

Movies data analysis using MapReduce

Tushar B. Kute,
<http://tusharkute.com>

What is MapReduce?

- MapReduce is a framework using which we can write applications to process huge amounts of data, in parallel, on large clusters of commodity hardware in a reliable manner.
- MapReduce is a processing technique and a program model for distributed computing based on java.
- The MapReduce algorithm contains two important tasks, namely Map and Reduce.

Map and Reduce

- Map takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key/value pairs).
- Secondly, reduce task, which takes the output from a map as an input and combines those data tuples into a smaller set of tuples. As the sequence of the name MapReduce implies, the reduce task is always performed after the map job.

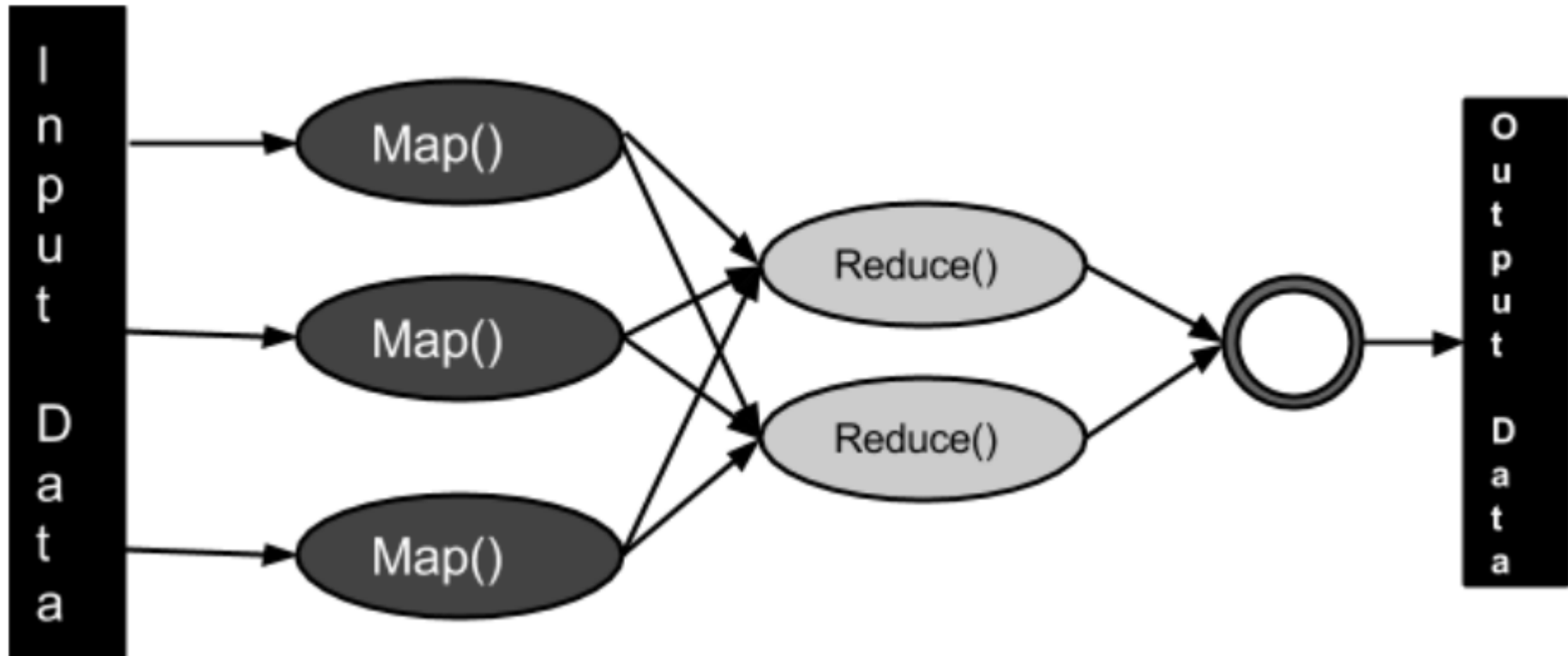
Map and Reduce

- The major advantage of MapReduce is that it is easy to scale data processing over multiple computing nodes.
- Under the MapReduce model, the data processing primitives are called mappers and reducers.
- Decomposing a data processing application into mappers and reducers is sometimes nontrivial. But, once we write an application in the MapReduce form, scaling the application to run over hundreds, thousands, or even tens of thousands of machines in a cluster is merely a configuration change.
- This simple scalability is what has attracted many programmers to use the MapReduce model.

The Algorithm

- MapReduce program executes in three stages, namely map stage, shuffle stage, and reduce stage.
- **Map stage:** The map or mapper's job is to process the input data. Generally the input data is in the form of file or directory and is stored in the Hadoop file system (HDFS). The input file is passed to the mapper function line by line. The mapper processes the data and creates several small chunks of data.
- **Reduce stage:** This stage is the combination of the Shuffle stage and the Reduce stage. The Reducer's job is to process the data that comes from the mapper. After processing, it produces a new set of output, which will be stored in the HDFS.

