho la']

```

Como podemos ver, teníamos un único ejemplo duplicado (el primero) que es el único que se ha eliminado. Esto nos indica que el comparador de elementos de los rdd comparan los strings *caracter a caracter*.

Pregunta TS1.8 ¿Cómo se podría obtener la misma salida pero utilizando una sola transformación y sin realizar la unión?

Se podría hacer que el filtro se quede con los elementos que empiecen por I o por E, haciendo el código del siguiente modo:

```

1 log = sc.parallelize(['E: e21', 'I: i11', 'W: w12', 'I: i11',
2                      'W: w13', 'E: e45'])
3
4 inferr = log.filter(lambda elem: elem[0] in ['I', 'E'])
5
6
7 print(inferr.collect())
8
9 ['E: e21', 'I: i11', 'I: i11', 'E: e45']

```

Pregunta TS1.9 ¿Cómo explica el funcionamiento de las celdas anteriores?

Comentamos una a una las celdas anteriores:

```

1 numeros = sc.parallelize([1,2,3,4,5])
2
3 print (numeros.reduce(lambda elem1,elem2: elem2+elem1))

```

Esta primera celda solo crea un RDD y luego hace la operación reduce, que combina los valores usando una función que le indicamos, en este caso la suma de los valores. Así, devolverá la suma de todos los valores ya que esta operación es conmutativa y asociativa.

```
1 #Tiene sentido esta operacion?
2 numeros = sc.parallelize([1,2,3,4,5])
3
4 print (numeros.reduce(lambda elem1,elem2: elem1+elem2))
```

Esta celda hace lo mismo que la anterior salvo que la operación que realiza es la resta. Como ya hemos visto en una pregunta anterior, esta operación no tiene sentido pues la resta no es una operación conmutativa por lo que el resultado no siempre será el mismo.

```
1 palabras = sc.parallelize(['HOLA', 'Que', 'TAL', 'Bien'])
2
3 pal_minus = palabras.map(lambda elemento: elemento.lower())
4
5 print (pal_minus.reduce(lambda elem1,elem2: elem1+"-"+elem2))
6 #y esta tiene sentido esta operacion?
7 # Que pasa si ponemos elem2+"-"+elem1
```

En esta celda se crea un RDD que tiene cadenas de caracteres, primero se pasan a minúsculas aplicándoles la función `.lower` y luego se concatenan todas las palabras a un solo string uniéndolas por guiones, obteniendo como salida:

```
1 hola-que-tal-bien
```

Si cambiásemos el orden en el que se concatenan los elementos como se nos indica en el comentario, la salida cambia y las palabras se van uniendo en el orden inverso, debido a cómo se realiza la operación reduce. La salida es:

```
1 bien-tal-que-hola
```

```
1 r = sc.parallelize([('A', 1), ('C', 4), ('A', 1), ('B', 1), ('B', 4)])
2 rr = r.reduceByKey(lambda v1,v2:v1+v2)
3 print (rr.collect())
```

En esta celda se crea un RDD que tiene como elementos Tuplas (`clave, valor`). A continuación, usando `reduceByKey` pasándole como argumento una función suma, lo que se hace es sumar los valores de las tuplas cuya clave sea la misma, obteniendo la salida esperada:

```
1 [('C', 4), ('A', 2), ('B', 5)]
```

La última celda es la siguiente:

```
1 r = sc.parallelize([('A', 1), ('C', 4), ('A', 1), ('B', 1), ('B', 4)])
2 rr1 = r.reduceByKey(lambda v1,v2:v1+v2)
3 print (rr1.collect())
4 rr2 = rr1.reduceByKey(lambda v1,v2:v1)
5 print (rr2.collect())
```

De nuevo, se crea un RDD cuyos elementos son tuplas (`clave, valor`), primero se realiza lo mismo que en la celda anterior, y luego se vuelve a aplicar sobre el RDD obtenido la función `reduceByKey` esta vez pasándole como función simplemente mantenerla clave que tiene. Sin embargo, esto no parece producir ningún efecto sobre el RDD, pues cuando se realiza la operación `collect` para ver los elementos del mismo, la salida es la misma que se produce en la celda anterior.

Como **conclusión** a esta pregunta, podemos decir que es importante cómo se aplica la función `reduce` sobre los RDD y que hay que tener cuidado con las operaciones que indicamos a esta función pues podrían no producir los resultados que se desean.

Pregunta `groupByKey`

Dada la siguiente celda

