
TP1 Le traitement Batch avec Hadoop HDFS

Objectifs du TP

Initiation au framework hadoop et au patron MapReduce, utilisation de docker pour lancer un cluster hadoop de 3 noeuds.

Outils et Versions

- Apache Hadoop Version : 2.7.2.
- Docker Version 17.09.1
- IntelliJ IDEA Version Ultimate 2016.1 (ou tout autre IDE de votre choix)
- Java Version 1.8.
- Unix-like ou Unix-based Systems (Divers Linux et MacOS)

Installation

Nous allons utiliser tout au long de ce TP trois contenaires représentant respectivement un noeud maître (Namenode) et deux noeuds esclaves (Datanodes).

Vous devez pour cela avoir installé docker sur votre machine, et l'avoir correctement configuré. Ouvrir la ligne de commande, et taper les instructions suivantes :

1. Télécharger l'image docker uploadée sur dockerhub:

docker pull hajjitarikensam/hadoop:ensam

2. Créer les trois contenaires à partir de l'image téléchargée. Pour cela : 2.1. Créer un réseau qui permettra de relier les trois contenaires:

docker network create --driver=bridge hadoop

3. Créer et lancer les trois contenaires (les instructions -p permettent de faire un mapping entre les ports de la machine hôte et ceux du contenaire):

docker run -itd --net=hadoop -p 50070:50070 -p 8088:8088 -p 7077:7077 -p 16010:16010 -- name hadoop-master --hostname hadoop-master hajjitarikensam/hadoop:ensam

docker run -itd -p 8040:8042 --net=hadoop --name hadoop-slave1 --hostname hadoop-slave1 hajjitarikensam/hadoop:ensam

docker run -itd -p 8041:8042 --net=hadoop --name hadoop-slave2 --hostname hadoop-slave2 hajjitarikensam/hadoop:ensam

4. Entrer dans le contenaire master pour commencer à l'utiliser :

docker exec -it hadoop-master bash

Le résultat de cette exécution sera le suivant :

```
root@hadoop-master:~#
```

Vous vous retrouverez dans le shell du namenode, et vous pourrez ainsi manipuler le cluster à votre guise. La première chose à faire, une fois dans le contenaire, est de lancer hadoop et yarn. Un script est fourni pour cela, appelé start-hadoop.sh. Lancer ce script.

./start-hadoop.sh

Le résultat devra ressembler à ce qui suit :

```
Toot@hadoop-master:~# ./start-hadoop.sh

Starting namenodes on [hadoop-master]
hadoop-master: Warning: Permanently added 'hadoop-master,172.22.0.2' (ECDSA) to the list of known hosts.
hadoop-master: starting namenode, logging to /usr/local/hadoop/logs/hadoop-root-namenode-hadoop-master.out
hadoop-slave1: Warning: Permanently added 'hadoop-slave1,172.22.0.3' (ECDSA) to the list of known hosts.
hadoop-slave2: Warning: Permanently added 'hadoop-slave2,172.22.0.4' (ECDSA) to the list of known hosts.
hadoop-slave2: starting datanode, logging to /usr/local/hadoop/logs/hadoop-root-datanode-hadoop-slave2.out
hadoop-slave1: starting datanode, logging to /usr/local/hadoop/logs/hadoop-root-datanode-hadoop-slave1.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: Warning: Permanently added '0.0.0.0' (ECDSA) to the list of known hosts.
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop-root-secondarynamenode-hadoop-master.out
starting yarn daemons
starting yarn daemons
starting resourcemanager, logging to /usr/local/hadoop/logs/yarn--resourcemanager-hadoop-master.out
hadoop-slave2: Warning: Permanently added 'hadoop-slave2,172.22.0.4' (ECDSA) to the list of known hosts.
hadoop-slave1: Warning: Permanently added 'hadoop-slave1,172.22.0.3' (ECDSA) to the list of known hosts.
hadoop-slave2: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-root-nodemanager-hadoop-slave2.out
(hadoop-slave1: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-root-nodemanager-hadoop-slave1.out
```

Premiers pas avec Hadoop

Toutes les commandes interagissant avec le système Hadoop commencent par hadoop fs. Ensuite, les options rajoutées sont très largement inspirées des commandes Unix standard.

• Créer un répertoire dans HDFS, appelé *input*. Pour cela, taper :

hadoop fs -mkdir -p input

Si pour une raison ou une autre, vous n'arrivez pas à créer le répertoire input, avec un message ressemblant à ceci: ls: `.': No such file or directory, veiller à construire l'arborescence de l'utilisateur principal (root), comme suit:

hadoop fs -mkdir -p /user/root

Nous allons utiliser le fichier purchases.txt comme entrée pour le traitement MapReduce. Ce fichier se trouve déjà sous le répertoire principal de votre machine master.

