

ENSICAEN - 3A INFO, IMAGE
TP 2 - WEB SÉMANTIQUE



LAGARRIGUE Lucie & VIMONT Ludovic
September 30, 2015

1 Structure

1.1 Indexation

Nous avons donc réalisé deux classes afin d'implémenter notre solution : la classe Document et la classe Word. Voici la structure de la classe Document :

```
1 public class Document {
2     private String documentName;
3     private HashMap<Word, Double> indexes;
4     private double sumPonderationSquare;
5
6     ...
7 }
```

On retrouve donc trois attributs :

- Le nom du document
- Une HashMap contenant tout les mots du document, que l'on va associer à la valeur du TF-IDF de ce mot.
- On stock également, la valeur de la somme des pondérations au carré, comme cette dernière ne dépend pas des éléments de la recherche, on peut la calculer en avance afin d'accélérer la phase de recherche.

Concernant la classe Word, voici sa structure :

```
1 public class Word {
2     private String name;
3     private HashMap<String, Integer> occurrences = new HashMap<>();
4     ...
5 }
```

On y voit donc deux attributs :

- Le nom du mot
- Le nom du document dans lequel le mot est présent et le nombre d'occurrences dans ce document donné.

Enfin, pour réaliser la phase d'indexation, nous utilisons les trois variables suivantes :

```
1 LinkedList<Document> documents = new LinkedList<Document>();
2 HashMap<String, Integer> allWords = new HashMap<String, Integer>();
3 ...
4 LinkedList<String> stopList = hfStopList.getWords();
```

- La première c'est bien sur tout simplement la liste de tout les documents que l'on désire indexer.
- La seconde, c'est une liste de tous les mots présents dans chaque document. Elle nous permet tout simplement de compter le nombre de fois où un mot est présent dans chaque document.
- Enfin, la dernière c'est une liste de mots contenus dans une stop-list (celle proposée sur la plateforme). Lors de l'ajout d'un mot dans un document, on vérifie que ce dernier n'est pas présent dans la stop-list.

Enfin, la dernière phase de l'indexation permet d'enregistrer les résultats calculés dans un fichier texte, dont voici un extrait :

```

1 | corpus/Automne\_GA.txt 159,58
2 | cache 1,61
3 | hameaux 2,3
4 | ...
5 | pauvres 1,61
6 | mourir 2,3
7 | ####

```

On commence par fournir le nom du document, ainsi que la valeur du coefficient de Salton correspondant au document. Ensuite, on va simplement lister les mots avec la valeur de leur pondération. Enfin, on termine le fichier par un ensemble de caractères afin de prévenir de l'analyse d'un nouveau document.

1.2 Recherche

Pour l'implémentation de la recherche, nous avons une classe, la classe Search :

```

1 | public class Search {
2 |     private Map<String, Double> documentCoefs;
3 |     private List<String> words;
4 |     private HashMap<String, Map<String, Double>> saltonCoefs;
5 |     private Map<String, Double> finalValues = null;
6 |     private String indexFileName;
7 |     ...
8 | }

```

La classe dispose des attributs suivants :

- Une Map contenant le nom du document et son coefficient
- Une liste des mots contenus dans la requête de l'utilisateur
- Une HashMap mettant en relation les documents, les mots et le coefficient de Salton pour un mot dans un document donné
- Une Map contenant les informations simplifiées qui seront affichées comme résultat à l'utilisateur
- Une variable stockant le nom du fichier contenant l'index, qui sera donné au constructeur de la classe

Parmi ces attributs, il n'y a que le nom du fichier que la classe ne construira pas.

La recherche est composé de quatre étapes principales :

- On extrait les mots de la requête de l'utilisateur et on les stocke.
- On ouvre et lit l'index pour rechercher ses mots. A chaque fois qu'on rencontre un document, on stocke les informations de ce dernier, même si le mot n'est pas compris dedans. Si l'un des mots de la requête est compris dans le titre, on augmente grandement la pondération du document. Cette valeur est arbitraire. Lorsqu'on rencontre un mot, on regarde s'il est dans la requête. Si oui, on le stocke en reliant le document dans lequel il est. On enregistre également sa pondération.
- Lorsqu'on a tous ces éléments, on calcule le coefficient de Salton pour chaque document et on le stocke dans une nouvelle liste.
- On peut ensuite trier les résultats. Pour cela, on utilise le nombre de mots correspondant dans chaque document, leur pondération, le coefficient de Salton. On a juste à mettre ensemble tout ce qui a été récupéré ou calculé précédemment. On peut donc ensuite afficher la liste des documents dans lesquels la requête a été trouvé dans un bon ordre.

