

Introduction à Map-Reduce - TP

Martin Kirchgessner, Vincent Leroy

2015

Votre compte-rendu sera une archive `Votre_Nom.zip` (ou `tar.gz`) contenant :

- Un fichier *texte* (`reponses.txt`, pas `votrenom.doc` et surtout pas `.pdf`) avec les réponses aux questions qui suivent, et
- Vos sources. Ecrivez un fichier par problème qui vous est soumis. Exemple : `Question3_1.java` pour la question 3.1.

Envoyez l'archive à `martin.kirchgessner@imag.fr`, avec [TPHadoop] au début du sujet. Mettez votre binôme en copie de cet e-mail.

Envoyez l'archive le jour de la dernière séance de TP. -2 points par jour de retard.

0 Kit de survie sous Eclipse

Usez et abusez des raccourcis clavier suivants :

- `Ctrl + Espace` pour l'autocomplétion, qui permet aussi de générer des méthodes (essayez dans le corps d'une classe)
- `Ctrl + Shift + I` "Quick fix", propose des actions pertinentes selon ce qu'il y a sous le curseur.
- `Ctrl + Shift + O` pour ajouter les imports manquants et retirer les inutiles.
- `Ctrl + clic` pour atteindre la définition du nom sous le curseur de la souris (utilisez `Alt + ←` pour revenir là où vous étiez).
- `Alt + Shift + R` pour renommer ce qui est sous le curseur (partout où c'est pertinent).

Voir aussi [http://eclipse-tools.sourceforge.net/Keyboard_shortcuts_\(3.0\).pdf](http://eclipse-tools.sourceforge.net/Keyboard_shortcuts_(3.0).pdf)

1 Prise en main

Nous avons préparé un projet Eclipse qui contient les bibliothèques nécessaires. Vous pourrez ainsi tester vos programmes localement, sur de petits fichiers, avant de les lancer sur le cluster. Ce qui est très pratique si vous avez besoin de debugger avec des points d'arrêt.

1. Téléchargez sur votre poste de travail `TPIntroHadoop.zip`
2. Dans les menus d'Eclipse, cliquez sur `File > Import ...`
3. Dans la catégorie **General**, sélectionnez **Existing Projects into Workspace** (et non **Archive**)
4. Choisissez **select archive file** et retrouvez l'archive téléchargée
5. Cliquez sur **Finish**

Le répertoire des sources contient `Question0_0.java`, qui peut vous servir d'exemple pour écrire chaque programme (mais pour certaines questions, attention à modifier les types de clés/valeurs qu'il déclare). Vous remarquerez que le **Mapper** et **Reducer** sont déclarés comme des classes internes : ce

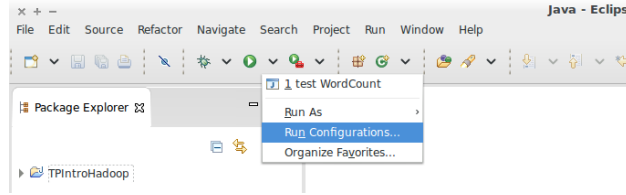


FIGURE 1 – Pour indiquer quels paramètres utiliser en test local, il faut créer une configuration d’exécution via le menu du bouton “Run”

n’est pas imposé par le framework, c’est même déconseillé. On se le permet dans le cadre du TP pour faciliter les corrections.

1.1 Exécution locale

Avant de lancer un programme MapReduce sur un cluster, on le teste *toujours* localement, sur un échantillon des données à traiter. Dans ce cas, le framework ne lit/écrit pas sur HDFS, mais sur votre système de fichiers local.

1. Créez un programme qui compte le nombre d’occurrences de chaque mot dans un fichier texte (cf. “WordCount” vu en cours).
2. Créez à la racine du projet un fichier texte, avec le (petit) contenu de votre choix.
3. La fonction `main` prend en argument le chemin du fichier d’entrée et de sortie. Avant de lancer le programme avec Eclipse, il faut donc lui indiquer ces chemins (vous pouvez utiliser des chemins relatifs à la racine du projet) - cf. Figure 1.