```
1 r = sc.parallelize([('A', 1), ('C', 2), ('A', 3), ('B', 4), ('B', 5)])
2 rr = r.groupByKey()
3 res= rr.collect()
4
5 print(rr.collect())
6
7
8 # Que operacion realizar al RDD rr para que la operacion sea como un reduceByKey
```

Lo primero que vemos es que la operación `rr.collect` devuelve una lista con tuplas (`clave, pyspark Re`). Estos iterables debemos pasarlos primero a una lista, y luego realizar sobre ellos la operación que quisiésemos hacer con el `reduceByKey`. Vemos que los pasamos a una lista utilizando `mapValues(list)` que convierte los valores a un tipo:

```
1 rrf = rr.mapValues(list).collect()
2 print(rrf)
3
4 [('C', [2]), ('A', [1, 3]), ('B', [4, 5])]
```

Y, a continuación, podemos realizar la operación que haríamos con el `reduce`. Por ejemplo, si queremos hacer la suma de los vectores utilizando `map` podríamos hacer todo en una línea de la siguiente manera:

```
1 rrf = rr.mapValues(list).map(lambda x : (x[0], sum(x[1])))
2 print(rrf.collect())
3
4 [('C', 2), ('A', 4), ('B', 9)]
```

Ahora, se nos pide simular el `groupByKey` usando `reduceByKey` y `map`. Para ello, usando `reduceByKey` podemos obtener las listas que obtendríamos tras aplicar el `mapValues(list)` que hemos obtenido anteriormente.

```
1 simul_group = r.reduceByKey(lambda v1,v2: [v1,v2])
2 print(simul_group.collect())
3
4 [('C', 2), ('A', [1, 3]), ('B', [4, 5])]
```

Pregunta sobre join

```
1 rdd1 = sc.parallelize([('A',1), ('B',2), ('C',3)])
2 rdd2 = sc.parallelize([('A',4), ('B',5), ('C',6)])
3
4 rddjoin = rdd1.join(rdd2)
```

El resultado de esto es

```
1 [('A', (1, 4)), ('B', (2, 5)), ('C', (3, 6))]
```

Es decir, un nuevo RDD en el que los valores son tuplas con los valores que hay en cada uno de los RDD. Se nos pide que, dado ese código, cambiemos las claves de los dos RDDs iniciales para ver qué RDD se crea finalmente. Si lo hacemos, cambiando por ejemplo los RDD del siguiente modo:

```
1 rdd1 = sc.parallelize([('A',1), ('B',2), ('D',3)])
2 rdd2 = sc.parallelize([('A',4), ('B',5), ('E',6)])
3
4 rddjoin = rdd1.join(rdd2)
```

El resultado obtenido es el siguiente:

```
1 [('A', (1, 4)), ('B', (2, 5))]
```

Es decir, que el `join` está haciendo la unión de los elementos cuyas claves están en la intersección del conjunto de claves. Como `D` y `E` no están en ambos RDD, no están en la intersección del conjunto de claves y por tanto no se obtienen en el RDD final.

Tipos de Join

Se nos pregunta que qué ocurre cuando sustituimos `join` por `leftOuter/rightOuter/fullOuter join`. El resultado es que estos tipos de `join` crean elementos en el nuevo RDD aunque sus claves no estén en ambos RDD iniciales. En concreto:

- `leftOuter` añade al nuevo RDD también las tuplas (`clave, valor`) cuya clave esté en el RDD sobre el que se llama la función, pero no estén en el RDD que se pasa como parámetro. Se añade un `None` en la tupla conjunta del RDD final:

