



**HAL**  
open science

# Indexation de Co-Occurrences dans des Corpus de Documents Structurés et Production de Cartes Sémantiques Interactives

Pierre Pompidor, Boris Carbonneill, Michel Sala

► **To cite this version:**

Pierre Pompidor, Boris Carbonneill, Michel Sala. Indexation de Co-Occurrences dans des Corpus de Documents Structurés et Production de Cartes Sémantiques Interactives. Document numérique - Revue des sciences et technologies de l'information. Série Document numérique, 2009, 12 (1), pp.53-79. 10.3166/DN.12.1.53-79 . lirmm-00394364

**HAL Id: lirmm-00394364**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00394364>**

Submitted on 21 Oct 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Indexation de co-occurrences dans des corpus de documents structurés et production de cartes sémantiques interactives

**Pierre Pompidor \*** — **Boris Carbonneill \*\*** — **Michel Sala \***

(\*) *LIRMM, UMR Université Montpellier II – CNRS, 34090 Montpellier*

(\*\*) *Société C6, 770 av. Alfred Sauvy, 34470 Perols*

---

**RÉSUMÉ.** Confrontés à la problématique de l'indexation de très grands corpus documentaires d'entreprises, nous avons mis au point une méthode simple mais efficace (en temps de calcul et de volumétrie), permettant de filtrer par document les co-occurrences les plus représentatives de ceux-ci. Le choix d'un contexte de co-occurrences a deux raisons. D'une part les requêtes portant sur des corpus spécialisés et composées par des experts, s'appuient sur peu de termes précisément choisis dont l'indexation des associations permet la construction de cartes sémantiques de navigation dans les concepts du corpus. Pour cela nous prenons en compte de la structure des documents en validant les contenus des paragraphes par ceux de leurs titres. Notre méthode s'appuie sur des mesures tf.idf successives effectuées dans le contexte d'un document et non d'un corpus, sur les contenus des paragraphes auxquels sont intégrés progressivement la hiérarchie des titres les introduisant. Puis nous exploitons simultanément une ontologie de contrôle et les requêtes des utilisateurs comportant les termes précédemment discriminés pour valider par le théorème de Bayes, les associations sémantiques ainsi déterminées, qui finalement permettent la production de cartes sémantiques.

**ABSTRACT.** This paper addresses the problem of indexing very large enterprise corpuses. We have designed a simple yet efficient (in terms of computation time and the size of the generated results) method allowing to filter, on a per-document basis, the most representative co-occurrences of the documents. The reason for using co-occurrences is twofold. First, queries composed by experts on specialized corpuses rely statistically on few chosen terms, for which we index the associations. Second, such co-occurrences facilitate the construction of semantic maps used to navigate the concepts of the corpus. Our main approach is to take into account the structure of the documents by validating the content of the paragraphs by their titles. Our method starts with successive tf.idf measures of paragraph contents taken in the context of a document, to which we progressively integrate the hierarchy of their introducing titles. We then simultaneously exploit a control ontology and the user queries containing the terms that we discriminated in the first step in order to validate, using Bayes' theorem, the semantic associations contained in a paragraph given the terms of its title.

**MOTS-CLÉS :** *Exploitation de la structure des documents, Indexation incrémentale de très grands corpus, Contexte de co-occurrences, Théorème de Bayes, Cartes sémantiques*

**KEYWORDS :** *Structured Document Retrieval, Very large and fast corpus indexing, Co-occurrences' context, Bayes' theorem, Semantic maps*

## 1. Introduction

Cet article se place dans le champ de l'indexation de grands corpus spécialisés (comprenant plusieurs dizaines de milliers de documents). En effet, nous travaillons en collaboration avec la société C6 qui propose des solutions de gestion électronique de documents, notamment pour la constitution de dossiers d'autorisation de mise sur marché de médicaments. Or ni la recherche « full text » qui génère énormément de bruit, ni l'exploitation de méta-données trop formelles (et rarement présentes), ne satisfont les utilisateurs (experts) de ces corpus. Par ailleurs le « coût » d'une indexation classique (tant en terme de temps que de volumétrie) est trop élevé pour des entreprises dont les besoins ne sont pas prioritairement ciblés sur la recherche d'information. Nous avons donc choisi deux orientations permettant de réduire les coûts d'indexation, en terme de volumétrie et en temps de calcul, (en sachant bien entendu que l'ordre d'indexation sera très important) :

- seul un nombre très limité d'associations sémantiques est retenu par document (nous préférons le terme d'association sémantique à celui de co-occurrences qui lui peut recouvrir des collocations ne faisant pas sens) ;
- l'insertion d'un nouveau document dans le corpus ne nécessite pas la ré-indexation des documents déjà indexés.

