

Mise en place d'un logiciel de recherche documentaire via la méthode booléenne en Java

Université Angers - Faculté des Sciences

Enseignant : Bernard Levrat

Groupe : Dorian Coffinet, Thibault Gauthier, Yassine Badih

I) Environnement de développement

L'ensemble du projet est codé dans un environnement Linux. Les machines utilisent un environnement Ubuntu. Nous avons choisi de coder le projet à l'aide du langage JAVA 7. Il nous semblait plus facile de le réaliser dans ce langage car chacun d'entre nous est à l'aise dans ce dernier. De plus nous avons utilisé un logiciel de "versionning" pour le projet à l'aide de Git et GitHub. Il nous a paru essentiel d'avoir ce type d'outils lorsqu'on travaille à plusieurs sur un projet.

II) Répartition des tâches

a) Développement itératif

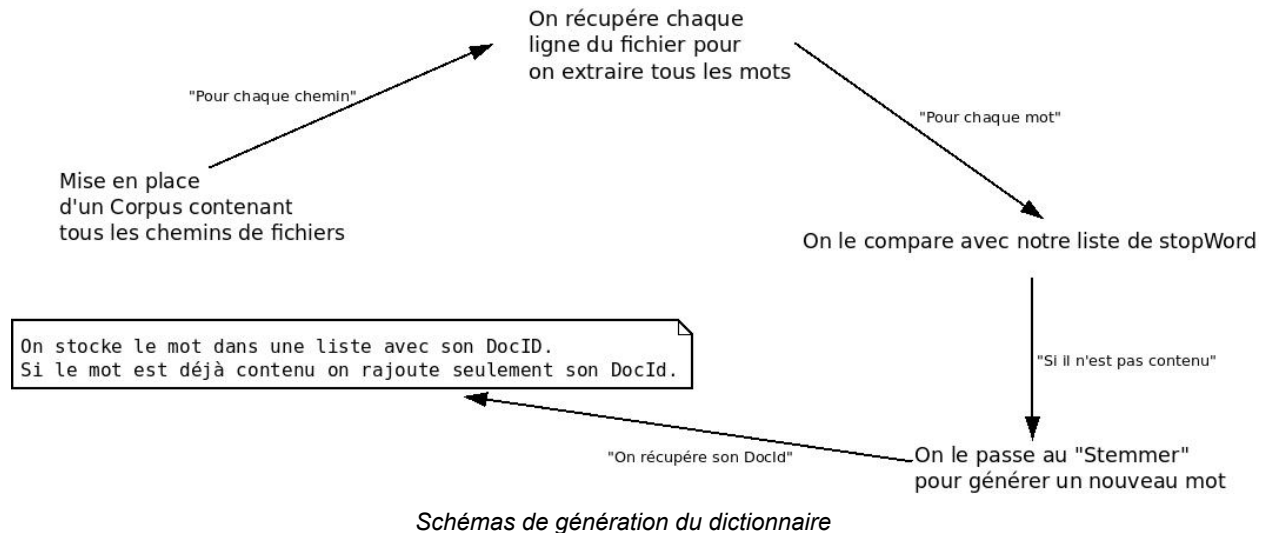
Pour arriver à réaliser le projet dans le temps imparti nous avons découpé le projet en plusieurs phases. Ces dernières sont scindées en différentes tâches :

1. Lecture des fichiers
 - Lecture de l'ensemble des textes du corpus
2. "Stemming" pour un fichier donné
 - Création du parser (Stopword)
 - Création du "Stemmer"
3. Gestion du dictionnaire
 - Représentation du dictionnaire
 - Écriture du dictionnaire
4. Gestion de la "Query" utilisateur
 - Création d'un algorithme booléen
 - Création d'un algorithme par défaut
5. Interface graphique
 - Création d'une interface graphique

Entre chaque phase nous avons fait un travail de mise en commun pour relier chaque code. Ceci nous a permis à chaque étape d'avoir une version stable et fonctionnelle.

