

# R

Massih-Reza Amini - Éric Gaussier  
Préface de Stephen Robertson

# Recherche d'information

## Applications, modèles et algorithmes

### Data mining, décisionnel et big data



# Recherche d'information

*Massih-Reza Amini, professeur d'informatique à l'université Grenoble Alpes, est titulaire d'une thèse en Informatique de l'université Paris 6. Ses recherches portent sur l'apprentissage automatique appliqué aux problèmes d'accès à l'information à large échelle. Il est co-auteur de nombreux articles scientifiques notamment de l'ouvrage Apprentissage machine paru aux éditions Eyrolles. Il dirige actuellement l'équipe AMA dont les recherches se situent en analyse de données, modélisation et apprentissage automatique.*

*Éric Gaussier, professeur d'informatique à l'université Grenoble Alpes, est titulaire d'un diplôme en Mathématiques Appliquées de l'école Centrale Paris et d'une thèse en Informatique de l'université Paris 7. Ses travaux de recherche s'inscrivent dans la Science des données, au carrefour de l'apprentissage statistique, de la recherche d'information et du traitement automatique des langues. Il est co-auteur de nombreux articles et brevets dans ces domaines et dirige actuellement le laboratoire d'Informatique de Grenoble.*

**Le premier ouvrage francophone sur les algorithmes qui sous-tendent les technologies de big data et les moteurs de recherche !**

Depuis quelques années, de nouveaux modèles et algorithmes sont mis au point pour traiter des données de plus en plus volumineuses et diverses. Cet ouvrage présente les fondements scientifiques des tâches les plus répandues en recherche d'information (RI), tâches également liées au data mining, au décisionnel et plus généralement à l'exploitation du big data.

La deuxième édition de cet ouvrage propose un exposé détaillé et cohérent des algorithmes classiques développés dans ce domaine, abordable par des lecteurs qui cherchent à connaître le mécanisme des outils quotidiens d'Internet. De plus, le lecteur approfondira les concepts d'indexation, de compression, de recherche sur le Web, de classification et de catégorisation, et pourra prolonger cette étude avec les exercices corrigés proposés en fin de chapitre.

Ce livre s'adresse tant aux chercheurs et ingénieurs qui travaillent dans le domaine de l'accès à l'information et employés de PME qui utilisent en profondeur les outils du webmarketing, qu'aux étudiants de Licence, Master, écoles d'ingénieurs ou doctorants qui souhaitent un ouvrage de référence sur la recherche d'information.

## Sommaire

**Représentation et indexation.** Prétraitements linguistiques. Les deux lois de base en recherche d'information. Représentation documentaire. Index inversé. **Recherche d'information.** Modèles de recherche. Expansion de requêtes. Mesures d'évaluation. **Recherche sur le Web.** Architecture de la Toile. Trois inventions à la base du Web. Collecte et indexation des pages sur la Toile. Nouvelles stratégies de recherche. Calcul approché de similarités en grande dimension. **Catégorisation de documents.** Formalisme. Sélection de variables. Modèles génératifs. Modèles discriminants. Mesures d'évaluation. **Partitionnement de documents.** Les étapes du partitionnement. Principaux algorithmes de partitionnement. Évaluation. Applications à l'accès à l'information. **Réseaux de neurones profonds.** Neurone formel. Quelques réseaux. Applications en RI. **Recherche de thèmes latents.** Analyse sémantique latente. Analyse sémantique latente probabiliste. Le modèle LDA. **Considérations pratiques.** Logiciels libres pour la recherche d'information. Logiciels libres pour la catégorisation et le partitionnement.

**Massih-Reza Amini - Éric Gaussier**

**Préface de Stephen Robertson**

# **Recherche d'information**

**Applications, modèles et algorithmes**

**2<sup>e</sup> édition**

**EYROLLES**

The logo for EYROLLES, featuring the word "EYROLLES" in a bold, sans-serif font. Below the text is a horizontal line with a small circle in the center, resembling a stylized underline or a decorative element.

ÉDITIONS EYROLLES  
61, bd Saint-Germain  
75240 Paris Cedex 05  
[www.editions-eyrolles.com](http://www.editions-eyrolles.com)

En application de la loi du 11 mars 1957, il est interdit de reproduire intégralement ou partiellement le présent ouvrage, sur quelque support que ce soit, sans l'autorisation de l'Éditeur ou du Centre Français d'exploitation du droit de copie, 20, rue des Grands Augustins, 75006 Paris.

© Groupe Eyrolles, 2013, 2017, ISBN : 978-2-212-67376-0

# Préface

---

La recherche d'information, autrefois vue comme un domaine de spécialité à l'intersection des techniques documentaires et de la science informatique, est devenue l'une des technologies majeures du  $\text{xxi}^{\text{e}}$  siècle. Chacun s'attend en effet aujourd'hui à pouvoir trouver en quelques secondes des informations diverses sur tout type de sujet : horaires des transports en commun, principes de la production d'électricité, nature des maladies infectieuses, pharmacie la plus proche fournissant des antalgiques, films à l'affiche du cinéma voisin, analyse critique des œuvres d'Erik Satie, fondements de l'existentialisme de Jean-Paul Sartre ou tout détail trivial de la vie courante. Chacun considère cela comme allant de soi et cet « allant de soi » est né du développement des moteurs de recherche sur le Web.

Les fondements technologiques des moteurs de recherche peuvent être décrits très simplement, même si de nombreuses connaissances, combinaisons de développements théoriques et de savoir-faire expérimentaux, ont été accumulées dans ce domaine. Créer un moteur de recherche médiocre est facile ; en créer un qui soit à la fois pertinent et rapide est une tout autre histoire, et cela quelle que soit la taille de la collection considérée (collection personnelle de courriers électroniques ou intégralité du corpus de la Bibliothèque nationale). De façon étonnante, ce sont les moteurs de recherche sur Internet qui ont tenu le haut du pavé ces quelque vingt dernières années. Pour toutes sortes de raisons, ils ont atteint un niveau de maturité qui semble bien en avance sur ce qui se pratique à des échelles plus réduites.

Cet ouvrage est une introduction fondamentale à la technologie de la recherche d'information et ses applications, pour la plupart liées au Web. Il combine traitement automatique des langues et modèles théoriques. Outre l'ordonnancement de documents en réponse à une requête, il couvre la classification supervisée (en catégories prédéfinies) et non supervisée (clustering). L'importance des concepts statistiques dans ce domaine est centrale, depuis les caractéristiques statistiques des langues (loi de Zipf) jusqu'aux modèles probabilistes de recherche d'information et aux modèles à thèmes latents. Cet ouvrage était nécessaire pour mettre à la portée d'un plus large public les fondamentaux de cette technologie moderne incontournable qu'est la recherche d'information.