Lancez le programme. Les traces d’exécution du framework apparaissent dans l’onglet “Console”. En cas de bug, vous y verrez l’exception à l’origine du plantage. Ouvrez le fichier généré et vérifiez vos résultats.

Quand le programme marche, à la fin des traces d’exécution se trouvent les compteurs tenus à jour par Hadoop. Ajoutez/retirez du contenu à votre fichier de texte, relancez votre programme et observez l’évolution des compteurs pour répondre aux questions suivantes.

1. A quoi correspondent les compteurs `Map input records`? Et `Map output records`?
2. Quel est le lien entre `Map output records` et `Reduce input records`?
3. A quoi correspond `Reduce input groups`?

1.2 Premier contact avec HDFS

Un cluster a été déployé pour les besoins du TP. Un nom d’utilisateur et un mot de passe vous ont été remis, ils vous permettent de vous connecter via `ssh` au noeud maître “NameNode” (qui est aussi un système Linux classique).

1. Ouvrez un shell sur le noeud maître avec `ssh votrelogin@152.77.78.100`¹
2. Exécutez `hdfs dfs` pour lister les commandes disponibles.
3. Exécutez `hdfs dfs -ls /` pour afficher le contenu de la racine d’HDFS.

¹Vous pouvez changer votre mot de passe, mais nous vous recommandons plutôt de copier votre clé publique dans `~/.ssh/authorized_keys` avec `ssh-copy-id`

4. Exécutez `hdfs dfs -ls` pour afficher le contenu de votre répertoire HDFS personnel.

Prenez le temps d'explorer ce que contient notre cluster.

Question : quel est le chemin, dans HDFS, vers votre répertoire personnel ?

1.3 Exécution sur le cluster

Le programme écrit en partie 1.1, s'il fonctionnait localement, devrait fonctionner sur le cluster sans modifications. Mais il faut d'abord en faire un paquet compilé.

1. Dans le navigateur d'Eclipse, faites un clic droit sur le répertoire `src`. Cliquez sur "Export..."
2. Ouvrez la catégorie "Java", puis double-cliquez sur "JAR file".
3. Parmi les ressources à exporter, **vérifiez que le répertoire `lib` n'est pas sélectionné.**
4. Entrez le nom du paquet, par exemple `tp.jar`.

Le paquet généré sera à la racine de votre espace de travail Eclipse : transférez-le avec `scp` vers votre répertoire personnel sur le noeud maître.

Nous allons maintenant compter le nombre d'occurrences de chaque mot dans les 5 tomes des *Misérables* de Victor Hugo. Si votre classe principale s'appelle `Question1_1`, lancez votre programme sur le noeud maître en tapant :

```
hadoop jar tp.jar Question1_1 /data/miserables wordcount
```

Question : dans les traces d'exécution (au début), retrouvez le nombre de "splits" lus sur HDFS. A quel compteur ce nombre correspond-il ?

En préparation de la question suivante, conservez les traces d'exécution de ce programme dans un fichier temporaire, et veillez à conserver le répertoire qui contient les résultats.

1.4 Combiner et nombre de reducers

Nous allons modifier légèrement le programme, et comparer les compteurs et résultats.

- Ajoutez `job.setCombinerClass(???)` dans le programme principal (à vous de compléter avec la bonne valeur en argument).
- Ajoutez aussi `job.setNumReduceTasks(3)`.
- Vérifiez que le programme fonctionne localement, puis exécutez-le sur les 5 tomes des *Misérables* (sans écraser le résultat de la partie 1.3).

Questions :

1. Quelle est la différence entre le répertoire de résultats obtenu ici, et celui de la partie 1.3 ? Pourquoi ?
2. Quels compteurs permettent de vérifier que le combiner a fonctionné ?
3. Quels compteurs permettent d'estimer le gain effectivement apporté par le combiner ? Comparez aux valeurs obtenues sans combiner pour justifier votre réponse.

Si vous avez le temps : Quel est le mot le plus utilisé dans *Les Misérables* ?