2 Limite(s) de notre structure

Un des principales limites de notre structure, c'est le stockage des mots par HashMaps dans la classe Document, lorsque l'on va parcourir le fichier afin de compter les occurrences d'un mot la HashMap n'étant pas prévu pour pouvoir récupérer une clé grâce à sa valeur, cela nous oblige à faire un parcours de cette dernière, afin de pouvoir mettre à jour le mot existant au lieu d'en recréer un nouveau.

La présence de HashMaps rend aussi le tri par valeur plus difficile : en effet, java impose de créer une classe qui va hériter de la classe Comparator.

3 Facultatifs

Nous avons donc, réalisé plusieurs choses en plus afin d'améliorer la performance de notre moteur de recherche. Pour l'indexation, nous avons utilisé les stop-lists afin de réduire l'impact des mots grammaticaux qui sont de toute évidence moins intéressants.

Afin de pouvoir améliorer la pertinence de nos recherches nous utilisons le titre du fichier, si ce dernier possède des termes de la requête dans son nom, on considère alors que le fichier aura sans doute de l'importance, on ajoute alors un poids de 500 (valeur totalement arbitraire, qu'on aurait pu faire varier en fonction de paramètres comme la taille du document), afin de renforcer l'importance de ce document par rapport aux autres.

Nous aurions pu également implémenter une interface web à notre moteur de recherche, pour ce faire nous aurions pu prévoir de fournir des arguments à un .jar qui aurait été appelé par exemple par une page web en php, notamment grâce à des fonctions comme exec. Toutefois, nous avons manqué de temps.

4 Captures d'écran

```
Entrez r pour rechercher un ou plusieurs mots
Entrez i pour lancer l'indexation
Entrez q pour quitter
> i
Entrez le nom du fichier 0 enregistrer l'index :
index
corpus/2 Nuit.txt
corpus/Les colchiques_GA.txt
corpus/Chanson d automne PV.txt
corpus/3 Je découvre le livre.txt
corpus/Ecoutez la chanson bien douce_PV.txt
corpus/2 Demain dès l'aube.txt
corpus/2 Harmonie du soir.txt
corpus/Saltimbanques_GA.txt
corpus/3 Le bulletin de santé.txt
corpus/Les effarés AR.txt
corpus/2 Il était une fois.txt
corpus/Le dormeur du val AR.txt
corpus/3 Saint irenee de Lyon.txt
corpus/2 La fin de Satan.txt
corpus/La courbe de tes yeux PL.txt
corpus/2 Ballade à la Lune.txt
corpus/Automne GA.txt
corpus/2 Clair de lune.txt
corpus/3 Soir sur la terrasse.txt
corpus/3 Francis Cabrel.txt
corpus/L' amoureuse PL.txt
corpus/2 Chanson d automne.txt
corpus/Le pont Mirabeau GA.txt
corpus/3 La musique et la nuit.txt
corpus/2 Descente de Géants.txt
corpus/3 OCM Madame Bovary.txt
corpus/2 La lune blanche.txt
corpus/3 Livres audio iPhone.txt
corpus/3 Le desert des soirs.txt
corpus/3 Les rives romanesques.txt
Fin de l'indexation en : 1091 ms.
Entrez r pour rechercher un ou plusieurs mots
Entrez i pour lancer l'indexation
Entrez q pour quitter
> |
```

```
Entrez r pour rechercher un ou plusieurs mots
Entrez i pour lancer l'indexation
Entrez q pour quitter
> r
Entrez le nom du fichier texte contenant l'index :
index.txt
Entrez la requête :
chanson violons
Résultats :
corpus/Ecoutez la chanson bien douce_PV.txt
corpus/3 La musique et la nuit.txt
corpus/2 Ballade à la Lune.txt
corpus/Chanson d automne PV.txt
corpus/2 Chanson d automne.txt
corpus/2 Clair de lune.txt
corpus/Automne_GA.txt
La recherche a pris : 24 ms.
Entrez r pour rechercher un ou plusieurs mots
Entrez i pour lancer l'indexation
Entrez q pour quitter
> |
```