**Stephen Robertson**  
septembre 2012



# Table des matières

<b>Notations</b> .....	<b>IX</b>
<b>Liste des tableaux</b> .....	<b>XI</b>
<b>Liste des figures</b> .....	<b>XIII</b>
<b>Liste des algorithmes</b> .....	<b>XVII</b>
CHAPITRE 1	
<b>Introduction</b> .....	<b>1</b>
1.1 Concepts étudiés dans ce livre . . . . .	3
1.2 Organisation du livre . . . . .	7
CHAPITRE 2	
<b>Représentation et indexation</b> .....	<b>9</b>
2.1 Prétraitements linguistiques . . . . .	10
2.1.1 Segmentation . . . . .	11
2.1.2 Normalisation . . . . .	13
2.1.3 Filtrage par un antidictionnaire . . . . .	16
2.2 Les deux lois de base en recherche d'information . . . . .	19
2.2.1 Loi de Heaps . . . . .	19
2.2.2 Loi de Zipf . . . . .	20
2.3 Représentation documentaire . . . . .	22
2.3.1 Modèle vectoriel . . . . .	22
2.3.2 Pondération des termes . . . . .	23
2.4 Index inversé . . . . .	26

2.4.1	Indexation dans des collections statiques . . . . .	28
2.4.2	Indexation dans des collections dynamiques . . . . .	31
2.5	Exercices . . . . .	34

### CHAPITRE 3

<b>Recherche d'information</b> . . . . .	<b>49</b>
3.1 Modèles de recherche . . . . .	50
3.1.1 Modèles booléens . . . . .	51
3.1.2 Modèles vectoriels . . . . .	53
3.1.3 Modèles probabilistes . . . . .	57
3.1.4 Une approche axiomatique de la RI . . . . .	71
3.2 Expansion de requêtes . . . . .	72
3.2.1 La méthode « boucle de rétropertinence » . . . . .	73
3.2.2 La méthode « boucle de rétropertinence en aveugle » . . . . .	75
3.3 Mesures d'évaluation . . . . .	75
3.3.1 Évaluation de résultats non ordonnés . . . . .	76
3.3.2 Évaluation de résultats ordonnés . . . . .	78
3.4 Exercices . . . . .	84

### CHAPITRE 4

<b>Recherche sur le Web</b> . . . . .	<b>105</b>
4.1 Architecture de la Toile . . . . .	106
4.2 Trois inventions à la base du Web . . . . .	106
4.2.1 Langage HTML . . . . .	107
4.2.2 Protocole de transfert hypertexte et adresses Web . . . . .	109
4.3 Collecte et indexation des pages sur la Toile . . . . .	110
4.3.1 Robot d'indexation . . . . .	111
4.3.2 Index distribués . . . . .	114
4.4 Nouvelles stratégies de recherche . . . . .	115
4.4.1 Modèle d'apprentissage automatique pour la RI . . . . .	116
4.4.2 PageRank . . . . .	119
4.5 Calcul approché de similarités en grande dimension . . . . .	122
4.5.1 MinHash . . . . .	123
4.5.2 Hachage local pour la recherche d'un document proche . . . . .	126
4.6 Exercices . . . . .	128

CHAPITRE 5	
<b>Catégorisation de documents</b> .....	<b>137</b>
5.1 Formalisme .....	138
5.2 Sélection de variables .....	140
5.2.1 Le seuillage sur la mesure <i>Document Frequency</i> (df) .....	141
5.2.2 L'information mutuelle ponctuelle (IMP) .....	141
5.2.3 L'information mutuelle (IM) .....	143
5.2.4 La mesure $\chi^2$ .....	144
5.3 Modèles génératifs .....	146
5.3.1 Modèle multivarié de Bernoulli .....	147
5.3.2 Modèle multinomial .....	150
5.4 Modèles discriminants .....	153
5.4.1 Modèle logistique .....	157
5.4.2 Séparateurs à vaste marge .....	158
5.5 Mesures d'évaluation .....	163
5.6 Exercices .....	165
CHAPITRE 6	
<b>Partitionnement de documents</b> .....	<b>175</b>
6.1 Définitions .....	176
6.2 Les étapes du partitionnement .....	177
6.3 Principaux algorithmes de partitionnement .....	182
6.3.1 Partitionnement à plat : méthodes de réallocation .....	182
6.3.2 Partitionnement hiérarchique .....	190
6.4 Évaluation .....	200
6.5 Applications à l'accès à l'information .....	202
6.6 Exercices .....	203
CHAPITRE 7	
<b>Réseaux de neurones profonds</b> .....	<b>215</b>
7.1 Neurone formel .....	217
7.2 Quelques réseaux .....	219
7.2.1 Perceptron .....	219
7.2.2 ADALINE .....	222
7.2.3 Perceptrons Multicouches (PMC) .....	223

7.3 Applications en RI . . . . .	228
7.3.1 Réduction de dimension . . . . .	228
7.3.2 Représentation vectorielle des mots . . . . .	230
7.4 Exercices . . . . .	231
CHAPITRE 8	
<b>Recherche de thèmes latents</b> .....	<b>239</b>
8.1 Analyse sémantique latente . . . . .	241
8.1.1 Décomposition en valeurs singulières . . . . .	241
8.1.2 L'analyse sémantique latente pour la RI . . . . .	243
8.1.3 Limitations . . . . .	245
8.2 Analyse sémantique latente probabiliste . . . . .	245
8.2.1 Remarques . . . . .	247
8.3 Le modèle LDA . . . . .	249
8.4 Exercices . . . . .	252
CHAPITRE 9	
<b>Considérations pratiques</b> .....	<b>257</b>
9.1 Logiciels libres pour la recherche d'information . . . . .	258
9.1.1 dpSearch . . . . .	258
9.1.2 Lucene/SolR . . . . .	258
9.1.3 MG . . . . .	258
9.1.4 Terrier . . . . .	260
9.1.5 Zettair . . . . .	260
9.2 Logiciels libres pour la catégorisation et le partitionnement . . . . .	260
9.3 Le passage à l'échelle ou le Big Data . . . . .	261
9.3.1 Traitement parallèle et distribué . . . . .	261
9.3.2 Traitement de flux de données . . . . .	262
<b>Bibliographie</b> .....	<b>263</b>
<b>Index</b> .....	<b>271</b>

# Notations

---

$\mathcal{C} = \{d_1, \dots, d_N\}$	Collection contenant $N$ documents
$\mathcal{V} = \{t_1, \dots, t_V\}$	Vocabulaire constitué de $V$ termes
$M$	Nombre total de mots dans $\mathcal{C}$
$M_y$	Nombre total de types de mots dans $\mathcal{C}$
$M_{Nor}$	Nombre de types de mots après racinisation
$\mathbf{X}_{V \times N}$	Matrice termes-documents
$\text{tf}_{t,d}$	Nombre d'occurrences du terme $t$ dans $d \in \mathcal{C}$
$\text{ntf}_{t,d}$	Nombre d'occurrences normalisé du terme $t$ dans le document $d$
$\text{tf}_{t,\mathcal{C}}$	Nombre d'occurrences du terme $t$ dans la collection $\mathcal{C}$
$\text{tf}_{t,q}$	Nombre d'occurrences du terme $t$ dans la requête $q$
$ d $	Nombre de termes différents dans $d$
$l_d$	Nombre total d'occurrences des termes du vocabulaire dans $d$ ; $l_d = \sum_{t \in \mathcal{V}} \text{tf}_{t,d}$
$l_{\mathcal{C}}$	Nombre total d'occurrences des termes du vocabulaire dans $\mathcal{C}$ ; $l_{\mathcal{C}} = \sum_{d \in \mathcal{C}} \sum_{t \in \mathcal{V}} \text{tf}_{t,d}$
$\text{df}_t$	Nombre de documents de la collection contenant le terme $t$ ( $\text{idf}_t = \log \frac{N}{\text{df}_t}$ )
$\mathbf{d} = (w_{id})_{i \in \{1, \dots, V\}}$	Représentation vectorielle du document $d$
$w_{id}$	Poids du terme d'indice $i$ du vocabulaire dans $d$

