

## **TP : MapReduce, YARN et API Java**

### **Résumé**

Dans ce TP nous allons écrire des programmes Java afin d'extraire des informations synthétiques à partir d'un cluster Hadoop.

L'objectif général est d'écrire des programmes qui implémentent le patron d'architecture MapReduce afin d'extraire des données stockées sur le système de fichiers HDFS d'un cluster. Ce TP présente l'API Java qui permet de réaliser des traitements qui s'exécutent parallèlement sur les nœuds du cluster Hadoop.

### **1. Introduction**

MapReduce est un Framework de traitement de données en clusters créé par Google. Son utilité est de permettre l'écriture de programmes de traitement et de génération de larges ensembles de données sur un cluster de machines. MapReduce est un composant central de Hadoop permettant le traitement résilient et distribué d'ensembles de données massives sur des clusters d'ordinateurs, au sein desquels chaque nœud possède son propre espace de stockage. Concrètement MapReduce propose deux fonctionnalités principales. Il répartit le travail sur les différents nœuds du cluster (map), puis les organise et réduit les résultats fournis par chaque nœud en une seule réponse qui correspond à la requête (reduce).

Afin de gérer les traitements écrits selon le patron d'architecture MapReduce, Hadoop offre son composant YARN. Ce dernier, permet aux programmeurs de lancer des traitements (jobs) sur des données du système de fichiers HDFS, et de suivre leurs avancements.

### **2. Exemple de traitement MapReduce**

Dans cet exemple nous allons effectuer des traitements afin d'extraire des informations synthétiques sur un fichier de données stocké sur le HDFS. Le fichier de données à utiliser contient des informations sur des ventes. La structure du fichier est la suivante :

Date	Prod	Prix	Paiement	Nom	Ville	Etat	Pays	Compte	DernConn

**L'objectif est de trouver le nombre des produits vendus dans chaque pays.**

- 1) **Le fichier des données de travail est déjà mis sur le HDFS.** Son chemin est le suivant :  
`storage/Sales/SalesJan2009.csv`
- 2) Lancer votre IDE et créer un nouveau projet avec le nom « Vente ». Il faut enregistrer ce projet dans le répertoire « TP\_Hadoop » de la machine **nodeX**.
- 3) Dans le projet « Vente » écrire les trois classes dont les codes sont indiqués dans les pages ci-après.
- 4) Afin de compiler et lancer le projet un script dédié pour la compilation et l'exécution du code Java est déjà mis en place. Vous pouvez le télécharger sur l'url suivant (**Attention ! : il faut le télécharger au sein de votre projet**):  
`$ wget https://azammouri.com/min-5a/Makefile`
- 5) Lancer la compilation et l'exécution via la commande :  
`$ make`

## Le Mapper

### VenteMapper.java

```
import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class VenteMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    private final static IntWritable one = new IntWritable(1);
    private Text Country = new Text();

    public void map(LongWritable Key, Text value, Mapper.Context context)
    throws IOException, InterruptedException{

        String InputLine = value.toString();

        String[] Splited_InputLine = InputLine.split(",");

        Country.set(Splited_InputLine[7]);
        context.write(Country, one);
    }
}
```

Pour chaque appel de la méthode map() une paire clé-valeur est passée en arguments. Le traitement de la méthode map() commence par diviser le texte qui est reçu comme entrée de la méthode.

Afin de diviser le texte, nous avons utilisé la méthode « split ». Le caractère selon lequel la séparation est faite est : « , ». Par la suite une paire de (clé, valeur) est construite. Cette paire consiste à associer à chaque ville une valeur de « 1 ».

## Le Reducer

### **VenteReducer.java**

```
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class VenteReducer extends Reducer<Text, IntWritable, Text, IntWritable>
{
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values, Context
context) throws IOException, InterruptedException {

        int Sum=0;
        for(IntWritable val: values) {
            Sum = Sum + val.get();
        }

        result.set(Sum);

        context.write(key, result);
    }
}
```

Les deux premiers types de données de la fonction `reduce()` sont « Text » et « IntWritable » qui représentent respectivement les types de la clé et la valeur.

La suite du mapper est de la forme : `<Nom_Pays, 1>`. Cette sortie du mapper devient une entrée du reducer. D'où l'utilisation des types « Text » et « IntWritable ».

Les deux derniers arguments de la fonction `reduce()` sont du type « Text » et « IntWritable » respectivement. Ceci correspond à la clé et à la valeur de sortie de la méthode `reduce()`.