L'utilisateur choisit le nom du fichier pour l'index pour la recherche et l'indexation. L'indexation renvoie à l'utilisateur la liste de tous les fichiers indexés. La recherche quant à elle affiche la liste des fichiers où les mots ont été trouvés.

5 Performances

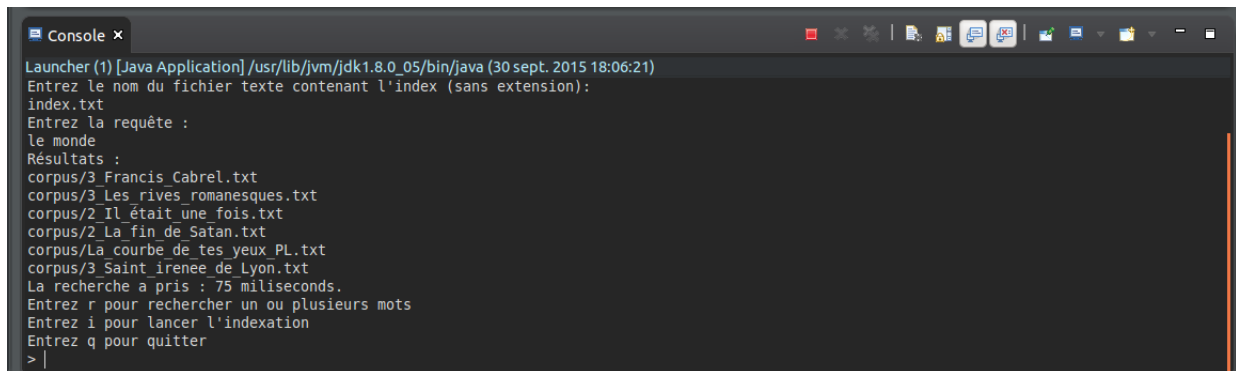
Nous avons décidé de mettre à l'épreuve notre moteur de recherche, nous avons rajouté une vingtaine de documents, notamment un livre de 2000 lignes (le fichier Saint irenee de Lyon). Pour faire les tests, nous avons utilisé la classe System de java :

```
1 | long startTime = System.currentTimeMillis();
2 | ...
3 | long stopTime = System.currentTimeMillis();
4 | System.out.println("Temps écoulé : " + (stopTime - startTime) + " ms.");
```

Nous avons alors mesuré le temps effectif passé sur l'indexation et la recherche grâce au nombre de clocks écoulés entre les actions. Concernant l'indexation, voici quelques chiffres :

- Pour les 10 fichiers de bases, l'indexation prend 100 ms.
- Pour 20 fichiers assez légers, elle prend environ 150 ms.
- Pour les 30 fichiers, elle explose en passant à 800 ms.

On se rend alors vite compte que le principe de l'indexation est lourd et que l'optimiser pourrait être très compliqué. A l'inverse une fois l'indexation effectué notre recherche est assez rapide, elle prend seulement 75 ms. Comme, on peut le voir sur la capture suivante.



```
Console x
Launcher (1) [Java Application] /usr/lib/jvm/jdk1.8.0_05/bin/java (30 sept. 2015 18:06:21)
Entrez le nom du fichier texte contenant l'index (sans extension):
index.txt
Entrez la requête :
le monde
Résultats :
corpus/3_Francis_Cabrel.txt
corpus/3_Les_rives_romanques.txt
corpus/2_Il_était_une_fois.txt
corpus/2_La_fin_de_Satan.txt
corpus/La_courbe_de_tes_yeux_PL.txt
corpus/3_Saint_irenee_de_Lyon.txt
La recherche a pris : 75 milliseconds.
Entrez r pour rechercher un ou plusieurs mots
Entrez i pour lancer l'indexation
Entrez q pour quitter
> |
```