



ÉCOLE CENTRALE DE LYON

MSO 3.1 - TECHNOLOGIES INFORMATIQUES DU BIG DATA
OPTION INFORMATIQUE
RAPPORT

TP1 – Hadoop map-reduce

Élèves :
Thibaud MIQUEL

Enseignant :
Stéphane DERRODE

Table des matières

1	Présentation du TP	3
1.1	Objectif du TP	3
1.2	Présentation d'Hadoop	3
1.3	Présentation de Map Reduce	3
2	Réponse aux questions	4
2.1	Exercice 1 - Questionner un fichier de ventes	4
2.1.1	1 - Quel est le nombre d'achats effectués pour chaque catégorie d'achat ? .	4
2.1.2	2 - Quelle est la somme totale dépensée pour chaque catégorie d'achat ? .	4
2.1.3	3 - Quelle somme est dépensée dans la ville de San Francisco dans chaque moyen de paiement ?	5
2.1.4	4 - Dans quelle ville la catégorie Women's Clothing a permis de générer le plus d'argent Cash ?	6
2.1.5	5 - Ajoutez une requête originale et complexe (i.e. nombre de STEPS > 1) de votre choix sur ce fichier. Dans le CR, prenez le temps de m'expliquez l'objectif de la requête, de présenter le code et le résultat obtenu : tableau, graphique.	7
2.2	Exercice 2 - Anagramme	8

Table des figures

1	Code pour la question 1	4
2	Résultat pour la question 1	4
3	Code pour la question 2	5
4	Résultat pour la question 2	5
5	Code pour la question 3	6
6	Résultat pour la question 3	6
7	Code pour la question 4	6
8	Résultat pour la question 4	7
9	Code pour la question 5	7
10	Résultat pour la question 5	8
11	Code pour la partie 2	9
12	Résultat pour la partie 2	9

1 Présentation du TP

1.1 Objectif du TP

L'objectif de ce TP est d'utiliser la bibliothèque MRJob pour coder des requêtes et les exécuter directement sur Hadoop en passant par Docker.

1.2 Présentation d'Hadoop

Hadoop est un framework Java open source utilisé pour le stockage et traitement des big data. Les données sont stockées sur des serveurs standard peu coûteux configurés en clusters.

Dans son système de fichiers distribué HDFS (Hadoop Distributed File System), les données sont divisées en datasets de taille réduite et enregistrées sur des clusters de serveurs standard. Ces serveurs de base étant construits avec des configurations matérielles très simples, ils sont peu coûteux et peuvent évoluer très facilement pour suivre l'augmentation des volumes de données.

Hadoop utilise le modèle de programmation fonctionnelle MapReduce pour exécuter un traitement parallèle entre les datasets. Lorsqu'une requête est transmise à la base de données, les jobs (et les données) spécifiés dans cette requête ne sont pas traités séquentiellement, mais répartis et exécutés simultanément sur différents serveurs. À l'issue du traitement, les résultats des différents jobs sont regroupés et renvoyés à l'application, ce qui améliore considérablement la vitesse de traitement.

YARN est l'abréviation de « Yet Another Resource Negotiator » (plus simplement, un négociateur de ressources). Cet élément assure la gestion et planification des ressources (clusters) Hadoop et décide de ce qui doit se passer dans chaque nœud de données. Le nœud maître central qui gère toutes les demandes de traitement est le « Resource Manager ». Le Resource Manager interagit avec les différents Node Managers : chaque DataNode esclave possède son propre Node Manager pour l'exécution des tâches.

1.3 Présentation de Map Reduce

MapReduce est un Framework de traitement de données en clusters, inventé par Google, dans lequel sont effectués des calculs parallèles, et souvent distribués, de données potentiellement très volumineuses.

Il s'agit d'un composant central du Framework logiciel Apache Hadoop, qui permet le traitement robuste et distribué d'ensembles de données non structurées massifs sur des clusters d'ordinateurs, au sein desquels chaque nœud possède son propre espace de stockage. Concrètement, le framework propose deux fonctionnalités principales. Il répartit le travail de calcul sur les différents nœuds du cluster (map), puis les organise et réduit les résultats fournis par chaque nœud en une seule réponse cohérente à une requête (reduce).

<https://github.com/mattwg/mrjob-examples>

2 Réponse aux questions

2.1 Exercice 1 - Questionner un fichier de ventes

2.1.1 1 - Quel est le nombre d'achats effectués pour chaque catégorie d'achat ?