2 Top-tags Flickr par pays, avec tri en mémoire

Nous allons maintenant analyser des meta-données de photos, en utilisant un extrait du jeu de données “Yahoo! Flickr Creative Commons 100M”. Le but est de trouver les tags les plus utilisés par pays. Pour identifier un pays à partir des coordonnées d’une photo, vous utiliserez la classe fournie `Country`.

L’archive que vous avez téléchargée au début du TP contient deux fichiers qui vont vous aider à vous familiariser avec ces données :

- `flickrSpecs.txt` décrit le format du fichier. Pour décoder les textes, par exemple les tags, utilisez `java.net.URLDecoder.decode(String s)`,
- `flickrSample.txt` est un extrait à utiliser pour les tests locaux.

A partir d’ici, la référence de l’API Hadoop pourra vous être utile :

<https://hadoop.apache.org/docs/r2.2.0/api/index.html>

Attention : dans cette référence, de nombreuses classes existent en double. Quand c’est le cas, choisissez toujours celle qui appartient au package `org.apache.hadoop.mapreduce`, qui correspond à l’API moderne, dite “YARN”. Les doublons appartiennent au package `org.apache.hadoop.mapred`, qui est celui de l’ancienne API, dite “MR1”.

2.1 Map et Reduce

Puisque l’on veut trouver les K tags les plus utilisés *par pays*, a priori l’identifiant d’un pays (les deux lettres retournées par `country.toString()`) est une bonne clé intermédiaire. Les valeurs associées seront les tags (les chaînes de caractères) qui ont été associés à une photo prise dans ce pays. En sortie, il ne restera qu’au plus K couples (*pays, tag*) par pays.

A vous d’implémenter le programme complet. Pour la fonction `reduce` :

1. Comptez le nombre d’occurrences de chaque tag avec `java.util.HashMap<String,Integer>`.
2. Puis trie les tags par nombre d’occurrences décroissants :
 - Créez (dans un fichier séparé) une classe `StringAndInt` qui implémente l’interface `Comparable<StringAndInt>` et contient deux champs, qui dans ce cas représenteront le tag et son nombre d’occurrences. La méthode `compareTo` ne prendra en compte que le nombre d’occurrences.
 - Utilisez la classe `java.util.PriorityQueue<StringAndInt>` afin de ne conserver que les K tags les plus populaires.

K sera un nouveau paramètre à passer au programme principal via la ligne de commande. Utilisez l’objet `Configuration` pour transmettre sa valeur du programme principal aux *reducers* (via l’objet `Context`).

Implémentez puis testez *localement* ce programme.

2.2 Combiner

Question : pour pouvoir utiliser un *combiner*, quel devrait être le type des données intermédiaires ? Donnez le type sémantique (que représentent ces clés-valeurs ?) et le type Java.

Avant de tester ce programme sur le cluster :

1. Modifiez la classe que l'on a ajouté en 2.1, de sorte qu'elle implémente aussi l'interface **Writable**. La référence de cette interface vous fournira quelques indices, toutefois :
 - encapsulez un **Text** pour la chaîne de caractères
 - dans leur exemple, la méthode statique **read** est inutile
 - dans leur exemple, il manque un constructeur sans arguments
2. Modifiez une copie du programme écrit dans la partie 2.1, en y ajoutant un *combiner*. Cette fois, on ne peut pas réutiliser le *reducer*, il faut implémenter une troisième classe (qui doit étendre **Reducer<K,V,K,V>**).

Exécutez cette variante (avec $K = 5$) sur `/data/flickr.txt`. Quels sont les tags les plus utilisés en France ?

Question : Dans le *reducer*, nous avons une structure en mémoire dont la taille dépend du nombre de tags distincts : on ne le connaît pas a priori, et il y en a potentiellement beaucoup. Est-ce un problème ?

3 Top-tags Flickr par pays, avec tri par Hadoop

Une fonction (parfois) importante du framework MapReduce est le tri des clés intermédiaires. Le framework permet au développeur, en plus de spécifier des types ad-hoc, de fournir ses propres fonctions de tri. On peut ainsi contrôler :

- la formation des groupes avant la phase *reduce*, c'est à dire quelles valeurs seront associées à quelle clé dans le *reducer*, et
- l'ordre des couples (*cle, valeur*) qui seront itérés par le *reducer*.