Outre l'**économie de moyen** (en termes de temps de calcul et de volumétrie), le second objectif de notre approche et la création de **cartes sémantiques interactives** qui vont permettre aux utilisateurs-experts des corpus indexés, de naviguer dans ceux-ci. En effet, toujours dans le cadre de l'exploitation de très gros corpus, l'entreprise C6 est souvent confrontée à cette demande de la part de ses clients (notamment des laboratoires pharmaceutiques). Or même si une majorité de documents est annotée par des méta-données, ces annotations ne recouvrent généralement que des informations contextuelles à la création de ceux-ci (noms des auteurs, dates de création et de révision...), et non pas une description structurée de leur contenu. Vu l'ampleur des champs sémantiques couverts par ces gros corpus, nous avons fait le choix que l'aide à l'exploration du corpus devait être guidée par les premières requêtes exprimées par les utilisateurs. Les champs sémantiques proposés aux utilisateurs dépendront alors des contextes sémantiques associés aux co-occurrences des termes présents dans leurs requêtes.

Pour arriver à ces fins, nous proposons un nouveau processus d'indexation d'associations sémantiques entre termes (qui vont ainsi définir des contextes sémantiques), assujetti à trois niveaux de contrôles successifs :

- les lemmes candidats aux associations sont ceux qui discriminent le mieux les paragraphes du document par rapport au document lui-même, et non pas ceux qui classiquement discriminaient le mieux les documents par rapport au corpus, cette discrimination s'appuyant sur la structure du document par l'intégration progressive du contenu de la hiérarchie de titres aux contenus des paragraphes ;
- les associations sémantiques candidates sont basées sur les lemmes discriminés lors de la première étape, puis pondérées suivant une

ontologie du domaine (ou à défaut un thésaurus généraliste), et les requêtes des utilisateurs du système de gestion documentaire ;

- enfin les associations candidates sont progressivement filtrées par rapport au contenu du titre introduisant la section du document dans laquelle elles apparaissent, et cela grâce au théorème de Bayes qui les contextualise.

L'indexation étant purement incrémentale (c'est à dire qu'un document n'est indexé qu'une seule fois), nous ne pouvons prétendre à des résultats optimaux. La cohérence globale de l'indexation d'un document dans le corpus, (notamment pour éviter qu'un terme polysémique soit mal indexé), est assurée indirectement par l'analyse des requêtes des utilisateurs, requêtes elles-mêmes analysées par rapport à la proximité de leurs termes dans une ontologie support. Dans ce contexte, aucun travail n'a encore porté simultanément sur des très grands corpus de documents textuels, (la société C6 traite des corpus contenant jusqu'à un million de pages), nécessitant une indexation incrémentale, et en utilisant l'heuristique de l'analyse du contenu des titres pour associer un contexte sémantique à une section du document.

Pour illustrer notre propos et en se focalisant sur un corpus recelant un grand nombre de documents centrés sur la toxicité des champignons, nous prendrons comme fil rouge le fragment de texte ci-après. Les termes en italique sont ceux inconnus de l'ontologie support qui sera présentée plus loin.

#### La toxicité des champignons

La toxicité des champignons ingérés est d'autant plus grave que l'apparition des symptômes est tardive. Ces symptômes sont le fait de deux syndromes principaux : les syndromes *phalloïdien* et *orellanien*.

#### Le syndrome *phalloïdien* et le syndrome *orellanien* des champignons à lamelles

##### Les amanites

La reconnaissance des amanites passe principalement par l'identification d'un anneau et d'une volve. Les principales espèces responsables du syndrome *phalloïdien* sont *amanita phalloides*, *amanita verna* et *amanita virosa*.

##### Les cortinaires

Les cortinaires sont caractérisés par la présence fréquente d'une *cortine*. C'est un genre qui comprend une très grande variété d'espèces dont certaines sont comestibles, d'autres d'une grande toxicité.

