



COLLÈGE  
ROSEMONT

# Valorisation de données

## Par : Abderrazak Sahraoui

# Sommaire

- Programmation Python pour Hadoop
- Installation MRJob
- Exemples de programmes MapReduce

# Programmation MapReduce en Python

## Why mrjob?

### Overview

mrjob is the easiest route to writing Python programs that run on Hadoop. If you use mrjob, you'll be able to test your code locally without installing Hadoop or run it on a cluster of your choice.

Additionally, mrjob has extensive integration with Amazon Elastic MapReduce. Once you're set up, it's as easy to run your job in the cloud as it is to run it on your laptop.

<https://mrjob.readthedocs.io/en/latest/index.html>

- Le framework **Hadoop** a été écrit en **Java**. Cependant, Il est tout a fait possible d'écrire des **job MapReduce** dans un langage autre que Java.
- **Python** est largement utilisé pour faire des programmes pour Hadoop. Python possède un module (bibliothèque) appelé **mrjob**. Cette bibliothèque permet de tester sur un ordinateur local (pc) les programmes Python écrits pour Hadoop.
- L'installation de **mrjob** sur un poste local se fait par le programme **pip**

# Installation mrjob

## pip

pip is the [package installer for Python](https://pip.pypa.io/en/latest/). You can use it to install packages from the [Python Package Index](https://pypi.org/) and other indexes.

<https://pip.pypa.io/en/latest/>

- L'installation de **mrjob** sur un poste local se fait par l'invite de commande de la façon suivante

```
CA: Invite de commandes
Microsoft Windows [version 10.0.19044.2604]
(c) Microsoft Corporation. Tous droits réservés.

C:\Users\asahraoui>py -m pip install mrjob==0.7.4
```

- On choisit ici une version compatible 0.7.4 avec notre plateforme big data HDP 2.5
- L'installation du module **pathlib** permet de gérer les chemins d'accès aux fichiers sur tout système d'exploitation

```
C:\Users\asahraoui>py -m pip install pathlib
```

# Exemple de job MapReduce en Python

```
1 """ Programme pour compter les mots d'un texte et afficher pour chaque mot le nombre d'occurence"""
2 from mrjob.job import MRJob
3 from mrjob.step import MRStep
4
5 class CompteMots(MRJob):
6     def steps(self):
7         return [
8             MRStep(mapper=self.mapper_get_mots,
9                   reducer=self.reducer_compte_mots)
10        ]
11
12     def mapper_get_mots(self, _, line):
13         line = line.strip()
14         line = line.lower()
15         line = line.replace(",", " ")
16         line = line.replace(".", " ")
17         line = line.replace(";", " ")
18         line = line.replace("?", " ")
19         line = line.replace("!", " ")
20         line = line.replace(":", " ")
21
22         mots = line.split()
23         for mot in mots:
24             yield mot, 1
25
26     def reducer_compte_mots(self, mot, uns):
27         yield mot, sum(uns)
28
29 if __name__ == '__main__':
30     CompteMots.run()
```

Pour executer ce programme, on utilise un terminal ou l'invite de commande et on soumet la ligne suivante :

```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL
PS C:\Users\asahraoui> cd documents\valorisation
PS C:\Users\asahraoui\documents\valorisation> Py CompteMots.py VictorHugo.txt
No config found, falling back on auto-configuration
```