Vous trouverez les explications détaillées dans la référence de la classe **Reducer**.

L'objectif de cette partie est de faire faire à Hadoop le tri des tags par nombre d'occurrences décroissant.

Question préliminaire : spécifiez les 2 jobs nécessaires (le second va utiliser le résultat du premier), en précisant le type des clés/valeurs en entrée, sortie et intermédiaires. Pour le second, spécifiez aussi les comparateurs à utiliser pour le tri et pour le découpage en groupes.

Avant de vous lancer dans l'implémentation, validez votre réponse avec les encadrants.

3.1 Passes MapReduce en chaîne

Commencez par implémenter et tester localement le premier *job*.

Cette première passe matérialise un fichier intermédiaire, qui ne sera utilisé que par une autre passe MapReduce. Dans ces cas-là, plutôt que d'écrire le fichier au format texte, on utilise le format binaire d'Hadoop. Cela permet d'utiliser des types sémantiques en entrée du **Mapper** qui va suivre, au lieu du couple **LongWritable, Text** imposé par le format texte.

Pour cela, dans la fonction qui prépare le `job` que vous venez d'écrire, changez l'*input format* afin d'utiliser `SequenceFileOutputFormat`.

Exécutez cette version-là sur le cluster (n'oubliez pas le *combiner*).

Question : un avantage de cette méthode est que le tri est réalisé sur le disque dur des *data nodes*, plutôt qu'en mémoire : on peut donc trier des quantités de données plus importantes². Dans notre application, cette méthode présente un autre avantage pour la fonction `reduce` finale ; lequel ?

Implémentation du deuxième job :

- notre classe `StringAndInt` doit maintenant implémenter une seule interface : `WritableComparable<StringAndInt>`
- afin de lire le fichier binaire généré par la première passe, utilisez `SequenceFileInputFormat`.
- pour faciliter la lecture du résultat, utilisez `TextOutputFormat` en sortie.
- les fonctions de tri sont à implémenter dans une classe dédiée, qui
 - doit étendre `WritableComparator`, et
 - doit avoir un constructeur sans arguments, qui appelle `super(StringAndInt.class, true)`.
- **attention** le tri des données intermédiaires est fait *avant* le découpage en groupes. La fonction de tri principale doit donc pre-placer les groupes.

Question : s'il existe des tags classés *ex aequo* dans le top-*K* d'un pays, a-t-on la garantie d'obtenir toujours le même résultat d'une exécution à l'autre ? Pourquoi ?

3.2 Placement des données

Les tâches de *map* comme de *reduce* sont assignées en priorité aux machines qui ont déjà, localement, les données. Le placement des blocs du fichier d'entrée est contrôlé par HDFS, mais pour les données intermédiaires il est possible d'indiquer à quelle machine assigner chaque couple (*cle*, *valeur*). On peut ainsi contrôler comment les couples intermédiaires seront répartis entre les *reducers*, ce qui est un moyen indirect de contrôler le placement des données en sortie du *reduce*.

Dans l'application qu'on vient d'écrire, il est utile d'écrire un *partitioner* pour le premier job, afin d'assurer la localité des données pour le second.

A vous de trouver lequel et de l'implémenter.

Il faut implémenter une classe qui étend `Partitioner`, et la fournir à `job.setPartitioner`.

Questions :

1. Le fichier complet pèse 44.4 Go, et nous utilisons des blocs de 64Mo. Avec combien de machines le premier *map* sera le plus rapide possible ?
2. Une fois que nous utilisons ce *partitioner*, le deuxième *job* génère-t-il des transferts de données sur le réseau ? Pourquoi ?
3. En pratique : si les deux *jobs* utilisent `job.setNumReduceTasks(3)`, d'après les compteurs que se passe-t-il ?

²C'est d'ailleurs le test classique des performances d'un cluster, cf. <http://sortbenchmark.org/>.