$s(q, d)$	Score de similarité entre une requête $q$ et un document $d$ (ce score est parfois noté $RSV(q, d)$ , où $RSV$ représente la <i>Retrieval Status Value</i> )
<b>I</b>	Matrice identité
$R_{d,q}$	Jugement de pertinence binaire assigné au document $d$ pour la requête $q$
$\mathcal{A}_q$	Ensemble des documents (ou <i>alternatifs</i> ) de $\mathcal{C}$ concernés par la requête $q$
$P(\cdot)$	Une distribution de probabilité
$P(t   d)$	Probabilité de présence du terme $t$ dans $d \in \mathcal{C}$
$P(t   \mathcal{C})$	Probabilité de présence du terme $t$ dans la collection
$\mathcal{Y} = \{1, \dots, K\}$	Ensemble de $K$ étiquettes de classes, apprises ou données
$\mathcal{Z} = \{z_1, \dots, z_K\}$	Ensemble de $K$ thèmes latents
$df_t(k)$	Nombre de documents de la classe d'étiquette $k$ contenant le terme $t$
$S = ((\mathbf{d}_j, c_j))_{i=1}^m$	Ensemble d'apprentissage contenant $m$ documents avec leur étiquette de classe
$S_k$	Ensemble des documents de l'ensemble $S$ appartenant à la classe d'étiquette $k$
$N_k(S)$	Nombre de documents de l'ensemble $S$ appartenant à la classe d'étiquette $k$
$\mathcal{D}$	Distribution de probabilité suivant laquelle les exemples d'apprentissage sont identiquement et indépendamment générés
$\mathcal{F}$	Une classe de fonctions $\mathcal{F} = \{f : \mathbb{R}^V \rightarrow \mathcal{Y}\}$
$\langle \cdot, \cdot \rangle$	Produit scalaire
$\hat{R}_m(f, S)$	L'erreur empirique de la fonction de prédiction $f$ sur l'ensemble d'apprentissage $S$ de taille $m$
$R(f)$	L'erreur de généralisation de la fonction de prédiction $f$
$\{G_1, \dots, G_k\}$	Ensemble de $K$ classes obtenues sur une collection de documents $\mathcal{C}$
$\{\mathbf{r}_1, \dots, \mathbf{r}_K\}$	$K$ vecteurs représentant les centres de gravité des classes; $\forall k, \mathbf{r}_k = \frac{1}{ G_k } \sum_{\mathbf{d} \in G_k} \mathbf{d}$
IMP	Information mutuelle ponctuelle
$IM(t, c)$	Information mutuelle du terme $t$ dans la classe $c$
$[\pi]$	Fonction indicatrice, égale à 1 si le prédicat $\pi$ est vrai et 0 sinon
$\mathbb{E}(X)$	Espérance mathématique de la variable aléatoire $X$

# Liste des tableaux

---

2.1	Les 42 mots les plus fréquents (et peu informatifs) présents dans la collection du Wikipédia français. . . . .	17
2.2	Quelques statistiques sur la collection du Wikipédia français. Par <i>collection pré-traitée</i> , on entend la collection <i>segmentée, normalisée et filtrée</i> . Les nombres moyens sont arrondis par défaut. . . . .	18
2.3	Différentes variantes du codage <i>tf-idf</i> , proposées dans SMART (Salton 1975). $Char_d$ correspond à la taille du document $d$ , en nombre de caractères le constituant. . . . .	25
3.1	Les trois méthodes d'estimation de la probabilité d'apparition d'un terme dans un document les plus répandues dans les modèles de langue. . . . .	66
3.2	Mesures de rappel et de précision sur un ensemble de 10 documents ordonnés d'après la réponse d'un moteur de recherche $\mathcal{M}$ pour une requête fictive $q$ . $R_{d_{rg},q}$ correspond à la valeur du jugement de pertinence associé au document $d$ situé au rang $rg$ de la liste ordonnée, par rapport à la requête $q$ . . . . .	80
3.3	Différentes paires de valeurs $(\tau_{rg}, r_{rg})$ calculées pour l'exemple jouet du tableau 3.2, où $\tau_{rg}$ est le taux de documents non pertinents ordonnés avant le rang $rg$ et $r_{rg}$ est le <i>rappel</i> au rang $rg$ (gauche), ainsi que la courbe <i>ROC</i> correspondante (droite). . . . .	82
4.1	Quelques exemples d'instructions à placer dans le fichier <i>robots.txt</i> pour indiquer aux robots d'indexation les parties d'un site à explorer ou à éviter lors de leurs collectes de pages. . . . .	114
5.1	Quelques statistiques sur la collection de Reuters-RCV2 français. . . . .	140

5.2	Tableau de contingence comptabilisant les nombres de documents appartenant, ou non, à la classe $c$ et contenant, ou non, le terme $t$ d'une base d'entraînement de taille $m$ . . . . .	142
5.3	Tableau de contingence représentant simultanément deux caractères observés $X$ et $Y$ sur une même population. . . . .	145
5.4	Comparaison entre les quatre méthodes de sélection de variables (df, IMP, IM et $\chi^2$ ). . . . .	146
5.5	Tableau de contingence indiquant les nombres de bonnes et de mauvaises prédictions du classifieur associé à une classe $k$ d'après le jugement d'un expert. . . . .	163
5.6	Comparaison entre les micro et macromoyennes des mesures $F_1$ des différents modèles génératifs et discriminants sur la base Reuters RCV-2 français (tableau 5.1). Les documents sont représentés dans l'espace vectoriel obtenu avec les 50 000 termes les plus informatifs d'après la mesure IM. Chaque expérience est répétée 10 fois en sélectionnant aléatoirement les bases d'entraînement et de test de la collection initiale avec les proportions mentionnées dans le tableau 5.1. Chaque performance reportée dans ce tableau est la moyenne des performances obtenues sur les bases de test ainsi échantillonnées. . . . .	165
6.1	Paramètres de la formule de Lance-Williams pour différentes mesures d'agrégation . . . . .	194
8.1	Exemple de quelques termes dans cinq thèmes latents de la collection de Reuters-RCV2 (tableau 5.1) obtenus avec la méthode PLSA. Le nombre de thèmes latents, $K$ , était fixé à 40. . . . .	248
9.1	Logiciels et distributions <i>open source</i> les plus utilisés en RI pour la recherche documentaire. . . . .	259