The MapReduce



Inserting Data into HDFS

- The MapReduce framework operates on $\langle \text{key}, \text{value} \rangle$ pairs, that is, the framework views the input to the job as a set of $\langle \text{key}, \text{value} \rangle$ pairs and produces a set of $\langle \text{key}, \text{value} \rangle$ pairs as the output of the job, conceivably of different types.
- The key and the value classes should be in serialized manner by the framework and hence, need to implement the Writable interface. Additionally, the key classes have to implement the Writable-Comparable interface to facilitate sorting by the framework.
- Input and Output types of a MapReduce job: (Input) $\langle k_1, v_1 \rangle \rightarrow \text{map} \rightarrow \langle k_2, v_2 \rangle \rightarrow \text{reduce} \rightarrow \langle k_3, v_3 \rangle$ (Output).

Data input and output

	Input	Output
Map	$\langle k1, v1 \rangle$	list ($\langle k2, v2 \rangle$)
Reduce	$\langle k2, \text{list}(v2) \rangle$	list ($\langle k3, v3 \rangle$)

Terminologies

- Mapper - Mapper maps the input key/value pairs to a set of intermediate key/value pair.
- NamedNode - Node that manages the Hadoop Distributed File System (HDFS).
- DataNode - Node where data is presented in advance before any processing takes place.
- MasterNode - Node where JobTracker runs and which accepts job requests from clients.
- SlaveNode - Node where Map and Reduce program runs.
- JobTracker - Schedules jobs and tracks the assign jobs to Task tracker.
- Task Tracker - Tracks the task and reports status to JobTracker.
- Job - A program is an execution of a Mapper and Reducer across a dataset.
- Task - An execution of a Mapper or a Reducer on a slice of data.

Example:

- Use movies dataset. Write a map and reduce methods to determine the average rating of movies. The input consists of series of lines, each containing movie number, user number, rating and timestamp. The map should emit movie number and list of ratings and reduce should return the average rating for each movie number.

The dataset:u.data

1	196	242	3	881250949
2	186	302	3	891717742
3	22	377	1	878887116
4	244	51	2	880606923
5	166	346	1	886397596
6	298	474	4	884182806
7	115	265	2	881171488
8	253	465	5	891628467
9	305	451	3	886324817
10	6	86	3	883603013
11	62	257	2	879372434
12	286	1014	5	879781125
13	200	222	5	876042340
14	210	40	3	891035994
15	224	29	3	888104457
16	303	785	3	879485318
17	122	387	5	879270459
18	194	274	2	879539794

Example:

- Movies.java

Compilation and Execution

- Let us assume we are in the home directory of a Hadoop user (e.g. /home/rashmi).
- Follow the steps given below to compile and execute the above program.
- **Step 1**
 - The following command is to create a directory to store the compiled java classes.
 - `$ mkdir movies`

Compilation and Execution

- **Step 2**

Download [hadoop-core-1.2.1.jar](http://mvnrepository.com/artifact/org.apache.hadoop/hadoop-core/1.2.1), which is used to compile and execute the MapReduce program. Visit the following link

<http://mvnrepository.com/artifact/org.apache.hadoop/hadoop-core/1.2.1>

To download the jar. Let us assume the downloaded folder is /home/rashmi/movies.

- **Step 3**

The following commands are used for compiling the wordcount.java program and creating a jar for the program.

```
$ javac -classpath hadoop-core-1.2.1.jar movies/Movies.java
```

```
$ jar -cvf snow.jar -C movies/ .
```

Compilation and Execution

- **Step 4**
 - The following command is used to create an input directory in HDFS.
 - `$hadoop fs -mkdir /input`
- **Step 5**
 - The following command is used to copy input dataset file on HDFS.
 - `$hadoop fs -put u.data /input`
- **Step 6**
 - The following command is used to verify the files in the input directory.
 - `$hadoop fs -ls /input`

Compilation and Execution

- **Step 7**
 - The following command is used to run the Snow application by taking the input files from the input directory.
 - `$hadoop jar movies.jar Movies /input /output`
 - Wait for a while until the file is executed. After execution, the output will contain the number of input splits, the number of Map tasks, the number of reducer tasks, etc. The output directory must *not* be existing already.

Compilation and Execution

- **Step 8**

- The following command is used to verify the resultant files in the output folder.
- `$hadoop fs -ls /output`

- **Step 9**

- The following command is used to see the output in part-r-00000 file. This file is generated by HDFS.
- `$hadoop fs -cat /output/part-r-00000`

Compilation and Execution

- **Step 10**

The following command is used to copy the output file from HDFS to the local file system for analyzing.

- `$hadoop fs -get /output/part-r-00000`

Output:

1	3
2	3
3	2
4	4
5	2
6	3
7	3
8	3
9	4
10	4
11	3
12	4
13	3
14	4
15	2

Thank you

This presentation is created using LibreOffice Impress 5.1.6.2, can be used freely as per GNU General Public License



@mitu_skillologies



/mITuSkillologies



@mitu_group

Web Resources

<http://mitu.co.in>

<http://tusharkute.com>

tushar@tusharkute.com

contact@mitu.co.in