```
1 rdd1 = sc.parallelize([('A',1), ('B',2), ('C',3)])
2 rdd2 = sc.parallelize([('A',4), ('A',5), ('B',6), ('D',7)])
3 rddjoin = rdd1.leftOuterJoin(rdd2)
4
5 [('A', (1, 4)), ('A', (1, 5)), ('B', (2, 6)), ('C', (3, None))]
```

- `rightOuter` hace lo mismo que el anterior pero en el sentido opuesto, es decir, usando el segundo RDD.
- `fullOuter` combina los dos anteriores.

Pregunta TS1.10 Borra la salida y cambia las particiones en parallelize. ¿ Qué sucede ?

Lo que ocurre si ponemos 10 particiones es lo siguiente (borrando anteriormente el contenido del directorio donde tenemos la salida):

```
1 numeros = sc.parallelize(range(0,1000),10)
2 numeros.saveAsTextFile('salida')
3
4 %ls -la salida/*
5
6
7
8 -rw-r--r-- 1 root root 290 Oct  4 10:42 salida/part-00000
9 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00001
10 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00002
11 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00003
12 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00004
13 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00005
14 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00006
15 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00007
16 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00008
17 -rw-r--r-- 1 root root 400 Oct  4 10:42 salida/part-00009
18 -rw-r--r-- 1 root root  0 Oct  4 10:42 salida/_SUCCESS
```

Como se podía esperar, se crea un archivo para cada una de las particiones de salida, según el número de particiones que le hayamos indicado.

2.1. Procesamiento El Quijote

Pregunta TS2.1 Explica la utilidad de cada transformación y detalle para cada una de ellas si cambia el número de elementos en el RDD resultante. Es decir si el RDD de partida tiene N elementos, y el de salida M elementos, indica si $N > M$, $N = M$ o $N < M$.

Dado el siguiente código:

```
1 charsPerLine = quijote.map(lambda s: len(s))
2 allWords = quijote.flatMap(lambda s: s.split())
3 allWordsNoArticles = allWords.filter(lambda a: a.lower() not in ["el", "la"])
4 allWordsUnique = allWords.map(lambda s: s.lower()).distinct()
5 sampleWords = allWords.sample(withReplacement=True, fraction=0.2, seed=666)
6 weirdSampling = sampleWords.union(allWordsNoArticles.sample(False, fraction=0.3))
```

Se nos pide indicar qué hace cada una de las transformaciones. Procedemos línea a línea:

1. `charsPerLine = quijote.map(lambda s: len(s))` esta línea obtiene un nuevo RDD en el que se cambian las líneas del quijote por la longitud de cada línea. No se modifica el número de elementos, se tiene $N = M$.

2. `allWords = quijote.flatMap(lambda s: s.split())` obtiene un nuevo RDD en el que se separa cada uno de los strings iniciales en cada una de sus palabras y luego cada palabra se convierte en un elemento del RDD. Por tanto, el RDD de salida tiene muchos más elementos que el de entrada, podemos decir que $M \gg N$.
3. `allWordsNoArticles = allWords.filter(lambda a: a.lower() not in ["el", "la"])` se toma el RDD de las palabras y se pasa un filtro que elimina todos los artículos `*el*` y `*la*` (sin distinción de mayúsculas/minúsculas) del documento. Al eliminar elementos, el número de elementos en la salida será menor que en la entrada por lo que $M < N$.
4. `allWordsUnique = allWords.map(lambda s: s.lower()).distinct()` se crea un RDD que tiene un elemento por cada palabra distinta (sin distinción mayúsculas/minúsculas) que haya en el RDD `allWords`. De nuevo, se eliminan elementos repetidos por lo que el número de elementos en la salida será menor que en la entrada, $M < N$.
5. `sampleWords = allWords.sample(withReplacement=True, fraction=0.2, seed=666)` se extrae aleatoriamente con reemplazamiento un 20 % de las palabras que hay en `allWords` para crear un nuevo RDD. Por supuesto, se tendrá que $M = 0,2N$ ($M < N$).
6.