Comme dans la très grande majorité des documents des corpus manipulés par la solution de gestion documentaire de la société C6, cet exemple fait apparaître des nouveaux termes inconnus dans l'ontologie et qui devront y être intégrés (ici des noms latins de champignons et de syndromes, souvent des molécules). Intuitivement en lisant ce texte, nous percevons que les mots essentiels du texte sont ceux relatifs au thème général du texte : *champignon*, *toxicité* ; aux noms des syndromes toxiques *syndrome phalloïdien* et *orellanien* ; aux noms vernaculaires des espèces citées (*amanite*, *cortinaire*) ; à l'anatomie des champignons : *lamelle*, *anneau*, *volve* et *cortine* ; et à un moindre degré : aux noms latins des espèces (*amanita phalloides*...) (ces noms seraient pertinents si le paragraphe relatif aux cortinaires en contenait). Les termes non pertinents sont ceux relatifs : à la classification des champignons (*espèce*, *genre* et *variété*), et plus encore, ceux structurant leurs descriptions (*reconnaissance*, *identification* et *présence*). Le but est de pouvoir indexer les trois

paragraphes de ce texte par des associations pertinentes (comme *champignon – toxicité* ou *anneau – volve (sachant/dans le contexte des amanites)*) en ignorant celles qui ne le sont pas (comme *reconnaissance – identification*).

Nous allons maintenant, après un bref état de l'art de l'indexation des documents textuels structurés (ch. 2), expliciter la notion de co-occurrence mise en oeuvre (ch. 3), pour ensuite détailler les objectifs et les étapes de notre méthode d'indexation (ch. 4 et 5), avant de conclure sur la génération des cartes sémantiques (ch. 6).

## **2. Petit état de l'art sur l'indexation de documents textuels structurés**

Les systèmes de recherche d'informations (RI), et plus spécialement ceux de recherche documentaire (RD), ont très longtemps utilisé des représentations de données très simples pour opérer des requêtes sur les textes, ou classer ceux-ci en différentes catégories. Dans cet état de l'art, nous ne commenterons seulement que les travaux dédiés aux tâches de recherche documentaire (trouver parmi un ensemble de documents ceux qui répondent le mieux à une requête), et non ceux dédiés à celles de classification (attribuer à chaque document une ou plusieurs catégories), ce qui ne correspond pas aux objectifs de notre travail, et d'autre part nous nous focaliserons que sur les travaux d'indexation guidée par la structure des documents.

Si les systèmes de RD ont très longtemps utilisé des représentations de données vectorielles pour opérer des requêtes sur les textes, à partir du début des années 1990, ces représentations ont commencé à prendre en compte la structure des documents pour mener des travaux sur deux axes : la "recherche de passages" et la "recherche de sous-structures". Les premiers nous intéressent peu, car ils se limitent généralement à découper un document en sous-documents, et à réappliquer à ces sous-parties les modèles habituels (souvent donc vectoriels) de la RI. La prise en compte « simultanée » du document et de ses sections pour opérer des recherches plus fines n'est introduite qu'à partir de 1994 par Wilkinson [Wilkinson, 1994]. Par la suite la RI a commencé à apporter quelques outils basés sur les réseaux bayésiens pour le traitement des documents formatés en XML (et principalement analysés lors du congrès INEX « Initiative for the Evaluation of XML retrieval »).

En effet, si les réseaux bayésiens ont fait leur apparition dans le domaine de la recherche documentaire depuis le système Inquery [Callan et al., 1992], ils ne l'ont été très longtemps que pour analyser des documents "plats", avant qu'un premier projet [Myaeng et al., 1998] ne les utilise pour vérifier que les sections ciblées par les requêtes soient bien représentatives du document (voir paragraphe 5.6). Enfin, beaucoup de travaux ont porté ces dernières années sur les indexations supportées par des réseaux bayésiens, s'appuyant sur la structure de documents semi-structurés (formatés en XML), et principalement dans des buts de classification [Piwowski, 2002] [Denoyer, 2004] [Denoyer et al., 2004] [Bratko et al., 2004]. De ces travaux, notre travail se détache principalement par le fait que notre système gère des corpus de documents textuels (et non formatés en XML), dont la sémantique portée par la structure est prioritairement prise en compte dans la recherche de contextes sémantiques (en permettant de produire des cartes sémantiques personnalisées suivant l'analyse des premières requêtes exprimées par les utilisateurs).