# Liste des figures

---

1.1	Les faits marquants de l'évolution de la RI, relatés dans cet ouvrage, à partir du premier système créé pendant la Seconde Guerre mondiale jusqu'en 2010. <i>SIGIR</i> est la conférence par excellence du domaine. <i>ECIR</i> et <i>CORLA</i> sont les conférences européenne et francophone en RI. <i>ICTIR</i> est une conférence internationale sur la théorie de la RI. . . . .	2
2.1	Constitution du vocabulaire, index inversé des termes et représentation des documents dans l'espace des termes pour une collection de documents donnée. . .	11
2.2	Illustration de la loi de Heaps sur la collection du Wikipédia français. . . . .	20
2.3	Illustration de la loi de Zipf sur la collection du Wikipédia français (gauche). Pour plus de visibilité, nous avons utilisé une double échelle logarithmique, où $\ln$ est le logarithme népérien. La droite qui interpole le mieux les points, au sens des moindres carrés, est d'équation $\ln(fc) = 17,42 - \ln(\text{rang})$ . À droite, nous avons reporté les mots de la collection dont le rang vaut au plus dix, ainsi que leurs fréquences. . . . .	21
2.4	Représentation par sac de mots. Dans le document $d$ , les indices des termes correspondent à leurs indices dans la liste des termes constituant le vocabulaire. . . . .	23
2.5	Création de l'index inversé pour une collection statique constituée de 3 documents, $\mathcal{C} = \{d_1, d_2, d_3\}$ et de 5 termes, $\mathcal{V} = \{t_1, t_2, t_3, t_4, t_5\}$ . On suppose ici que l'ordre alphabétique est induit par l'ordre entre les indices des termes, i.e. le terme $t_1$ est avant le terme $t_2$ , etc. . . . .	29
2.6	Illustration de l'indexation distribuée. Nous supposons ici que l'indexation de la collection a lieu lorsque cette dernière ne subit pas de modifications dans le temps. . . . .	31
2.7	Illustration de la stratégie de mise à jour directe sur place. . . . .	33
2.8	Illustration de la stratégie de fusion. . . . .	34
3.1	Les différentes étapes de la recherche d'information . . . . .	50

3.2	Résultats de recherche booléenne pour trois requêtes formées par les termes <i>maison</i> , <i>appartement</i> et <i>loft</i> et les opérateurs ET, OU et SAUF. . . . .	53
3.3	Forme typique de la répartition des mots au sein d'une collection. . . . .	69
3.4	Boucle de rétropertinence sur la base de la figure 3.1 . . . . .	74
3.5	Description schématique des mesures rappel et précision. Le sous-ensemble des documents <i>pertinents</i> par rapport à un besoin d'information est illustré par des hachures et l'ensemble des documents retournés par le système comme pertinents par rapport à ce même besoin d'information est montré par des points. . . . .	77
3.6	Courbes précision-rappel (en pointillé) et précision interpolée (en trait plein) obtenues à partir de l'exemple du tableau 3.2. . . . .	81
4.1	Différentes étapes de la recherche sur la Toile avec les deux éléments ( <i>a</i> ) de collecte et d'indexation et ( <i>b</i> ) d'interaction et de recherche, séparés par des pointillés.	107
4.2	Exemple d'un document HTML (en haut à gauche), de sa structure interne (à droite) et de la page interprétée et affichée par un navigateur (en bas à gauche). Dans la structure du document, la portée de chaque élément est montrée par la couleur du rectangle le contenant. . . . .	110
4.3	Graphes représentatif de 659 388 pages du Wikipédia . . . . .	111
4.4	Schématisme de la procédure d'apprentissage d'une fonction de score $f$ sur une base d'entraînement constituée de trois requêtes $\{q_1, q_2, q_3\}$ ainsi que des jugements de pertinence associés sur une collection de documents $\mathcal{C}$ donnée (figure tirée de Usunier (2006)). . . . .	117
4.5	Un graphe dirigé représentant 6 pages Web ainsi que leurs liens et composé de deux sous-graphes <i>gauche</i> (nœuds $p_1, p_2$ et $p_3$ ) et <i>droite</i> (nœuds $p_4, p_5$ et $p_6$ ). . . . .	120
4.6	Stockage du vocabulaire dans un tableau de taille fixe pour chaque entrée. . . . .	129
4.7	Stockage du vocabulaire dans une chaîne de caractères. . . . .	130
5.1	Catégorisation de documents en deux phases : <i>a</i> ) entraînement d'un classifieur à partir d'une collection de documents étiquetés et <i>b</i> ) prédiction des étiquettes de classe des documents d'une base test avec le classifieur appris. . . . .	139
5.2	Performances micromoyennes des mesures $F_1$ (section 5.5) des modèles génératifs Naive-Bayes multivarié de Bernoulli et multinomial en fonction de la taille du vocabulaire sur la base Reuters RCV2 français (tableau 5.1). L'axe des abscisses est en échelle logarithmique. . . . .	153
5.3	Trois fonctions de coût pour un problème de catégorisation à deux classes en fonction du produit $c \times h$ , où $h$ est la fonction apprise. Les fonctions de coût sont l'erreur de classification : $\mathbb{1}[c \times h \leq 0]$ , le coût exponentiel : $e^{-c \times h}$ et le coût logistique : $\frac{1}{\ln(2)} \times \ln(1 + \exp(-c \times h))$ . Des valeurs positives (négatives) de $c \times h$ impliquent une bonne (mauvaise) classification. . . . .	156
5.4	Hyperplans pour un problème de classification linéairement séparable en dimension 2. Les vecteurs de support sont encerclés. . . . .	160

5.5	Hyperplans linéaires pour un problème de classification non linéairement séparable. Lorsqu'un exemple d'apprentissage est mal classé, il est considéré comme vecteur support. Sa distance à l'hyperplan de sa classe est $\frac{-\xi}{\ \lambda\ }$ . . . . .	162
6.1	Illustration des notions de <i>classes</i> et de <i>prototypes</i> . . . . .	178
6.2	Processus de construction d'un dendrogramme. . . . .	191
6.3	Exemple de dendrogramme obtenu avec une méthode non monotone. . . . .	195
7.1	Illustration d'un neurone formel. . . . .	217
7.2	Architecture d'un perceptron inspiré du système perceptif et composé de quatre composantes principales : la rétine, les fonctions d'association, les poids synaptiques et l'unité à seuil . . . . .	220
7.3	Illustration de la règle de mise à jour de l'algorithme du perceptron (Eq. 7.5) avec l'exemple $(\mathbf{d}_3, -1)$ choisi, qui est mal classé par l'hyperplan de vecteur normal $w^{(l)}$ . . . . .	222
7.4	Illustration des solutions trouvées par les algorithmes du perceptron (en pointillés) et de l'adaline (en trait plein) pour un problème de classification linéairement séparable. . . . .	223
7.5	Architecture d'un perceptron multicouches à une couche cachée (de profondeur 2). Sur cet exemple, les paramètres des biais sont introduits par des poids liés à deux unités supplémentaires associés à la couche d'entrée et à la première couche cachée ayant respectivement les valeurs fixées $x_0 = 1$ et $z_0 = 1$ . . . . .	224
7.6	Un réseau auto-associatif de profondeur 2. Avec des fonctions linéaires comme fonctions d'activation pour les unités de la couche cachée, la représentation trouvée sur la couche cachée de ce réseau est équivalente à celle trouvée par l'analyse en composantes principales (Bouillard et Kamp 1988). . . . .	229
7.7	Réseau auto-associatif de profondeur 4. Pour chaque entrée, la représentation apprise correspond aux poids des unités de la couche de code délimitée ici en pointillés. . . . .	230
8.1	Représentation sac de mots et à base de thèmes latents. L'appartenance des termes aux thèmes fait apparaître les notions de polysémie et de synonymie. . .	241
8.2	Représentations graphiques du modèle PLSA pour les cas de paramétrisation (a) asymétrique (algorithme 22) et (b) symétrique. Les rectangles ou plateaux dans chaque figure représentent le nombre de répétitions. . . . .	246
8.3	La distribution bêta pour différentes valeurs des paramètres $a$ et $b$ . . . . .	250
8.4	Modèle graphique associé au modèle LDA. $ d $ représente le nombre de termes dans un document et $N$ le nombre de documents de la collection. Chaque rectangle représente aussi le nombre de répétitions dans le processus de génération. . . . .	251