```
weirdSampling = sampleWords.union
2 (allWordsNoArticles.sample(False, fraction=0.3))
```

primero extrae sin reemplazamiento un 30 % de elementos que tenemos en el RDD que no tiene artículos (en este caso, tendríamos $M < N$ pues reducimos elementos) y a continuación une estos elementos extraídos al RDD que tiene una muestra del 20 % que hemos obtenido en el paso anterior. En esta unión aumenta el tamaño del RDD por lo que en este paso $M > N$.

Ahora, explicamos las funciones en general, su utilidad y si cambia el número de elementos del RDD resultante:

- `map` sirve para realizar una transformación sobre el RDD pasándole como parámetro una función. El número de elementos final del RDD es el mismo, pues a cada elemento se le aplica una función pero sigue siendo un elemento, no se elimina ninguno, por lo que $M = N$.
- `flatMap` realiza una transformación y, a continuación, si tenemos vectores como puntos, hace que cada elemento de los vectores sea un punto del nuevo RDD, por lo que en este caso el número se mantiene (si cada vector de salida tiene dimensión 1), o aumenta, por lo que $M \geq N$.
- `filter` selecciona un conjunto de elementos según una función que da una condición implícita. Todos los elementos pueden cumplir esa condición o solo un subconjunto de ellos, por lo que $M \leq N$.
- `distinct` toma los elementos del RDD que sean únicos (diferentes a todos los demás). Todos podrían ser diferentes o podría haber elementos iguales, por lo que $M \leq N$.
- `sample` toma una muestra del tamaño indicado. En este caso, el RDD no se modifica solo que se toman ejemplos del mismo, por lo que podemos decir que el tamaño del RDD se mantiene, $M = N$.

- `union` une dos RDDs en uno

Pregunta TS2.2 Explica el funcionamiento de cada acción anterior

Tenemos el siguiente código:

```
1 numLines = quijote.count()
2 numChars = charsPerLine.reduce(lambda a,b: a+b) # also charsPerLine.sum()
3 sortedWordsByLength = allWordsNoArticles.takeOrdered(20, key=lambda x: -len(x))
4 numLines, numChars, sortedWordsByLength
```

Este código hace lo siguiente

1. `numLines = quijote.count()` cuenta el número de líneas que tiene el quijote
2. `numChars = charsPerLine.reduce(lambda a,b: a+b)` utiliza un RDD anterior que tiene el número de caracteres que tiene cada línea y los suma, obteniendo así el número de caracteres total.

3. La siguiente línea

```
1 sortedWordsByLength = allWordsNoArticles.takeOrdered(20, key=lambda x: -len(x))
```

toma el RDD que tiene todas las palabras y toma las 20 palabras que tienen MAYOR longitud

A continuación, se pide que se implemente la acción `count` que nos dice el número de elementos que tiene un RDD usando lo siguiente:

- Usando `map` y `reduce`. Lo podemos hacer de la siguiente forma

```
1 numlines = quijote.map(lambda s: 1)
2 total_lines = numlines.reduce(lambda a,b: a+b )
3
4 5534
```

Hacemos un `map` para transformar cada línea en un 1, y luego los sumamos todos con `reduce`.

- Usando solo `reduce`. Lo podemos hacer pero debemos definir una función para distinguir si los elementos que se van a sumar son strings o son ya parte de la suma total. El código completo es el siguiente:

```
1 def f(a,b):
2     int_a = a if type(a) == int else 1
3     int_b = b if type(b) == int else 1
4
5     return int_a + int_b
6
7 total_lines = quijote.reduce(f)
8
9 5534
```


Pregunta TS2.3 Explica el proposito de cada una de las operaciones anteriores

```
1 import requests
2 import re
3 allWords = allWords.flatMap(lambda w:
4     re.sub("";|:|\.|,|-|-|'|'\s"", " ", w.lower()).split(" "))
5     .filter(lambda a: len(a)>0)
6 allWords2 = sc.parallelize(
7     requests.get("[url]el_quijote_ii.txt").iter_lines())
8 allWords2 = allWords2.flatMap(lambda w:
9     re.sub("";|:|\.|,|-|-|'|'\s"", " ", w.decode("utf8").lower())
10    .split(" ")).filter(lambda a: len(a)>0)
```

Esta celda se dedica a importar la librería `requests` que sirve para hacer peticiones web, y la librería `re` para trabajar con expresiones regulares, y a continuación crea dos RDD diferentes: uno en el que elimina todos los caracteres que no forman parte de palabras como tal, y otro en el que lee la segunda parte del Quijote, obtiene todas sus líneas y luego las separa y le quita los caracteres que no son parte de las palabras igual que hizo con la primera parte.