# Exemple de fichier texte à traiter

```
C: > Users > asahraoui > Documents > Valorisation > VictorHugo.txt
1  JeanThe Project Gutenberg EBook of Les misérables Tome I, by Victor Hugo
2
3  ...
4
5  L'hôte alors se pencha à son oreille, et lui dit d'un accent qui le fit
6  tressaillir:
7
8  --Allez-vous en.
9
10 Le voyageur était courbé en cet instant et poussait quelques braises
11 dans le feu avec le bout ferré de son bâton, il se retourna vivement,
12 et, comme il ouvrait la bouche pour répliquer, l'hôte le regarda
13 fixement et ajouta toujours à voix basse:
14
15 --Tenez, assez de paroles comme cela. Voulez-vous que je vous dise votre
16 nom? Vous vous appelez Jean Valjean. Maintenant voulez-vous que je vous
17 dise qui vous êtes? En vous voyant entrer, je me suis douté de quelque
18 chose, j'ai envoyé à la mairie, et voici ce qu'on m'a répondu.
19 Savez-vous lire?
20
21 En parlant ainsi il tendait à l'étranger, tout déplié, le papier qui
22 venait de voyager de l'auberge à la mairie, et de la mairie à l'auberge.
23 L'homme y jeta un regard. L'aubergiste reprit après un silence:
24
25 --J'ai l'habitude d'être poli avec tout le monde. Allez-vous-en.
26
27 L'homme baissa la tête, ramassa le sac qu'il avait déposé à terre, et
28 s'en alla. Il prit la grande rue. Il marchait devant lui au hasard,
29 rasant de près les maisons, comme un homme humilié et triste. Il ne se
```

<https://www.gutenberg.org/cache/epub/17489/pg17489.txt>

# Exemple de résultats sur un petit extrait

```
"demi-heure" 1
"derri\u00e8re" 1
"des" 2
"dessinait" 1
"deux" 1
"devant" 2
"devin\u00e9" 1
"dise" 2
"distinguer" 1
"dit" 4
"doigt" 1
"donne" 2
"doublant" 1
"douceur" 1
"dout\u00e9" 1
"du" 7
"ebook" 1
"effet" 1
"elle" 1
"en" 8
"encore" 1
"entour\u00e9" 1
"entourait" 1
"entra" 1
"entre" 2
"entrer" 2
"envoy\u00e9" 1
"escoublon)" 1
"esp\u00e8ce" 1
"est" 2
"et" 25
"eux" 1
"faim" 1
"faisait" 3
"fatigu\u00e9" 1
"fatigue" 2
```

## Exemple 2 : Fréquences triées

```
1  """ Programme pour compter les mots d'un texte et et les afficher en ordre croissant selon la fréquence """
2  from mrjob.job import MRJob
3  from mrjob.step import MRStep
4
5  class CompteMots(MRJob):
6      def steps(self):
7          return [
8              MRStep(mapper=self.mapper_get_mots,
9                     reducer=self.reducer_compte_mots),
10             MRStep(reducer=self.reducer_frequence_mots) ]
11
12     def mapper_get_mots(self, _, line):
13         line = line.strip()
14         line = line.lower()
15         line = line.replace(",", " ")
16         line = line.replace(".", " ")
17         line = line.replace(";", " ")
18         line = line.replace("?", " ")
19         line = line.replace("!", " ")
20         line = line.replace(":", " ")
21
22         mots = line.split()
23         for mot in mots:
24             yield mot, 1
25
26     def reducer_compte_mots(self, mot, uns):
27         yield str(sum(uns)).zfill(3), mot
28
29     def reducer_frequence_mots(self, mot, frequencies):
30         for frequency in frequencies :
31             yield frequency, mot
32
33 if __name__ == '__main__':
34     CompteMots.run()
```



# Fréquences triées

```
"cabaret"      "006"  
"feu"          "006"  
"l'homme"      "006"  
"ne"           "006"  
"pas"          "006"  
"quelque"      "006"  
"sur"          "006"  
"du"           "007"  
"rue"          "007"  
"dans"         "008"  
"en"           "008"  
"les"          "008"  
"lui"          "008"  
"que"          "008"  
"se"           "008"  
"vous"         "008"  
"par"          "009"  
"son"          "009"  
"qui"          "010"  
"un"           "011"  
"une"          "013"  
"\u00e0"        "017"  
"il"           "020"  
"et"           "025"  
"le"           "029"  
"la"           "033"  
"de"           "039"
```