# Liste des algorithmes

---

1	Algorithme d'indexation par bloc à base de tri . . . . .	30
2	Algorithme de fusion de deux listes inversées de termes avec l'opérateur ET . .	52
3	Modèle vectoriel de recherche - implémentation du score cosinus . . . . .	56
4	Modèle d'indépendance binaire . . . . .	60
5	Algorithme de descente du gradient pour l'ordonnancement . . . . .	118
6	Algorithme de PageRank . . . . .	121
7	Sélection de variables avec la mesure IM . . . . .	144
8	Modèle multivarié de Bernoulli, phase d'apprentissage . . . . .	148
9	Modèle multivarié de Bernoulli, phase de test . . . . .	150
10	Modèle multinomial, phase d'apprentissage . . . . .	152
11	Modèle multinomial, phase de test . . . . .	152
12	Modèle logistique, phase d'apprentissage . . . . .	157
13	Algorithme des $k$ plus proches voisins pour la catégorisation . . . . .	168
14	Algorithme d'AdaBoost . . . . .	169
15	Algorithme de la méthode à une passe . . . . .	183
16	Algorithme des k-moyennes . . . . .	185
17	Algorithme de partitionnement hiérarchique agglomératif . . . . .	197
18	Algorithme EM . . . . .	209
19	Algorithme de partitionnement hiérarchique agglomératif pour le lien simple	213
20	Algorithme de Perceptron . . . . .	221
21	Perceptron multicouches . . . . .	227
22	Modèle PLSA . . . . .	245
23	Modèle LDA . . . . .	250



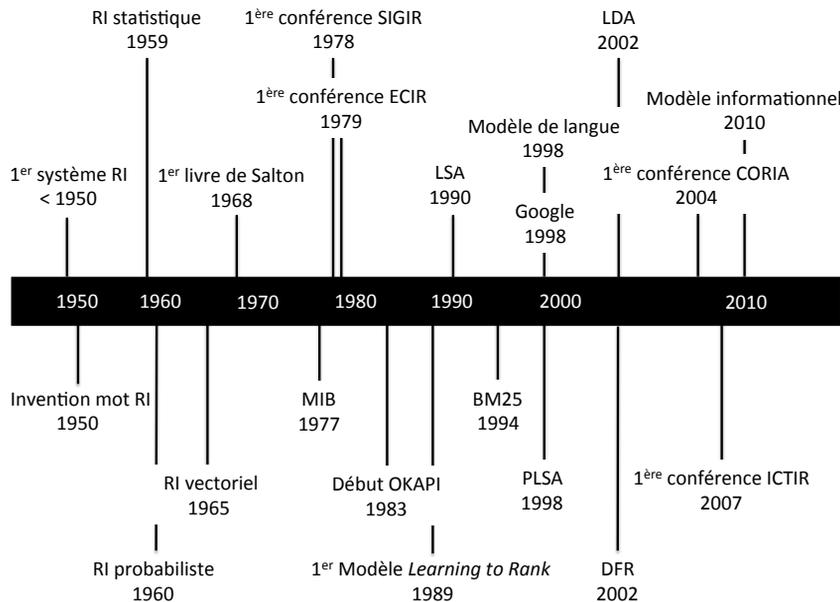
# Chapitre 1

## Introduction

---

*L'idée d'utiliser les ordinateurs pour traiter et trouver l'information enfouie dans une collection de documents et correspondant à une demande précise remonte à la fin de la Seconde Guerre mondiale. À cette époque, les militaires américains cherchaient un moyen pour indexer l'ensemble des documents scientifiques allemands afin d'y déceler toutes les informations et données pertinentes pour leurs recherches. Après la fin de la guerre, le domaine de la recherche d'information (RI) s'est développé grâce notamment à l'activité des libraires, des cabinets de juristes et d'autres professionnels travaillant sur des collections de documents spécialisées. Dès cette époque, la définition académique donnée à la tâche de RI était : la recherche de données non structurées, généralement des documents, satisfaisant une demande d'information spécifique. Par données non structurées, on entendait les données qui ne possèdent pas de structure interne précise ou qui sont sémantiquement ambiguës.*

Cette définition différencie bien cette forme de recherche de la recherche de données *relationnelles* ou *structurées*, stockées dans des bases de données telles que celles utilisées par les sociétés pour maintenir l'inventaire de leurs produits ou les dossiers de leur personnel. Dans ce cas, les données relationnelles ont une structure et une sémantique bien claires. La recherche dans ces bases s'effectue aussi en retournant les données qui satisfont exactement une expression d'algèbre relationnelle. De ce fait, dans les systèmes de gestion de bases de données, aucune erreur n'est tolérée. En revanche, pour les systèmes de recherche d'information, le but est plutôt d'ordonner les données suivant une mesure de similarité par rapport à la requête qui généralement est un ensemble de mots-clés, ne respectant aucune expression régulière ni algébrique. Le souci est alors de disposer des mesures de similarité qui feront que les données les plus pertinentes par rapport à la requête vont se trouver en tête de la liste retournée, au-dessus des données moins pertinentes ou non pertinentes. Au final, le point central en RI est de construire un



**Figure 1.1** - Les faits marquants de l'évolution de la RI, relatés dans cet ouvrage, à partir du premier système créé pendant la Seconde Guerre mondiale jusqu'en 2010. *SIGIR* est la conférence par excellence du domaine. *ECIR* et *CORIA* sont les conférences européenne et francophone en RI. *ICTIR* est une conférence internationale sur la théorie de la RI.

système qui respectera au mieux l'ordre qu'il faut attribuer aux données, du plus pertinent au moins pertinent par rapport à la requête de base.