Pour la question 1 :

```
Q1.py > MRachatcateg > reducer
1  #!/usr/bin/env python3
2
3  #-*- coding: utf-8 -*-
4  from mrjob.job import MRJob
5  import sys
6
7  class MRachatcateg(MRJob):
8      def mapper(self, _, line):
9          for word in [(line.split("\t"))[3]]:
10             yield(word, 1)
11      def reducer(self, word, counts):
12          yield(word, sum(counts))
13
14  if __name__ == '__main__':
15      MRachatcateg.run()
```

FIGURE 1 – Code pour la question 1

```
Streaming final output from hdfs:///user/root/tmp/mrjob/Q1.root.20220305.161437.723220/out
"Baby" 230293
"Books" 229787
"CDs" 230039
"Cameras" 229320
"Children's Clothing" 230469
"Computers" 229059
"Consumer Electronics" 229761
"Crafts" 229749
"DVDs" 230274
"Garden" 230073
"Health and Beauty" 229667
"Men's Clothing" 230430
"Music" 230150
"Pet Supplies" 229222
"Sporting Goods" 229932
"Toys" 229964
"Video Games" 230237
"Women's Clothing" 230050
Removing HDFS temp directory hdfs:///user/root/tmp/mrjob/Q1.root.20220305.161437.723220...
```

FIGURE 2 – Résultat pour la question 1

2.1.2 2 - Quelle est la somme totale dépensée pour chaque catégorie d'achat ?

Pour la question 2 :

```
Q2.py > ...
3  -*- coding: utf-8 -*-
4  from mrjob.job import MRJob
5  import sys
6
7  class MRsommetotcateg(MRJob):
8      def mapper(self, _, line):
9          for word in [(line.split("\t"))[3]]:
10             yield(word, float((line.split("\t"))[4]))
11      def reducer(self, word, counts):
12          yield(word, sum(counts))
13
14  if __name__ == '__main__':
15      MRsommetotcateg.run()
```

FIGURE 3 – Code pour la question 2

```
Streaming final output from hdfs:///user/root/tmp/mrjob/Q2.root.20220305.162328.951440/ou
"Baby" 57491808.439999565
"Books" 57450757.910000086
"CDs" 57410753.04000101
"Cameras" 57299046.64000095
"Children's Clothing" 57624820.93999975
"Computers" 57315406.31999985
"Consumer Electronics" 57452374.130001
"Crafts" 57418154.50000002
"DVDs" 57649212.13999978
"Garden" 57539833.109999985
"Health and Beauty" 57481589.56000032
"Men's Clothing" 57621279.04
"Music" 57495489.70000011
"Pet Supplies" 57197250.240000114
"Sporting Goods" 57599085.88999987
"Toys" 57463477.10999886
"Video Games" 57513165.579998754
"Women's Clothing" 57434448.969999254
Removing HDFS temp directory hdfs:///user/root/tmp/mrjob/Q2.root.20220305.162328.951440...
```

FIGURE 4 – Résultat pour la question 2

2.1.3 3 - Quelle somme est dépensée dans la ville de San Francisco dans chaque moyen de paiement ?

Pour la question 3 :

```
Q3.py > ...
1  #!/usr/bin/env python3
2
3  -*- coding: utf-8 -*-
4  from mrjob.job import MRJob
5  import sys
6
7  class MRsommeSFmoyen(MRJob):
8      def mapper(self, _, line):
9          splitted_line=line.split("\t")
10         if splitted_line[2]=="San Francisco":
11             yield(splitted_line[-1], splitted_line[4])
12         def reducer(self, word, counts):
13             yield(word,sum(counts))
14
15     if __name__ == '__main__':
16         MRsommeSFmoyen.run()
```

FIGURE 5 – Code pour la question 3

```
Streaming final output from hdfs:///user/root/tmp/mrjob/Q3.root.20220305.163608.016365/outp
"Amex" 1984573.8900000043
"Cash" 1982297.8699999987
"Discover" 2009176.530000004
"MasterCard" 2007667.0500000059
"Visa" 2011855.1999999892
Removing HDFS temp directory hdfs:///user/root/tmp/mrjob/Q3.root.20220305.163608.016365...
```

FIGURE 6 – Résultat pour la question 3

2.1.4 4 - Dans quelle ville la catégorie Women's Clothing a permis de générer le plus d'argent Cash ?

Pour la question 4 :

```
Q4.py > MRmaxWomClothinCash > steps
2
3  -*- coding: utf-8 -*-
4  from mrjob.job import MRJob
5  from mrjob.step import MRStep
6  import sys
7
8  class MRmaxWomClothinCash(MRJob):
9      def mapper(self, _, line):
10         splitted_line = line.split("\t")
11         if splitted_line[3]=="Women's Clothing" and splitted_line[5]=="Cash":
12             yield(splitted_line[2], float(splitted_line[4]))
13
14         def reducer_int(self, word, counts):
15             yield None, (word,sum(counts))
16
17         def reducer_find_max(self, _, pairs):
18             # each item of word_count_pairs is (count, word),
19             # so yielding one results in key=counts, value=word
20             yield max(pairs)
21
22     def steps(self):
23         return [
24             MRStep(mapper=self.mapper,
25                   reducer=self.reducer_int),
26             MRStep(reducer=self.reducer_find_max)
27         ]
28
29     if __name__ == '__main__':
30         MRmaxWomClothinCash.run()
```