III) Fonctionnement du projet

a) Génération du dictionnaire



Le dictionnaire crée est stocké sous forme de fichier .txt. Il se présente sous la forme : motStemmer[DocID1, DocID2, DocID3].

Pour un corpus de 364 documents d'un poids total d'environ 110Mo nous obtenons un dictionnaire avec 128383 entrées d'un poids d'environ 15Mo.

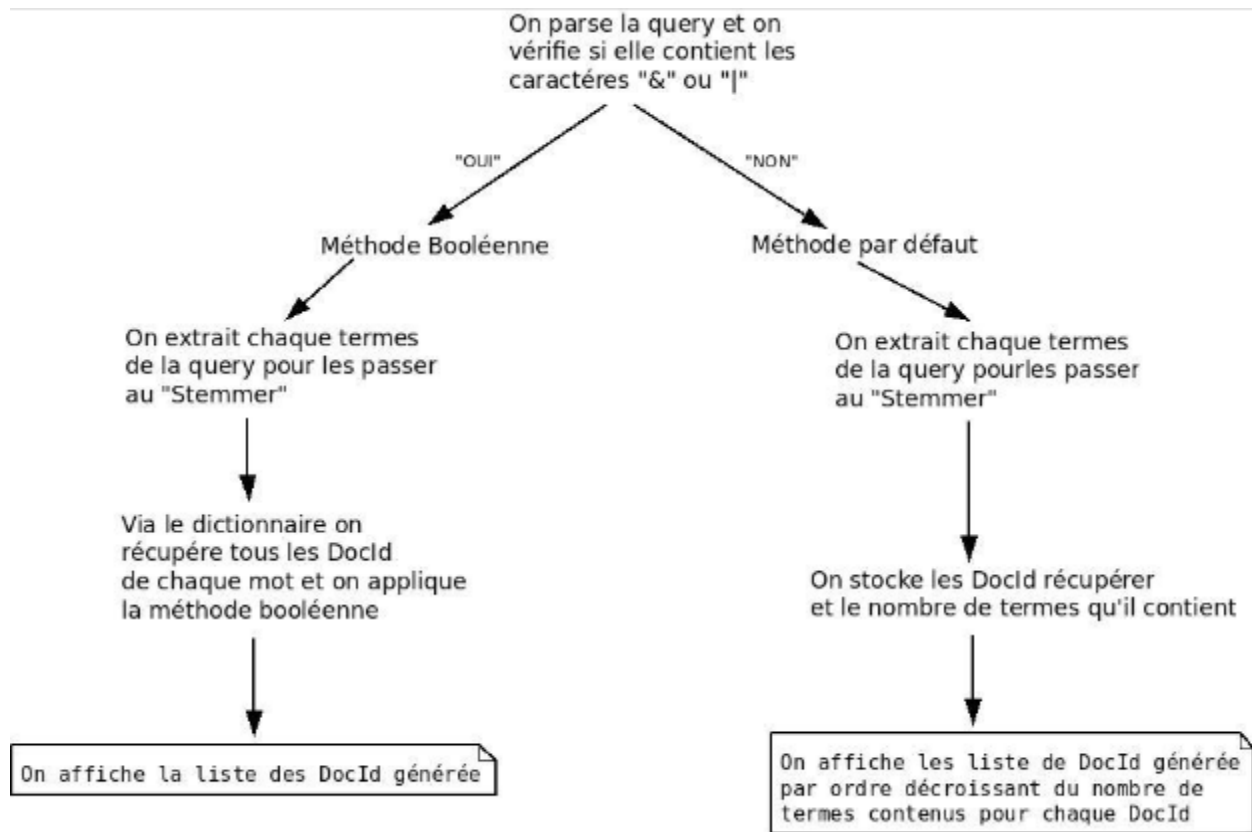
b) Les méthodes utilisées

Nous avons utilisé deux méthodes pour la recherche documentaire via notre programme:

1. Méthode Booléenne : Cette méthode reprend les fonctionnalités vu en cours pour faire la recherche. Nous avons implémenté un algorithme pour la recherche avec "ou" et un algorithme avec le "et". A la fin de la recherche nous affichons l'ensemble des documents classés par DocId croissants.
2. Méthode par défaut : Cette méthode est utilisé si il n'y a pas de mot "et" ou "ou" dans la query. L'algorithme opère ainsi :
 - Il va chercher pour chaque mots tous les DocID qui lui sont associés dans le dictionnaire => Contruction de hash avec en clef le mot et en valeur les DocID.

- Il crée une hash(H_{docID}) prenant le DocID en clef et en valeur le nombre de fois ou le DocID apparaît en valeur des hash créés à l'étape précédente.
- On trie H_{docID} par leur valeur de manière croissante.

c) Déroulement du traitement d'une Query



Schémas de traitement d'une query

Le schéma explique les deux scénarios possibles lorsque l'utilisateur lance une recherche. Le résultat de la recherche booléenne avec un "|" (un OR) est le même que le résultat de la méthode par défaut, seul l'ordre d'affichage est différent. La méthode par défaut affiche les documents par ordre décroissant du nombre de termes de la query contenus de le fichier tandis que la recherche booléenne affiche les documents par ordre alphabétique.

IV) Améliorations et tests

a) Optimisation du code

L'optimisation s'est portée sur l'articulation entre la lecture des fichiers du corpus, la phase de parsing et la phase de stemming. Au début le temps de génération du dictionnaire prenait environ 30 minutes. C'était du au fait que pour chaque fichiers on les passaient entièrement au parser puis au stemmer, ensuite nous passons le résultat pour qu'il soit ajouté au dictionnaire. De plus, pour chaque fichiers nous avons un nouveau parser et stemmer. Pour améliorer le temps d'exécution nous avons donc fait un seul parser et stemmer. Ainsi nous lisons lignes par lignes le fichier puis nous passons chaque mots au stemmer qui va l'ajouter au dictionnaire. Ce processus nous permet maintenant de générer le dictionnaire en environ 20 secondes. Nous allons donc 90 fois plus vite, ce qui est un gain de temps considérable.

b) Tests sur les algorithmes

Pour voir la véracité de nos différents algorithmes de recherche nous avons effectué quelques tests. Pour les illustrer nous allons utiliser les mots: "Shannon Wright". Voici les résultats obtenu avec les différents algorithmes.

Query	Méthode Booléen(M_B)	Méthode par défaut (M_D)
Shannon	∅	44
Wright	∅	158
Shannon Wright	∅	174
Shannon Wright	174	∅
Shannon&Wright	28	∅

Tableau de comparaison de résultats sur les mots Shannon et Wright

Ici on constate :

- Les résultats de "Shannon Wright" et "Shannon|Wright" sont identiques. Cela est du au fait que la recherche par défaut va chercher tous les documents ou sont présents "Shannon" ou "Wright" ce qui se ramène au "ou" mais la méthodes par défaut affiche les résultats différemment.

- Concentrons nous sur le résultat de “*Shannon|Wright*” qui est de 174. Pour vérifier si le résultat est bon il faut faire le calcul suivant :
 $(M_D(\text{Shannon}) + M_D(\text{Wright})) - M_B(\text{Shannon\&Wright}) = (158 + 44) - 28 = 174.$
Les deux nombres de résultat sont identiques.

c) Amélioration à apporter

Ici nous avons codé une application avec une méthode booléenne. On peut améliorer cette méthode avec l'acceptation de l'opérateur “non”. L'acceptation de bi-words. On peut aussi implémenter une méthode vectorielle pour améliorer l'efficacité de la recherche.

Annexes :

Vous pouvez récupérer le projet ici :

- http://github.com/geeknbar/projet_recherche_documentaire.
- Il comprend le code ainsi que la JavaDoc.

Diagramme UML disponible à la page suivante