Jusqu'au début des années 1990, le domaine de la RI est resté relativement cloisonné au monde professionnel. Plusieurs études avaient d'ailleurs montré à cette époque que le grand public préférerait s'informer auprès d'autres personnes plutôt que d'utiliser des systèmes automatiques de RI. Cette tendance s'est néanmoins inversée (et de façon irréversible) avec l'arrivée d'Internet et l'utilisation massive des moteurs de recherche. La Toile est maintenant un entrepôt universel qui a permis un partage sans précédent d'informations de toutes sortes à travers le globe. Ce développement a mis l'accent sur l'utilisation des techniques émergentes issues de la RI pour accélérer et rendre plus efficace la fouille à large échelle sur la Toile. De nos jours, la recherche d'information est très diversifiée et couvre maintenant d'autres problématiques liées à la fouille de données dans les grandes collections. Outre ses propres outils, la RI s'appuie aussi sur un ensemble de techniques issues d'autres domaines comme la statistique, l'apprentissage automatique ou le traitement automatique du langage naturel. La figure 1.1 montre les faits marquants de l'évolution de la RI, à partir du premier système créé pendant la Seconde Guerre jusqu'en 2010. Nous relatons une partie de cette évolution dans ce livre.

## 1.1 Concepts étudiés dans ce livre

Cet ouvrage présente les fondements scientifiques des tâches les plus répandues en RI à un niveau accessible aux étudiants de master et aux élèves ingénieurs. Notre souci principal a été de proposer un exposé cohérent des algorithmes classiques développés dans ce domaine, à destination des lecteurs qui cherchent à connaître le mécanisme des outils d'Internet qu'ils emploient tous les jours. Cette étude ne se limite pas à l'application initiale de RI et aborde aussi les problèmes connexes dans lesquels de nombreuses avancées techniques ont été réalisées ces dernières années. Nous allons nous intéresser plus particulièrement aux concepts décrits dans les paragraphes suivants.

### Indexation, représentation et compression

Les constructions du *dictionnaire* et de *l'index inversé*, ainsi que la représentation vectorielle des documents, constituent le point de départ dans toute manipulation ou recherche en RI. Dans une collection donnée, construire le dictionnaire ou le vocabulaire correspond à extraire une liste de termes utiles, caractéristiques des documents présents dans la collection. Cette liste servira à représenter les

documents dans un espace vectoriel. L'autre concept fondamental en RI est la constitution de l'index inversé. Il s'agit ici de construire, pour chaque terme du dictionnaire, la liste des index de documents contenant ce terme. Cette liste, aussi appelée *liste inversée*, rend l'appariement entre les requêtes et les documents de la collection plus efficace. Pour les très grandes collections de données, un problème majeur est le stockage de l'index et du dictionnaire dans la mémoire ou sur le disque. Le défi dans ces cas est de trouver un moyen de compression simple et rapide des données.

### **Recherche d'information**

Cette tâche constitue le cœur de cet ouvrage. Un système de recherche transcrit un besoin d'information donné sous la forme d'une requête constituée de mots-clés. Lorsque l'utilisateur examine le résultat de la recherche, il voit les documents triés par ordre décroissant de pertinence. Si la requête est une expression booléenne, l'utilisation de l'index inversé permet de trouver facilement et en un temps minimal tous les documents qui satisfont cette requête. En revanche, les systèmes booléens purs ne permettent pas de retrouver les documents similaires au besoin d'information de l'utilisateur et ne contenant pas exactement les termes de la requête. Plusieurs modèles ont été développés pour pallier ce problème, depuis les modèles vectoriels jusqu'aux modèles probabilistes. De même, plusieurs stratégies, qui consistent à étendre la requête afin d'y inclure des termes similaires mais non mentionnés originellement par l'utilisateur, ont vu le jour afin d'enrichir ces différents modèles.

### **Recherche sur le Web**

La Toile (ou le Web) est un entrepôt dynamique et distribué de documents qui, par sa taille, par le manque de supervision dans la génération et la suppression de documents, ainsi que par la diversité du type de ces derniers, rend la recherche bien plus difficile que la recherche traditionnelle effectuée sur des collections classiques. Les premiers moteurs de recherche sur la Toile reproduisaient néanmoins directement les méthodes de RI classiques, le défi principal étant de gérer des index inversés de très grandes tailles. La prise en compte, vers la fin des années 1990, d'une des caractéristiques essentielles du Web, à savoir les liens hypertextes reliant les documents entre eux, a permis, d'une part, de réaliser une meilleure indexation des pages Web et, d'autre part, de donner un score de notoriété à chaque page sur la base de la topologie de la Toile. Cela a conduit à la première

génération des moteurs de recherche vraiment adaptés au Web, dont Google fut le prototype. De nos jours, d'autres éléments sont pris en compte et les modèles utilisés reposent sur des techniques récentes d'apprentissage automatique.

### **Classification de documents**

Un système de classification de documents a pour but de catégoriser automatiquement une collection de documents suivant un ensemble de classes prédéfinies. Un exemple de tels systèmes est le catégoriseur de messages électroniques incorporé dans la plupart des boîtes e-mail et qui place les courriers suspects automatiquement dans le dossier des indésirables. Les systèmes de classification sont généralement conçus avec des techniques issues de *l'apprentissage statistique* et opèrent en deux phases. La première est la phase d'entraînement, lors de laquelle les paramètres du système sont réglés sur une base d'apprentissage contenant des documents avec leurs classes respectives. Le système apprend l'association entre les documents et leurs classes. C'est lors de la seconde phase, dite de *test*, que le système assigne une classe à chaque nouveau document entrant. Habituellement, les paramètres des systèmes d'apprentissage sont mis à jour périodiquement pendant le laps de temps où il n'y a pas de traitement à effectuer sur des documents arrivants.

### **Partitionnement de documents**

La tâche de partitionnement de documents (ou *document clustering* en anglais) est un autre problème important en RI. Le but est ici de réunir dans le même groupe les documents qui sont similaires par rapport à un critère donné. En pratique, les résultats de partitionnement indiquent non seulement la structure d'une collection, mais ils sont aussi souvent utilisés dans d'autres tâches de RI comme la navigation ou la recherche distribuée. La tâche de partitionnement est depuis longtemps un sujet de recherches intensives dans différents domaines, comme la *statistique*, *l'apprentissage machine*, *l'analyse de données* et la RI. Grâce à la capacité croissante des ordinateurs à traiter de grandes quantités de données dans des laps de temps réduits, ce paradigme a connu un développement rapide ces dernières années, aussi bien d'un point de vue théorique que pratique. On distingue trois types d'approches en partitionnement : les méthodes à base de similarité qui partitionnent les données en  $k$  groupes, où chaque groupe optimise un critère de partitionnement fondé sur une mesure de similarité ; les approches à base de projection qui se fondent sur la recherche des vecteurs propres d'une matrice de similarité

entre documents pour projeter puis regrouper les documents dans un espace de dimension réduite ; et, finalement, les méthodes à base de densité qui modélisent la distribution de probabilité des différents groupes avant d'assigner les documents aux différentes partitions. Nous présenterons en détail les méthodes à base de similarité, qui comptent parmi les plus utilisées.

