

Recherche Documentaire

Racine Mattieu - Jamet Jason - Grandjean Guillaume

Faculté des Sciences

23 Mars 2015

Recherche Documentaire I

The logo for 'Finefound' is displayed in a light gray rectangular box. The word 'Finefound' is written in a serif font, with each letter having a different color: 'F' is blue, 'i' is red, 'n' is yellow, 'e' is green, 'f' is blue, 'o' is red, 'u' is purple, 'n' is yellow, and 'd' is green.

Le meilleur moteur de recherche Documentaire !

Sommaire

- 1 Introduction
- 2 Fonctionnalités implémentées
 - Indexation
 - Traitement des requêtes
 - Pertinence des documents
 - Interface de recherche

Introduction

- Cahier des charges
 - Développement d'un moteur de recherche
 - Recherches correspondant à un corpus de documents
 - Interface pour fournir sa requête
- Outils utilisés
 - Php - Mysql
 - Pecl : collection de bibliothèques PHP (pspell / stem)
 - QuickHash

Sommaire

1 Introduction

2 Fonctionnalités implémentées

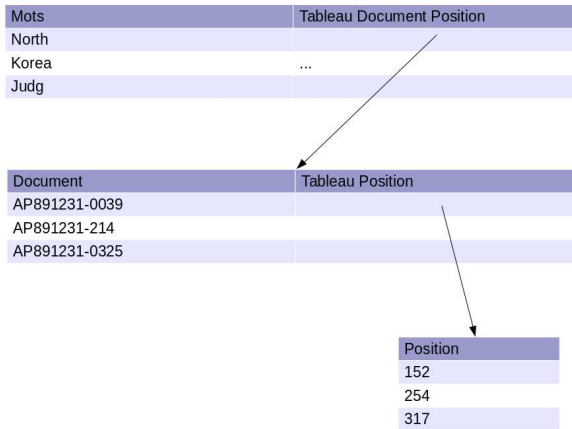
- Indexation
- Traitement des requêtes
- Pertinence des documents
- Interface de recherche

Indexation

- Tableau en mémoire
- Base de données
- Table de Hachage
- Nouvelle indexation

Indexation

Tableau en mémoire



Indexation

Base de données

chochoi document

	chochoi	position
id	: bigint(20)	
id_word	: bigint(20)	
id_document	: bigint(20)	
position	: int(11)	

  chochoi.word

	id : bigint(11)
	word : varchar(30)

Indexation

Table de Hachage

Indice	Liste
0	North
1	Korea
2	Judg
3	Democracy
4	Caesar
5	Glower
6	Amid
...	...
N	Cheer

Exemple avec le mot « Squat »

==> Présence dans la table de Hachage ?

==> Si oui Ajout dans le String de la requête Word.

==> Ajout de la position dans le String de la requête Position.

==> Ajout de la position dans le String de la requête Document.

Indexation

Nouvelle indexation



- ① Ajout d'un nouveau document au répertoire
- ② Lancer #new_index#
 - Vide les tables
 - Traite les nouveaux fichiers

Traitement des requêtes

- StopWords (a - this ...)
- KeyWords (and - or - not)
- Stemmer (judges => judg)
- Pspell (democrcy => democracy)
- Regex

Requête

And his dictatorship and makes democracy

Tableau de requête

```
array(2) [0]=> array(2) [0]=> array(2) [0]=> string(3) "his"  
[1]=> string(12) "dictatorship" [1]=> array(2) [0]=> string(5)  
"makes" [1]=> string(9) "democracy" [1]=> array(1) [0]=>  
string(3) "and"
```

Résultat

AP891216-
0001&word=dictatorship,make,democraci&stopWord=his

Pertinence des documents

- Parcours de l'index
- Calcul du taux de pertinence

Pertinence des documents

Parcours de l'index

	
Complexité de la requête	Requête et traitement simple
Execution lente	Execution et traitement rapide
Pagination possible	Pagination impossible

Calcul du taux de pertinence

Calcul du poids

TF : fréquence du terme

$$\frac{\text{“ nombre d'occurrences du terme ”}}{\text{“ nombre de mots dans le document ”}}$$

IDF : Importance des termes moins fréquents

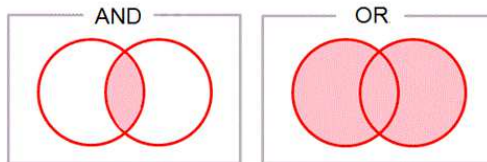
$$\log \left(\frac{\text{“ nombre total de documents ”}}{\text{“ nombre de documents où apparaît le terme ”}} \right)$$

Prise en compte de la position des mots

Calcul du taux de pertinence

Keywords

Inutilité des Keywords



Interface de recherche

- Conviviale
- Simple d'utilisation