```
1 allWords.take(10)
2 allWords2.take(10)
```

Esta celda simplemente muestra los 10 primeros elementos de cada uno de los RDD creados.

```
1 words = allWords.map(lambda e: (e,1))
2 words2 = allWords2.map(lambda e: (e,1))
3
4 words.take(10)
```

En esta celda se transforman los elementos de cada RDD en tuplas (`clave, valor`) cada una con un 1 como valor. Luego, se muestran las 10 primeras de uno de los RDD.

```
1 frequencies = words.reduceByKey(lambda a,b: a+b)
2 frequencies2 = words2.reduceByKey(lambda a,b: a+b)
3 frequencies.takeOrdered(10, key=lambda a: -a[1])
```

Esta celda usa `reduceByKey` para crear un RDD en el que se suman los valores de los elementos que tienen la misma clave. En este caso, lo que se está haciendo es crear un RDD en el que se ha contado cuántas veces aparece cada palabra en el texto. Luego, se muestran las 10 palabras que más salen usando `takeOrdered` con una función que ordena los elementos de mayor a menor.

```
1 res = words.groupByKey().takeOrdered(10, key=lambda a: -len(a))
2 res # To see the content, res[i][1].data
```

En esta celda se agrupan en listas los valores que compartan clave. Es decir, si tenemos dos parejas (`palabra1, 1`), (`palabra1, 1`), obtendremos como resultado (`palabra1, [1, 1]`). El vector en realidad queda como un `ResultIterable` de `pyspark`, pero podemos convertirlo a una lista como hicimos anteriormente.

```

1 joinFreq = frequencies.join(frequencies2)
2 joinFreq.take(10)

```

En esta celda, creamos un nuevo RDD que tiene (palabra, (valor1, valor2)) donde valor1 es la frecuencia de esa palabra en el primer texto del Quijote y valor2 es la frecuencia de esa palabra en el segundo texto. Luego, se muestran 10.

```

1 joinFreq.map(lambda e: (e[0], (e[1][0] - e[1][1])/(e[1][0] + e[1][1])))
2   .takeOrdered(10, lambda v: -v[1]),
3 joinFreq.map(lambda e: (e[0], (e[1][0] - e[1][1])/(e[1][0] + e[1][1])))
4   .takeOrdered(10, lambda v: +v[1])

```

En esta celda, se transforma el RDD anterior para convertirlo en un RDD que tiene (palabra, valor). Ahora, sin embargo, el número que toma valor es

$$\frac{\text{Apariciones palabra texto 1} - \text{Apariciones palabra texto 2}}{\text{Apariciones palabra texto 1} + \text{Apariciones palabra texto 2}}$$

Indica, respecto al total de veces que aparece la palabra en ambos textos, cuántas veces más aparece en el primer texto que en el segundo. Es decir, si en el primero aparece muchas más que en el segundo, este valor será más cercano a 1, y si en el segundo aparece más veces, este valor será cercano a -1. Así, podemos ver qué palabras aparecen más en el primero que en el segundo y viceversa. Eso es lo que hace la celda después usando takeOrdered.

Pregunta TS2.4 ¿Cómo puede implementarse la frecuencia con groupByKey y transformaciones?

Podemos realizarlo usando el siguiente código, cuyos comentarios indican qué se está haciendo en cada paso:

