Raphaël Gavache

# TF-IDF

Les deux textes sont dans ~/input
Il y a 4 jobs MapReduce:

1. **WordCount**
On transforme les majuscules en minuscules puis en on retire les caractères spéciaux ()/:;$^… en filtrant avec la regex "[^\\w\\s]"
 input -> output1
 Calcule le total de chaque mot par texte et écris :
 *old@callwild.txt*  -> *"3"*

2. **WordCountDoc**
 output1 -> output2
 Calcule le nombre total de mot dans chaque texte :
 and@callwild.txt  -> "1464;31778"

3. **DocCountWord**
 output2 -> output3
 Calcule le tf-idfs d'un mot
 the@callwild.txt  -> 0.07168

4. **SortKeyByValue**
 output3 -> output4
 Ordonne par ordre décroissant les tf-idfs avec une clef custom dans un job Map. Renvoie
 customKey  ->the@callwild.txt -> 0.07168481339291334

## Résultat - Top 20

| Mot | Texte | Tf-Idf |
|---|---|---|
| the | callwild | 0.07168 |
| the | robinson | 0.04850 |
| and | callwild | 0.04178 |
| i | robinson | 0.04233 |
| and | robinson | 0.03960 |
| to | robinson | 0.03543 |
| of | robinson | 0.02899 |
| of | callwild | 0.02740 |
| he | callwild | 0.02555 |
| was | callwild | 0.02187 |
| to | callwild | 0.02124 |
| a | callwild | 0.02064 |
| a | robinson | 0.01860 |

| his | callwild | 0.01765 |
|---|---|---|
| my | robinson | 0.01762 |
| in | callwild | 0.01687 |
| was | robinson | 0.01658 |
| in | robinson | 0.01594 |
| that | robinson | 0.01554 |
| it | robinson | 0.01517 |

## Screens