FIGURE 7 – Code pour la question 4

```
job output is in hdfs:///user/root/tmp/mrjob/Q4.root.20220305.171305.735244...
Streaming final output from hdfs:///user/root/tmp/mrjob/Q4.root.20220305.171305.735244...
"Winston\u2013Salem"      107966.709999999992
Removing HDFS temp directory hdfs:///user/root/tmp/mrjob/Q4.root.20220305.171305.735244...
Removing temp directory /tmp/Q4.root.20220305.171305.735244...
```

FIGURE 8 – Résultat pour la question 4

2.1.5 5 - Ajoutez une requête originale et complexe (i.e. nombre de STEPS > 1) de votre choix sur ce fichier. Dans le CR, prenez le temps de m'expliquez l'objectif de la requête, de présenter le code et le résultat obtenu : tableau, graphique...

Ma requête consiste à renvoyer pour chaque ville et pour chaque catégorie présente dans ces villes le moyen de paiement ayant la somme de dépense la plus faible selon cette granularité d'information.

Pour cela on effectuera donc 2 étapes de map-reduce.

La première consiste à mapper tous les triplets de : city, catégorie d'achat et moyen de paiement avec leur montant puis la phase de reduce consiste à sommer ces montants.

La 2eme etapes consistes à mapper tous les duets de city et catégorie d'achat avec leur duet de moyen de paiement et de montant d'achat puis la phase de reduce consiste à prendre pour le premier duet le 2eme duet avec le montant minimal.

```
Q5.py > MRminCategCityPayment > steps
3  # coding: utf-8
4  from mrjob.job import MRJob
5  from mrjob.step import MRStep
6  import sys
7
8  class MRminCategCityPayment(MRJob):
9      def mapper_1(self, _, line):
10         splitted_line = line.split("\t")
11         yield((splitted_line[2], splitted_line[3], splitted_line[5]), float(splitted_line[4]))
12
13     def reducer_1(self, word, counts):
14         yield(word, sum(counts))
15
16     def mapper_2(self, word, counts):
17         yield((word[0], word[1]), (word[2], counts))
18
19     def reducer_2(self, word, counts):
20         yield(word, min(counts))
21
22     def steps(self):
23         return [
24             MRStep(mapper=self.mapper_1,
25                   reducer=self.reducer_1),
26             MRStep(mapper=self.mapper_2,
27                   reducer=self.reducer_2)
28         ]
29
30 if __name__ == '__main__':
31     MRminCategCityPayment.run()
```

FIGURE 9 – Code pour la question 5


```
[ "Virginia Beach", "Garden" ] [ "Amex", 102227.51999999997 ]
[ "Virginia Beach", "Health and Beauty" ] [ "Amex", 127597.80000000003 ]
[ "Virginia Beach", "Men's Clothing" ] [ "Amex", 108735.87000000002 ]
[ "Virginia Beach", "Music" ] [ "Amex", 116526.46000000012 ]
[ "Virginia Beach", "Pet Supplies" ] [ "Amex", 106946.26999999989 ]
[ "Virginia Beach", "Sporting Goods" ] [ "Amex", 119690.86999999995 ]
[ "Virginia Beach", "Toys" ] [ "Amex", 119357.82000000001 ]
[ "Virginia Beach", "Video Games" ] [ "Amex", 116171.34999999999 ]
[ "Virginia Beach", "Women's Clothing" ] [ "Amex", 112379.27000000002 ]
[ "Washington", "Baby" ] [ "Amex", 111320.35000000002 ]
[ "Washington", "Books" ] [ "Amex", 110236.47999999998 ]
[ "Washington", "CDs" ] [ "Amex", 112097.25000000006 ]
[ "Washington", "Cameras" ] [ "Amex", 104351.55000000005 ]
[ "Washington", "Children's Clothing" ] [ "Amex", 123122.27999999997 ]
[ "Washington", "Computers" ] [ "Amex", 113802.11999999997 ]
[ "Washington", "Consumer Electronics" ] [ "Amex", 126873.76999999995 ]
[ "Washington", "Crafts" ] [ "Amex", 118089.37999999998 ]
[ "Washington", "DVDs" ] [ "Amex", 109665.42999999995 ]
[ "Washington", "Garden" ] [ "Amex", 115971.02000000005 ]
[ "Washington", "Health and Beauty" ] [ "Amex", 109512.59 ]
[ "Washington", "Men's Clothing" ] [ "Amex", 110329.37000000018 ]
[ "Washington", "Music" ] [ "Amex", 103795.06999999995 ]
[ "Washington", "Pet Supplies" ] [ "Amex", 120491.55000000003 ]
[ "Washington", "Sporting Goods" ] [ "Amex", 109360.71999999996 ]
[ "Washington", "Toys" ] [ "Amex", 119453.29000000002 ]
[ "Washington", "Video Games" ] [ "Amex", 119790.58000000002 ]
[ "Washington", "Women's Clothing" ] [ "Amex", 116062.88999999996 ]
[ "Wichita", "Baby" ] [ "Amex", 107603.76000000008 ]
[ "Wichita", "Books" ] [ "Amex", 107206.35000000003 ]
[ "Wichita", "CDs" ] [ "Amex", 118838.21999999994 ]
[ "Wichita", "Cameras" ] [ "Amex", 111703.85 ]
[ "Wichita", "Children's Clothing" ] [ "Amex", 107897.08999999992 ]
[ "Wichita", "Computers" ] [ "Amex", 117490.87000000002 ]
[ "Wichita", "Consumer Electronics" ] [ "Amex", 108186.13999999996 ]
[ "Wichita", "Crafts" ] [ "Amex", 118452.23000000003 ]
[ "Wichita", "DVDs" ] [ "Amex", 108189.82000000004 ]
[ "Wichita", "Garden" ] [ "Amex", 113218.40000000002 ]
[ "Wichita", "Health and Beauty" ] [ "Amex", 116416.51999999999 ]
[ "Wichita", "Men's Clothing" ] [ "Amex", 102608.17 ]
[ "Wichita", "Music" ] [ "Amex", 98765.84999999996 ]
```

