

# Utilisez les concepts de base, HDFS et MapReduce pour résoudre un problème de traitement de gros volumes de données :

Il est maintenant temps de se familiariser avec Hadoop Map Reduce.

Généralement, le code Map Reduce est écrit en java. Mais avec une fonctionnalité appelée Hadoop Streaming, vous pouvez écrire vos mappeurs et réducteurs dans n'importe quelle langue. Nous utiliserons python pour écrire le travail de réduction de la carte.

**Application :** Analyse des données de vente pour la chaîne de magasins à l'aide de Cloudera Platform.

- Nous allons utiliser un programme mapper et le reducer pour calculer les ventes totales de chaque magasin aux États-Unis.
- Analyser et déterminer les revenus générés dans différents magasins d'un pays / d'une région.

## Outils :

- Utilisez les concepts de base, HDFS et MapReduce de Big Data Hadoop Platform pour résoudre le problème.
- Traitement de gros volumes de données en parallèle en divisant le travail en un ensemble de tâches indépendantes.
- Utilisez le fichier de processus MapReduce pour le diviser en morceaux et le traiter en parallèle.

## Travail demandé :

1. Créer un répertoire appelé myinput
2. Copier le fichier purchases.txt dans le répertoire myinput
3. Afficher les dernières lignes du fichier

### ***Solution :***

- Créer un répertoire dans HDFS, appelé myinput. Pour cela, taper:

*hadoop fs -mkdir myinput*

- Pour copier le fichier purchases.txt dans HDFS sous le répertoire myinput, taper la commande:

*hadoop fs -put purchases.txt myinput/*

- Pour afficher le contenu du répertoire myinput, la commande est:

*hadoop fs -ls myinput*

4. Le dossier du projet contient les fichiers **mapper.py** et **reducer.py**.

- **Afficher le fichier mapper.py**

**Description :** le mapper lit le morceau du fichier et divise le fichier à l'aide du délimiteur de tabulation. Ensuite, il imprime uniquement l'emplacement du magasin et le coût de l'article comme sortie.

Tester ce mapper en local sur les 50 premières lignes du fichier purchases.txt en tapant l'instruction suivante, directement à partir de votre répertoire code:

*head -50 ../data/purchases.txt | ./mapper.py*

- **Afficher le code réduire**

**Description :** le reducer.py prend l'entrée triée et continue de vérifier si la nouvelle clé est égale à la clé précédente. Lorsque le changement se produit, il imprime le nom du magasin et le total des ventes du magasin.

Tester ce Reducer sur le disque local, en utilisant cette instruction.

```
head -50 ../data/purchases.txt |./mapper.py |sort |./reducer.py
```

**Lancer un Job sur Hadoop** Lancer un job sur Hadoop implique qu'on fera appel au mapper puis au reducer sur une entrée volumineuse, et qu'on obtiendra à la fin un résultat, directement sur HDFS. Pour faire cela, l'instruction à exécuter sur notre machine virtuelle est du type :

Une fois que vous avez écrit le mapper et le reducer, vous pouvez démarrer le travail de réduction de carte à l'aide de la commande:

```
hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-cdh5.13.0.jar -mapper code/mapper2.py -reducer code/reducer.py -file code/mapper2.py -file code/reducer.py -input myinput/purchases.txt -output joboutput
```

**Remarque 1 :** Nous utilisons Hadoop Streaming qui permet de créer et lancer des jobs MapReduce avec tout type d'exécutable ou script en tant que mapper et reducer. Ici nos scripts sont en Python, mais les mappers et reducers pourraient être des classes Java, des utilitaires unix, des scripts R, etc. La manière standard est d'écrire des programmes MapReduce en Java via l'API Java MapReduce.

**Remarque 2 :** Le répertoire d'entrée doit contenir un seul fichier. Le répertoire de sortie ne doit pas exister avant l'exécution de l'instruction.

**Remarque 3 :** Après avoir exécuté cela, le calcul prend un certain temps. Vous pouvez observer le fonctionnement des mappeurs et des réducteurs dans le terminal de sortie. Le terminal affichera une variété d'informations comme le pourcentage d'achèvement du travail.

- **Sortie :** La sortie de notre programme est enregistrée dans le fichier joboutput / part-00000

Vous pouvez afficher la sortie dans le même terminal en exécutant la commande :

```
hadoop fs -ls joboutput
```

Au lieu de visualiser l'ensemble de la sortie, vous pouvez afficher une plus petite partie de notre sortie en utilisant la commande :

```
hadoop fs -cat joboutput / part-00000 | Moins
```

Une fois votre tâche terminée, vous pouvez maintenant supprimer ce fichier de sortie de travail de notre système HDFS. La raison en est que nous ne pouvons pas générer différentes charges de travail de sortie avec le même nom de fichier de sortie. Ceci peut être réalisé en exécutant la commande suivante

```
hadoop fs -rm -r -f joboutput
```

Pendant ce temps, vous pouvez exporter les données dans un fichier texte et l'enregistrer dans votre répertoire si vous le souhaitez pour une utilisation future.

## Application sur des données réelles

Ecrire vos propres Mappers et Reducers pour :

- 1) Donner le chiffre d'affaire par magasin
- 2) Donner la liste des ventes par catégorie de produits.
  - Quelle est la valeur des ventes pour la catégorie Toys?
  - Et pour la catégorie Consumer Electronics?
- 3) Quel est le nombre total des ventes et la valeur totale des ventes de tous magasins confondus?

Tester les solutions :

```
head -50 Downloads/purchases.txt | code/mapper3.py
```

```
head -50 Downloads/purchases.txt | code/mapper3.py | sort |code/reducer3.py
```

Pour aller loin :

Utiliser un cluster AWS-EMR (Elastic Map Reduce) d'Amazon pour exécuter un Job Map Reduce de votre choix sur un vrai cluster distribué. [<https://aws.amazon.com/fr/emr/>]

**Afficher le dernier kilo-octet du fichier sur stdout.**

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

**Pour afficher le contenu** de votre fichier texte purchase.txt  
qui est présent dans votre répertoire hadoop.

```
hadoop fs -cat hadoop/purchases.txt
```

**# Tester vos codes en local :**

```
head -50 Downloads/purchases.txt |code/mapper2.py |sort |code/reducer.py
```

**# Exécuter map-reduce sous hadoop :**

```
hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-  
streaming-2.6.0-mr1-cdh5.13.0.jar -mapper code/mapper2.py -reducer  
code/reducer.py -file code/mapper2.py -file code/reducer.py -input  
myinput/purchases.txt -output joboutput
```

**Remarque :** nombre de map = nombre de block (en fonction du volume du fichier)  
Généralement le nombre de reducer est au voisinage de 10% du nombre de mapper

**# Afficher le fichier de sortie :**

```
hadoop fs -tail output/part-r-00000
```

**# consulter la liste des tâches map-reduce**

```
mapred job -list all
```

**# consulter le statut du job**

```
mapred job -status job_1670364086924_0001
```

**#☹ Forcer l'arrêt une tâche**

```
mapred job -status job_1670364086924_0002
```

**Remarque :** vous pouvez consulter ces informations via la page web.

--- all applications