```

1 # Usamos GroupByKey y sobre la salida, mapeamos cada elemento a
2 # elemento = (clave, [1,...,])
3 res = words.groupByKey().map(lambda x : (x[0], list(x[1])))
4 # Mapeamos de nuevo los elementos a (elemento, suma([1,...,1]))
5 res = res.map(lambda x: (x[0], sum(x[1])))
6 res.take(5)
7
8 [('el', 1232), ('hidalgo', 14), ('don', 370), ('mancha', 26), ('saavedra', 1)]

```

Pregunta TS2.5 ¿Cuál de las dos siguientes celdas es más eficiente? Justifique la respuesta.

Se nos presentan esta primera celda:

```

1 joinFreq.map(lambda e: (e[0], (e[1][0] - e[1][1])/(e[1][0] + e[1][1])))
2   .takeOrdered(10, lambda v: -v[1]),
3 joinFreq.map(lambda e: (e[0], (e[1][0] - e[1][1])/(e[1][0] + e[1][1])))
4   .takeOrdered(10, lambda v: +v[1])

```

Y esta segunda

```

1 result = joinFreq.map(lambda e: (e[0], (e[1][0] - e[1][1])/(e[1][0] + e[1][1])))
2 result.cache()
3 result.takeOrdered(10, lambda v: -v[1]), result.takeOrdered(10, lambda v: +v[1])

```

Y queremos saber cual de las dos es más eficiente. En este caso, la respuesta es sencilla:

La segunda celda es más eficiente, pues no estamos realizando dos veces la transformación de cada uno de los elementos (con sus correspondientes sumas, restas y divisiones) sino que estamos haciendo el RDD persistente y luego de ese mismo, tomando los elementos según un orden u otro. Por eso, deberíamos utilizar en general siempre la segunda opción.

Además, si nos fijamos en el [enlace a la persistencia de los RDD](#), nos indica que si lo guardamos en caché las acciones reutilizan el mismo RDD y así las operaciones pueden llegar a ir hasta 10 veces más rápido.

Antes de guardar el fichero, utilice coalesce con diferentes valores ¿Cuál es la diferencia?

Vamos a probar esto. Primero, en la [documentación de coalesce](#) leemos que sirve para devolver un nuevo RDD hecho en menos particiones. Vemos también que, en principio, el RDD que vamos a usar tiene 2 particiones:

```
1 allWords2.getNumPartitions()
2
3 2
```

A continuación, intentamos usar `coalesce` para modificar estas particiones incrementándolas.

```
1 for i in range(1,6):
2     print("Usamos coalesce con numpartitions={}".format(i))
3     !rm -rf palabras_parte2
4     allWords2 = allWords2.coalesce(numPartitions=i)
5     .saveAsTextFile("palabras_parte2")
6     !ls palabras_parte2
7
8 Usamos coalesce con numpartitions=1
9   part-00000  _SUCCESS
10 Usamos coalesce con numpartitions=2
11   part-00000  part-00001  _SUCCESS
12 Usamos coalesce con numpartitions=3
13   part-00000  part-00001  _SUCCESS
14 ...
```

Vemos que, a partir de 2 particiones, el número de particiones que se ven reflejadas en disco no aumenta más, esto ocurre porque `coalesce` con esos parámetros solo sirve para reducir. Podríamos usar como parámetro `shuffle = True`, lo cual haría que se barajasen de nuevo los elementos y así sí que se pueden crear nuevas particiones. La salida del código anterior poniendo este parámetro es:

```
1 Usamos coalesce con numpartitions=1
2 part-00000  _SUCCESS
3 Usamos coalesce con numpartitions=2
4 part-00000  part-00001  _SUCCESS
5 Usamos coalesce con numpartitions=3
6 part-00000  part-00001  part-00002  _SUCCESS
7 Usamos coalesce con numpartitions=4
```

```
8 part-00000 part-00001 part-00002 part-00003 _SUCCESS
9 Usamos coalesce con numpartitions=5
10 part-00000 part-00001 part-00002 part-00003 part-00004
11 _SUCCESS
```

Sin embargo, esto sería equivalente a usar la función `repartition` que hemos visto anteriormente, que como se nos indica no tiene en cuenta que los datos tengan que moverse en el disco para rehacer las particiones.