FIGURE 10 – Résultat pour la question 5

On peut voir que ce moyen de paiement est quasi systématiquement "Amex"

2.2 Exercice 2 - Anagramme

Étant donné un fichier de mots, écrire un script map-reduce qui détecte les mots ayant exactement les mêmes lettres (mais dans un ordre différent).

Cet algorithme devra fonctionner sous environnement Hadoop, avec plusieurs mapper et reducer en parallèle. L'algorithme ne devra pas tenir compte de la présence éventuelle de majuscules dans les mots.

```

anagram.py > ...
1  from mrjob.job import MRJob
2
3  class MRAnagram(MRJob):
4
5      def mapper(self, _, line):
6          # On convertit les mots en liste de lettre minuscule puis on ordonne cette liste
7          letters = list(str(line).lower())
8          letters.sort()
9
10         # on utilise comme clé la liste de lettre ordonnée
11         yield letters, line
12
13     def reducer(self, _, words):
14         # on renvoie les mots ayant au moins 2 anagrammes en matchant par liste ordonnée
15         anagrams = [w for w in words]
16         if len(anagrams) > 1:
17             yield len(anagrams), anagrams
18
19
20 if __name__ == "__main__":
21     MRAnagram.run()

```

FIGURE 11 – Code pour la partie 2

```

2  ["sprouts", "stupors"]
2  ["sports", "strops"]
2  ["prossy", "spyros"]
3  ["stroup", "sprout", "stupor"]
2  ["sporty", "sproty"]
5  ["sprot", "prost", "strop", "ports", "sport"]
3  ["pours", "roups", "porus"]
2  ["pros", "spor"]
2  ["pottur", "prutot"]
2  ["port", "trop"]
2  ["pour", "roup"]
3  ["pyro", "pory", "ropy"]
2  ["pro", "por"]
4  ["tossup", "spouts", "stoups", "uptoss"]
3  ["posts", "stops", "spots"]
2  ["poss", "sops"]
2  ["outputs", "putouts"]
2  ["spouty", "outspy"]
4  ["stoup", "potus", "spout", "pouts"]
3  ["typos", "topsy", "potsy"]
6  ["post", "tops", "opts", "pots", "spot", "stop"]
2  ["soupy", "pousy"]
2  ["opus", "soup"]
3  ["pows", "swop", "wops"]
2  ["posy", "opsy"]
3  ["ops", "pos", "sop"]
2  ["output", "putout"]
2  ["toup", "pout"]
3  ["opt", "pot", "top"]
2  ["pow", "wop"]
2  ["po", "op"]
2  ["torr", "rort"]
4  ["rousts", "tussor", "trouss", "stours"]
2  ["strows", "worsts"]
3  ["sorts", "ssor", "ssort"]
2  ["sorus", "sours"]
3  ["trouts", "tutors", "strout"]
2  ["torts", "trots"]
7  ["routs", "sutor", "roust", "torus", "tours", "stour", "ru
4  ["worst", "strow", "works", "trows"]

```

FIGURE 12 – Résultat pour la partie 2