### 3. Contextes de co-occurrences et associations sémantiques

Notre processus d'indexation s'intéresse moins aux termes qu'aux associations sémantiques entre termes. Déjà en effet, sur un moteur de recherche généraliste, ce sont les requêtes comprenant deux mots qui sont les plus nombreuses, (mots par requête sur AOL.com en août 2006 : 1 mot : 27,5% ; 2 mots : 29% ; 3 mots : 18,7% ; 4 mots : 11,1% ; ...), et ce phénomène s'accroît sur les moteurs de recherche des bases documentaires spécialisées, utilisées par des experts exprimant des requêtes précises. Mais si les moteurs de recherche généralistes privilégient la confiance (popularité mesurée par le « pagerank ») à la fréquence des occurrences des termes recherchés, dans une base de données documentaire spécialisée cette mesure de confiance n'est plus pertinente. Par ailleurs, les moteurs de recherche adaptables à des corpus spécifiques comme Apache Lucene [LUCENE] se focalisent sur des calculs des fréquences optimisés par différentes variantes de la formule de discrimination *tf.idf* (décrite plus avant). Or présenter à l'utilisateur les documents qui maximisent la fréquence de deux termes sans tenir compte de leurs co-occurrences soulève les problèmes suivants :

- la fréquence de l'un peut-être disproportionnée à l'autre, ou leurs présences peuvent être dissociées dans deux parties distinctes du document (et donc ces termes ne sont pas réellement co-occurents) ;
- et plus globalement cette présentation n'étaye pas la construction de cartes sémantiques en créant des liens sémantiques faibles.

Depuis quelques temps déjà, des travaux sont menés sur le calcul des fréquences de co-occurrences et de la définition des contextes sous-jacents (documentaire, positionnel ou syntaxique), et améliorés en considérant les dépendances syntaxiques entre les unités linguistiques [Besançon, 2002]. Nous nous plaçons dans un contexte documentaire, la phrase étant notre unité de base pour la génération des associations candidates qui sont ensuite contextualisées sur les lemmes des titres introduisant les paragraphes dans lesquels se situent ces phrases (en utilisant le théorème de Bayes).

### 4. Objectifs et étapes du processus

Nos objectifs globaux sont :

- de discriminer les lemmes par la mesure *tf.idf* appliquée dans le cadre de chaque document en utilisant la structure de celui-ci, (hiérarchie des titres organisant les sections du document, une section étant une sous-arborescence de paragraphes) ;
- d'indexer les paragraphes par un nombre réduit d'associations sémantiques composées d'au moins un lemme discriminant, en éliminant les associations sémantiques vides ou pauvres de sens par l'exploitation d'une ontologie et des requêtes des utilisateurs-experts ;
- et dans le futur, d'insérer de nouveaux termes dans l'ontologie initiale pour la spécialiser, voire de créer des versions propres à chaque corpus.

Il faut donc retenir que les corpus manipulés étant particulièrement volumineux, et la recherche d'associations sémantiques augmentant significativement la complexité des calculs, la sélection d'association sémantique est opérée document

par document, le maintien de la cohérence globale étant dévolue à l'ontologie révisée par l'exploitation effectuée par les utilisateurs sur le corpus.

Les **problèmes de volumétrie** sont cruciaux dans des corpus de centaines de milliers de pages, chaque document atteignant fréquemment les 3000 mots, la suppression des mots creux en éliminant à peu près les trois quarts. En moyenne, chaque document comprend une cinquantaine de paragraphes organisant chacun une quinzaine de phrases, chaque phrase conservant en moyenne 5 termes après la phase de pré-traitement. La conservation globale de toutes les co-occurrences possibles reviendrait donc (toujours en moyenne) à conserver 10 co-occurrences par phrase, soit 150 co-occurrences par paragraphe et donc 7500 co-occurrences par documents ce qui en nombre de termes multiplierait par 5 son volume !

Notre but est de ne conserver en moyenne que  $n$  (par défaut 3), associations sémantiques pertinentes par paragraphe, en ramenant le poids maximal de l'indexation d'un document à 150 associations sémantiques. Le fonctionnement général du processus d'indexation est illustré par l'algorithme suivant :

\* Pour chaque document du corpus :

*Phase d'amorce par indexation simple en utilisant la structure du document :*

Pré-traitement (lemmatisation et suppression des mots creux) (voir section 5.1)

Comptage de chaque lemme dans tous les paragraphes

\* Pour chaque paragraphe :

\* Pour tous les lemmes du paragraphe (simples ou agrégés)

- **Calculs successifs du tf.idf** attribuant un poids à ce lemme **en intégrant progressivement** dans la section le contenant les lemmes des titres introductifs (voir section 5.2)

*(une section est une sous-arborescence de paragraphes, initialement chaque paragraphe feuille forme une section)*

- **Conservation du maximum de cette mesure** et transformation de celle-ci en une probabilité de « pertinence » ( $P_{\text{pertinence}}(\text{lemme})$ )

- **Mémorisation des lemmes « probablement pertinents », génération des associations sémantiques** ayant au moins un de ceux-ci (sect. 5.3)

Enregistrement des requêtes ciblant ce document dans un entrepôt de données.