## Le Driver

### VenteDriver.java

```
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class VenteDriver {

    public static void main(String[] args) throws Exception {

        Configuration Conf = new Configuration();
        Job Traitement = Job.getInstance(Conf, "Total Ventes par Pays");

        Traitement.setJarByClass(VenteDriver.class);
        Traitement.setMapperClass(VenteMapper.class);
        Traitement.setReducerClass(VenteReducer.class);
        Traitement.setOutputKeyClass(Text.class);
        Traitement.setOutputValueClass(IntWritable.class);

        FileInputFormat.addInputPath(Traitement, new
Path("storage/Sales/SalesJan2009.csv"));

        FileOutputFormat.setOutputPath(Traitement, new
Path("storage/Sales/Resultats_Traitement_1"));

        System.exit(Traitement.waitForCompletion(true) ? 0 : 1);
    }
}
```

Si tout se passe bien, l'exécution de la commande « make » donnera le résultat présenté sur la figure ci-après :

```

hadoop@node-master: ~/TP_Hadoop/Ventes
hadoop@node-master:~/TP_Hadoop/Ventes$ make
javac -cp /home/hadoop/TP_Hadoop/hadoop-common-3.1.2.jar:/home/hadoop/TP_Hadoop/hadoop-mapreduce-client-core-3.1.2.jar:/home/hadoop/TP_Hadoop/hadoop-hdfs-3.1.2.jar src/VenteDriver.java src/VenteReducer.java src/VenteMapper.java -d bin
Note: src/VenteMapper.java uses unchecked or unsafe operations.
Note: Recompile with -Xlint:unchecked for details.
jar cfe projet.jar VenteDriver -C bin ./
hadoop jar projet.jar
2020-12-18 22:34:28,244 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2020-12-18 22:34:29,815 INFO client.RMProxy: Connecting to ResourceManager at node-master/192.168.0.22:8032
2020-12-18 22:34:30,649 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2020-12-18 22:34:30,693 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/hadoop/.staging/job_1608150350331_0002
2020-12-18 22:34:31,240 INFO input.FileInputFormat: Total input files to process : 1
2020-12-18 22:34:31,436 INFO mapreduce.JobSubmitter: number of splits:1
2020-12-18 22:34:31,841 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1608150350331_0002
2020-12-18 22:34:31,843 INFO mapreduce.JobSubmitter: Executing with tokens: []
2020-12-18 22:34:32,317 INFO conf.Configuration: resource-types.xml not found
2020-12-18 22:34:32,322 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'
2020-12-18 22:34:32,516 INFO impl.YarnClientImpl: Submitted application application_1608150350331_0002
2020-12-18 22:34:32,632 INFO mapreduce.Job: The url to track the job: http://node-master:8088/proxy/application_1608150350331_0002/
2020-12-18 22:34:32,639 INFO mapreduce.Job: Running job: job_1608150350331_0002
2020-12-18 22:34:46,060 INFO mapreduce.Job: Job job_1608150350331_0002 running in uber mode : false
2020-12-18 22:34:46,072 INFO mapreduce.Job: map 0% reduce 0%
2020-12-18 22:34:54,233 INFO mapreduce.Job: map 100% reduce 0%
2020-12-18 22:35:02,372 INFO mapreduce.Job: map 100% reduce 100%
2020-12-18 22:35:03,410 INFO mapreduce.Job: Job job_1608150350331_0002 completed successfully
2020-12-18 22:35:03,588 INFO mapreduce.Job: Counters: 53
  File System Counters
    FILE: Number of bytes read=883
    FILE: Number of bytes written=433275
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=123768
    HDFS: Number of bytes written=661
    HDFS: Number of read operations=8
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=6106
    Total time spent by all reduces in occupied slots (ms)=5774
    Total time spent by all map tasks (ms)=6106
    Total time spent by all reduce tasks (ms)=5774
    Total vcore-milliseconds taken by all map tasks=6106
  
```

Vous pouvez visualiser le résultat en exécutant la commande ci-après :

```

hadoop@node-master: ~/TP_Hadoop/Ventes
hadoop@node-master:~/TP_Hadoop/Ventes$ hdfs dfs -cat storage/Sales/Resultats_Traitement_1/part-r-00000
2020-12-18 22:43:15,341 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Argentina      1
Australia     38
Austria       7
Bahrain       1
Belgium       8
Bermuda       1
Brazil        5
Bulgaria       1
CO            1
Canada       76
Cayman Isls   1
China         1
Costa Rica    1
Country       1
Czech Republic 3
Denmark      15
Dominican Republic 1
  
```

### 3. Travail demandé

Vous allez travailler sur le fichier « [purchases.txt](#) », déjà existant sur votre cluster ([storage/purchases.txt](#)). La structure du fichier est donnée ci-après.

Date	Heure	Ville	Achat	Prix	Type de paiement

#### Objectif 1 :

En utilisant le patron d'architecture MapReduce, affichez le montant total des ventes par ville.

Un compte-rendu de votre travail doit être rendu avant le **13 Novembre 2023**.

Votre compte-rendu doit contenir vos codes avec des commentaires qui permettent de comprendre vos démarches.