

UNIVERSITÉ DE BORDEAUX I

PROGRAMMATION LARGE
ECHELLE

Compte rendu du TP 6
Map Reduce Patterns TopK

Réalisé par :

GAMELIN Antoine
LASSOUANI Sofiane

TABLE DES MATIÈRES

1	Introduction	1
1.1	Introduction	1
2	Exercice sur le Pattern TopK	2
2.1	Introduction :	2
2.2	Exercice 1 :	2
2.2.1	CityWritable	2
2.2.2	Mapper	2
2.2.3	Reducer	2
2.3	Exercice 2 :	3
2.3.1	Combiner	3
2.4	Exercice 3 :	3
3	Conclusion	4

INTRODUCTION

1.1 Introduction

Le but de ce TP est de mettre en oeuvre le pattern Top K vue en cours. Ce pattern consiste à sortir les K plus (grande/petite) valeur d'un ensemble d'entrée.

EXERCICE SUR LE PATTERN TOPK

2.1 Introduction :

Cette première partie du TP est consacré au Top K, qui sert à retrouver les k villes les plus peuplées. Nous utiliserons le fichier **worldcitiespop.txt**, de type CSV, qui contient les noms de ville avec leurs populations et des informations annexes.

2.2 Exercice 1 :

2.2.1 CityWritable

Pour bien mener ce TP nous avons décidé de créer un nouveau type de donnée qui implémente writable. Cette structure permet d'avoir un accès plus rapide aux différents champs du fichier csv et qu'on ne soit pas obligé de split à chaque fois la chaîne de caractère.

2.2.2 Mapper

Concernant le Mapper, nous avons utilisé la fonction setup afin d'initialiser un arbre binaire TreeMap en java. L'avantage du TreeMap par rapport au SortedArray, c'est que lorsqu'on travaille avec de nombreuses données il est plus rapide à l'exécution.

Le TreeMap reçoit comme clef le nombre de population et en valeur, un CityWritable avec toutes les informations.

Le mapper, ajoute chaque ville au TreeMap, et quand le nombre d'élément est supérieur à k, nous supprimons à chaque fois l'élément avec la plus petite population.

Dans ce TP le mapper n'envoie aucun message, c'est le clean up qui lors de son appel retourne les 10 plus grandes villes.

Afin de tester notre job nous avons initialisé la variable k à 10 sans passer par la configuration.

2.2.3 Reducer

Le reducer récupère le résultat des différents mappers et fonctionne comme le mapper il possède un TreeMap et supprime à chaque fois la ville avec la plus faible population quand la taille du TreeMap dépasse k.

2.3 Exercice 2 :

Dans cette partie il nous est demandé d'écrire un combiner, qui va alléger la charge au reducer. En effet sur un datanode, plusieurs mapper peuvent être exécuté, est au lieu d'envoyer les directement toutes les données au reducer, nous utilisons un combiner pour envoyer moins de données sur le reseaux.

2.3.1 Combiner

Le but du combiner est presque égale au reducer, la seule différence, c'est qu'en sortie au lieu de renvoyer un `<NullWritable,Text>` nous renvoyons un `<NullWritable, CityWritable>` pour que le réducteur puisse interpréter le résultat.

2.4 Exercice 3 :

Dans cette partie nous avons ajouté un paramètre au lancement du programme permettant de définir le nombre k de ville plus grande.

CONCLUSION

Ce tp nous a permis de mettre en pratique le pattern Top K étudié en cours .

Ce pattern est très simple d'implémentation et permet de traiter les fichiers très rapidement grâce à Hadoop.

Nous avons cependant rencontrés un petit problème, le reducer ne s'exécutait pas car l'entête de la classe réduire n'était pas identique que la méthode réduit.