### **Réseaux de neurones profonds**

Un réseau de neurones est un ensemble d'unités élémentaires, appelées neurones artificiels, reliées entre elles et permettant d'effectuer différentes transformations non linéaires. Depuis 2010, nous sommes témoins d'une percée importante de l'application de ces réseaux de neurones profonds dans plusieurs domaines phares en informatique comme la vision par ordinateur ou le traitement automatique du langage naturel, et il est attendu que l'apprentissage profond ait aussi un grand impact dans le domaine de la recherche d'information.

### **Recherche de thèmes latents**

Depuis le début des années 1990, on assiste à une recherche intensive sur l'extraction de *thèmes cachés* ou *latents* dans des collections de documents ou d'images. Ces études sont motivées par une meilleure prise en compte du discours dans la représentation des documents ou des images. L'approche sac de mots, très populaire et répandue en RI, consiste à caractériser les documents avec les termes qui les constituent, sans prendre en compte l'ordre ou le sens de ces termes. Ainsi, les documents qui expriment les mêmes concepts mais avec des termes différents (phénomène de synonymie) auront des représentations très différentes et ne seront pas tous considérés comme pertinents pour une requête fondée sur ces concepts. De façon similaire, les termes ayant des sens différents dans différents contextes (phénomène de polysémie) sont en général traités comme équivalents et rapprochent des documents qui ne devraient pas l'être. Les modèles latents ont été conçus pour pallier ces problèmes en exhibant les thématiques sous-jacentes et le vocabulaire qui leur est associé. Ces techniques visent ainsi à mieux représenter les documents textuels ou les images.

## 1.2 Organisation du livre

Les chapitres de ce livre décrivent les concepts présentés précédemment. À la fin de chaque chapitre, nous donnons aussi un ensemble d'exercices (corrigés) relatifs aux sujets développés. L'enchaînement des idées est le suivant :

- ▶ Dans le **chapitre 2**, nous décrivons la chaîne complète d'indexation qui nous permet de construire l'ensemble du vocabulaire en partant d'une collection de documents donnée. Nous décrivons aussi les algorithmes les plus répandus pour la construction de l'index inversé dans les collections *statiques*, *distribuées* et *dynamiques*. Dans la partie exercice, nous étudions, entre autres, un ensemble de techniques permettant de comprimer le vocabulaire et l'index inversé afin que, pour des collections de très grandes tailles, ces éléments puissent être stockés en mémoire et accessibles en temps raisonnable.
- ▶ C'est dans le **chapitre 3** que nous nous intéressons aux modèles de recherche d'information les plus courants. Nous présentons les avantages et les inconvénients de chacun, en mettant l'accent sur les éléments qui ont motivé leur évolution, depuis les modèles booléens et vectoriels jusqu'aux modèles probabilistes récents (de langue et fondés sur l'information). Beaucoup d'entre eux se retrouvent dans d'autres tâches comme la traduction automatique ou la reconnaissance de la parole.
- ▶ Le **chapitre 4** décrit les modèles et stratégies développés pour la recherche sur la Toile. Nous présentons d'abord les techniques usuelles pour collecter et indexer les documents entreposés sur différents serveurs connectés à la Toile, avant d'exposer les éléments spécifiques aux moteurs de recherche sur le Web : calcul d'un score de notoriété pour chaque page et déploiement de méthodes d'apprentissage d'ordonnancement à partir des données de clics des utilisateurs. Ce chapitre se termine avec la présentation d'une technique de calcul approché de similarités entre documents au sein de grandes collections à l'aide des tables de hachage.
- ▶ Le **chapitre 5** présente d'abord les techniques de sélection de variables utilisées en RI pour réduire la taille du vocabulaire et donc la dimension de l'espace de représentation. Nous nous intéressons ensuite aux premières techniques, qualifiées de *génératives*, introduites en RI pour la catégorisation de documents. Ces techniques visent tout d'abord à modéliser les distributions de probabilité générant les documents avant d'utiliser ces distributions pour prendre une décision sur les classes des documents. Dans la dernière partie de ce chapitre,

nous présentons les modèles *discriminants* issus principalement du domaine de l'apprentissage machine pour cette tâche.

- ▶ Nous exposons les deux approches les plus utilisées pour le partitionnement de documents dans le **chapitre 6**, à savoir les méthodes de réallocation et les méthodes de partitionnement agglomératif. Comme nous l'avons mentionné plus haut, le partitionnement a des applications dans différents domaines. Certaines d'entre elles sont également passées en revue dans ce chapitre.
- ▶ Le **chapitre 7** présente les modèles formels des réseaux de neurones avec une emphase sur la retro-propagation, qui est l'algorithme le plus connu pour apprendre les paramètres d'un réseau profond. Ce chapitre se termine par la présentation de deux applications des réseaux de neurones très utilisées en RI, à savoir la réduction de dimension et la représentation vectorielle des mots.
- ▶ C'est dans le **chapitre 8** que nous décrivons les modèles latents les plus utilisés pour l'extraction de thèmes dans les images et les documents. Nous nous intéressons tout d'abord aux premières tentatives réelles de modélisation des phénomènes de synonymie et polysémie à travers la décomposition en valeurs singulières, avant d'aborder les versions plus récentes qui étendent le cadre de base. Ces différentes versions seront illustrées aux travers des résultats qu'elles fournissent sur des collections de textes standards.
- ▶ Nous consacrons finalement quelques pages en fin d'ouvrage aux logiciels libres de recherche d'information, de catégorisation et de partitionnement.

Enfin, nous souhaitons qu'au-delà de sa vocation pédagogique, ce livre, qui s'adresse non seulement aux étudiants de master et aux élèves d'école d'ingénieurs, mais aussi aux doctorants en RI et aux ingénieurs de ce domaine, puisse donner à ses lecteurs le désir d'en savoir plus.

# Chapitre 2

## Représentation et indexation

---

*Dans ce chapitre, nous allons nous intéresser aux mécanismes amenant à la représentation documentaire et à l'indexation et qui constituent les briques de base de presque tous les algorithmes classiques en accès à l'information. Ces mécanismes comprennent un ensemble de prétraitements permettant la constitution du vocabulaire à partir d'une collection de documents, la représentation des documents dans l'espace des termes, la construction de l'index inversé des termes ainsi que la compression du vocabulaire et de l'index qui favorise l'accélération des traitements et l'accès aux données. La figure 2.1 illustre ces différents mécanismes.*

Dans la première partie de ce chapitre, nous allons passer en revue l'ensemble des prétraitements qui aboutissent à la création du vocabulaire (section 2.1). Nous présenterons ensuite les deux lois de base en accès à l'information décrivant l'évolution de la taille du vocabulaire par rapport à la taille de la collection de départ (loi de Heaps) et la distribution des termes par rapport à leur fréquence d'apparition dans une collection donnée (loi de Zipf). Dans les sections 2.3 et 2.4, nous présenterons successivement le modèle vectoriel, qui est communément utilisé pour la représentation documentaire, et les algorithmes classiques pour la constitution de l'index inversé dans des environnements statiques et dynamiques. La dernière section sera consacrée aux exercices concernant les différentes notions présentées dans ce chapitre. Nous allons aussi y traiter en particulier deux algorithmes de compression du vocabulaire et de l'index. La compression est dans ces cas sans perte de données et a pour but de stocker ces dernières de façon optimale ou d'accélérer leur transfert du disque d'une machine à sa mémoire vive. L'étude de ces algorithmes a son importance puisque les systèmes de recherche classiques parcourent plus rapidement les index compressés que leur version non compressées, et ceci parce que le transfert des parties compressées de l'index du disque et leur décompression dans la mémoire vive s'effectuent plus rapidement que le transfert direct des parties non compressées. En exemple, nous détaillerons l'impact des différents algorithmes de traitement et de compression présentés dans ce chapitre sur la collection du Wikipédia français constituée de 1 349 539 documents textuels français de plus de 20 mots visibles sur le site à la date du 27 mai 2011<sup>1</sup>.