• Charger le fichier purchases dans le répertoire input que vous avez créé :

hadoop fs -put purchases.txt input

• Pour afficher le contenu du répertoire input, la commande est:

hadoop fs -ls input

• Pour afficher les dernières lignes du fichier purchases:

```
hadoop fs -tail input/purchases.txt
```

Le résultat suivant va donc s'afficher:

```
[root@hadoop-master:~# hadoop fs -tail input/purchases.txt
        17:59
                 Norfolk Toys
                                 164.34 MasterCard
                 17:59
                         Chula Vista
2012-12-31
                                          Music
                                                  380.67 Visa
2012-12-31
                 17:59
                         Hialeah Toys
                                          115.21
                                                  MasterCard
                 17:59
2012-12-31
                                          Men's Clothing
                         Indianapolis
                                                          158.28
                                                                  MasterCard
                 17:59
                                         414.09
2012-12-31
                         Norfolk Garden
                                                  MasterCard
2012-12-31
                 17:59
                         Baltimore
                                          DVDs
                                                  467.3
                                                          Visa
2012-12-31
                 17:59
                         Santa Ana
                                          Video Games
                                                          144.73
                                                                  Visa
2012-12-31
                 17:59
                         Gilbert Consumer Electronics
                                                          354.66
                                                                  Discover
                                                  124.79
2012-12-31
                 17:59
                         Memphis Sporting Goods
                                                          Amex
                                                          MasterCard
2012-12-31
                 17:59
                         Chicago Men's Clothing
                                                  386.54
2012-12-31
                 17:59
                         Birmingham
                                          CDs
                                                  118.04
                                                          Cash
2012-12-31
                 17:59
                         Las Vegas
                                          Health and Beauty
                                                                   420.46
                                                                           Amex
2012-12-31
                 17:59
                         Wichita Toys
                                          383.9
                                                  Cash
2012-12-31
                 17:59
                         Tucson Pet Supplies
                                                  268.39
                                                          MasterCard
                                         Women's Clothing
2012-12-31
                 17:59
                         Glendale
                                                                   68.05
                                                                           Amex
2012-12-31
                 17:59
                         Albuquerque
                                          Toys
                                                  345.7
                                                          MasterCard
2012-12-31
                 17:59
                         Rochester
                                          DVDs
                                                  399.57
                                                          Amex
2012-12-31
                 17:59
                         Greensboro
                                          Baby
                                                  277.27
                                                          Discover
2012-12-31
                 17:59
                                          Women's Clothing
                                                                  134.95
                                                                           MasterCard
                         Arlington
2012-12-31
                 17:59
                         Corpus Christi
                                         DVDs
                                                  441.61 Discover
root@hadoop-master:~#
```

Nous présentons dans le tableau suivant les commandes les plus utilisées pour manipuler les fichiers dans HDFS:

Instruction	Fonctionnalité
hadoop fs –ls	Afficher le contenu du répertoire racine
hadoop fs –put file.txt	Upload un fichier dans hadoop (à partir du répertoire courant
	linux)
hadoop fs –get file.txt	Download un fichier à partir de hadoop sur votre disque local
hadoop fs -tail file.txt	Lire les dernières lignes du fichier
hadoop fs –cat file.txt	Affiche tout le contenu du fichier
hadoop fs -mv file.txt newfile.txt	Renommer le fichier
hadoop fs –rm newfile.txt	Supprimer le fichier
hadoop fs –mkdir myinput	Créer un répertoire
hadoop fs –cat file.txt \ less	Lire le fichier page par page

Interfaces web pour Hadoop

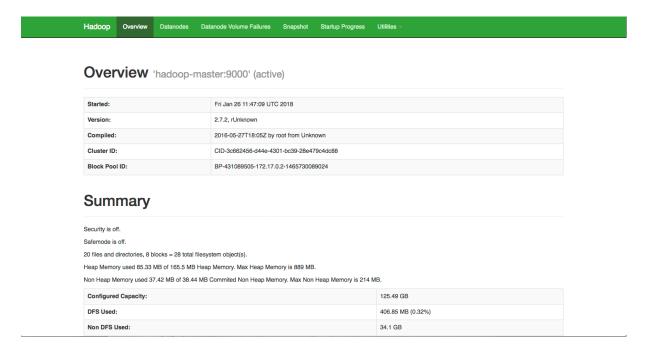
Hadoop offre plusieurs interfaces web pour pouvoir observer le comportement de ses différentes composantes. Vous pouvez afficher ces pages en local sur votre machine grâce à l'option -p de la commande docker run. En effet, cette option permet de publier un port du contenaire sur la machine hôte. Pour pouvoir publier tous les ports exposés, vous pouvez lancer votre contenaire en utilisant l'option -P.

En regardant le contenu du fichier start-container.sh fourni dans le projet, vous verrez que deux ports de la machine maître ont été exposés :

Le port 50070 : qui permet d'afficher les informations de votre namenode.

Le port 8088 : qui permet d'afficher les informations du resource manager de Yarn et visualiser le comportement des différents jobs.

Une fois votre cluster lancé et prêt à l'emploi, vous pouvez, sur votre navigateur préféré de votre machine hôte, aller à : http://localhost:50070 . Vous obtiendrez le résultat suivant :



Vous pouvez également visualiser l'avancement et les résultats de vos Jobs (Map Reduce ou autre) en allant à l'adresse : http://localhost:8088