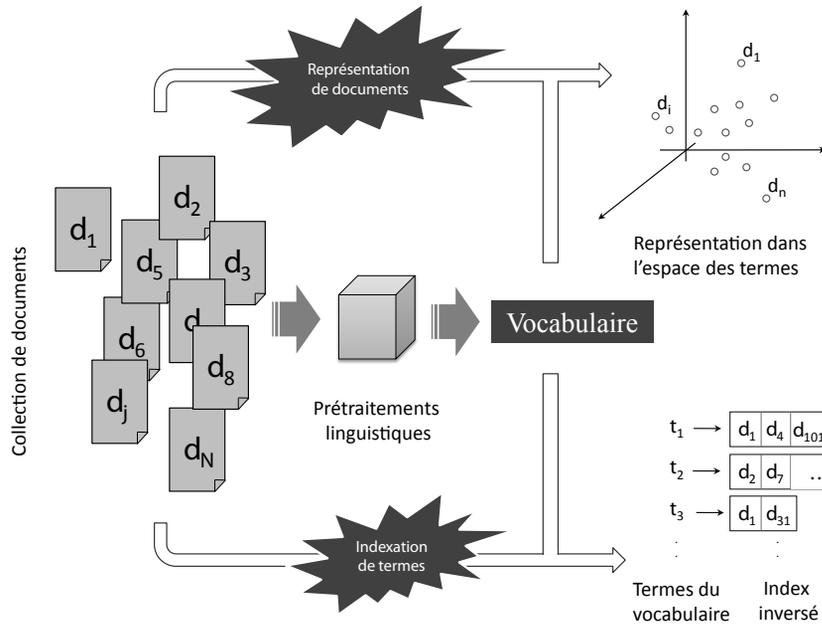
## 2.1 Prétraitements linguistiques

Dans cette partie, nous allons présenter l'ensemble des prétraitements linguistiques conduisant à la constitution de l'index inversé et de la représentation documentaire à partir d'une collection de documents donnée.

Ces prétraitements comprennent la *segmentation* des séquences de caractères présents dans les documents en mots distincts, la *normalisation* des mots, opération qui consiste à fournir une *forme canonique* pour chaque mot, et le *filtrage* qui vise

---

1. Voir <http://dumps.wikimedia.org/frwiki/latest/frwiki-latest-pages-articles.xml.bz2-rss.xml> pour les derniers transferts disponibles.



**Figure 2.1** - Constitution du vocabulaire, index inversé des termes et représentation des documents dans l'espace des termes pour une collection de documents donnée.

à supprimer les mots les plus fréquents<sup>2</sup>. Ces processus sont schématisés par le cube grisé dans la figure 2.1.

### 2.1.1 Segmentation

La segmentation (en anglais *tokenisation*) consiste à séparer une suite de caractères en éléments sémantiques, ou mots. Un *type* de mot est la classe de tous les mots ayant la même séquence de caractères et un *terme* est un type de mots que l'on garde pour former le vocabulaire.

Dans l'exemple suivant :

*Au sens étymologique, l'information est ce qui donne une forme à l'esprit.*

nous avons 14 mots :

*Au, sens, étymologique, l', information, est, ce, qui, donne, une, forme, à, l', esprit*

2. Nous réservons le nom *terme* pour désigner les mots normalisés restant après l'étape du filtrage.

mais seulement 13 types, puisqu'il y a deux instances de  $\{l'\}$ . Après filtrage par un antidictionnaire (section 2.1.3), il ne restera plus que 6 termes  $\{\textit{sens, étymologique, information, donne, forme, esprit}\}$  dans le vocabulaire.

Cette étape est difficile et cruciale puisqu'une mauvaise segmentation a un impact négatif direct sur le résultat de la recherche. En effet, si dans une collection certains termes ne sont pas indexés, les requêtes composées de ces termes ne pourront jamais être appariées avec les documents de la collection qui les contiennent.

Une bonne segmentation dépend de la prise en compte des spécificités de la langue des textes traités. Pour certaines langues asiatiques, comme le chinois, les mots dans un texte ne sont pas séparés par des espaces et la segmentation est dans ce cas une tâche ardue qui constitue un défi scientifique majeur. Pour les langues indo-européennes, la tâche de segmentation est plus aisée puisque l'espace et les signes de ponctuation donnent une première indication de séparation entre les différents éléments lexicaux. Néanmoins, chaque langue de ce groupe linguistique a sa spécificité propre et une simple segmentation par des espaces et des signes de ponctuation conduit généralement à des résultats d'indexation médiocres. Par exemple, pour le français, nous avons :

- les composés lexicaux à trait d'union comme *chassé-croisé, peut-être, rendez-vous*, etc.
- les composés lexicaux à apostrophe comme *jusqu'où, aujourd'hui, prud'homme*, etc.
- les expressions idiomatiques comme *au fait, poser un lapin, tomber dans les pommes*, etc.
- les formes contractées comme *Gad'zarts (les gars des Arts et Métiers), M'sieur, j'*, etc.
- les sigles et les acronymes comme *K7, A.R., CV, càd, P-V*, etc.

Dans ce cas, lorsque nous indexons un texte en français, nous ne souhaitons pas, par exemple, que le mot *aujourd'hui* soit séparé en *aujourd* et *hui* alors que nous voulons que les mots *un homme* et *l'homme* soient indexés sous le même terme, *homme*. Ce problème devient même extrême avec l'allemand, où les noms composés s'écrivent sans espace ; par exemple le mot *Rindfleischetikettierungsüberwachungsaufgabenübertragungsgesetz* doit, a priori, être séparé en les dix mots qui le constituent - *Rind (boeuf) + Fleisch (viande) + Etikettierung (étiquetage) + Überwachung (surveillance) + Aufgabe (tâche) + Übertragung (transmission) + Gesetz (droit)* - avant d'être indexé. Pour les langues européennes, différents logiciels de segmentation commerciaux existent. Le but de certains de ces logiciels est de réaliser une analyse plus poussée du texte en associant aux mots d'une phrase leur fonction grammaticale. La segmentation dans ces cas est une étape préliminaire à cette analyse.