A partir du moment où un nombre suffisant de requêtes concernant le document couvre tous les lemmes probablement pertinents indexés :

\* Pour chaque paragraphe :

\* Pour chaque phrase du paragraphe :

\* Pour toutes les associations sémantiques pré-sélectionnées :

- **Calcul de la probabilité de pertinence par rapport à l'ontologie support et aux requêtes des utilisateurs** :  $P_{\text{pertinence}}(\text{association})$

(l'ontologie est préalablement décomposée en composantes) (s. 5.4)

- **Mesure de la probabilité de pertinence des lemmes par rapport à cette association** :  $P_{\text{pertinence}}(\text{lemme}/\text{association})$  (section 5.5)

\* Pour tous les lemmes du titre introduisant le paragraphe :

Calcul de la probabilité de pertinence de l'association par rapport

aux lemmes du titre introduisant le paragraphe par le **théorème de**

**Bayes** (section 5.6) :  $P(\text{association}/\text{lemme}) = P(\text{lemme}/$

$\text{association}).P(\text{association}) / P(\text{lemme})$

**Indexation du paragraphe par les  $n$  premières associations (suivant les probabilités calculées précédemment)**

## 5. Détail des étapes

### 5.1 Prétraitement (lemmatisation et élimination des mots creux)

Nous procédons à une phase classique de lemmatisation des termes, puis à l'élimination des mots creux suivant les critères suivants :

- sont conservés : les substantifs, les groupes nominaux composés d'un nom et d'un adjectif, et ceux composés d'un nom suivi d'un terme inconnu (cette heuristique permet la conservation du nombre de collocations pertinentes dans les documents scientifiques, (ici *syndrome principal, phalloïdien*), et tous les groupes de mots inconnus et contigus (*amanita phalloïdes, amanita verna* et *amanita virosa*) ;
- tous les autres termes sont supprimés.

Nous paramétrons TreeTagger [TreeTagger] pour améliorer ce pré-traitement sur des syntagmes plus complexes. Sur notre exemple, les syntagmes *champignon à lamelles* ou *espèce de champignon* seront prochainement considérés globalement.

### 5.2 Calculs successifs du *tf.idf* attribuant un poids à chaque lemme du document dans une section, et corrélation d'une probabilité de pertinence

Le texte illustratif comprend plusieurs sections de texte que nous voulons caractériser par des associations sémantiques. Ces sections sont formées initialement par les paragraphes et ne comprennent pas les titres ou sous-titres qui les introduisent. Nous allons opérer cette intégration progressivement.

Après une première mesure *tf.idf* est calculée pour tous les lemmes de chaque paragraphe, une seconde mesure est opérée en absorbant le contenu des titres les plus proches de ces paragraphes qui deviennent des sections, et cela ainsi de suite jusqu'à ce que le titre de plus haut niveau soit absorbé par la section qui lui était la plus proche. Ainsi le paragraphe 2 :

La reconnaissance des amanites passe principalement par l'identification d'un anneau et d'une volve. Les principales espèces responsables du syndrome phalloïdien sont *amanita phalloïdes, amanita verna* et *amanita virosa*.

est ensuite étendu aux deux sections suivantes :

Les amanites :

La reconnaissance des amanites passe principalement par l'identification d'un anneau et d'une volve. Les principales espèces responsables du syndrome phalloïdien sont *amanita phalloïdes, amanita verna* et *amanita virosa*.

Puis :

Le syndrome phalloïdien et le syndrome orellanien des champignons à lamelles :

Les amanites :

La reconnaissance des amanites passe ...

Pour identifier les lemmes discriminant le plus fortement les sections par rapport au document, nous utilisons la fonction **Term Frequency Inverse Document Frequency** (*tf.idf*). Elle est issue du monde de la RI [Salton et al., 1988] et donne :

- plus de poids aux mots apparaissant souvent au sein d'un même document (ici dans une section, ces mots sont plus représentatifs de celle-ci) ;
- moins de poids aux mots appartenant à plusieurs documents (ici sections) en reflétant le fait que ces mots ont un faible pouvoir de discrimination.

En opérant cette mesure sur des paragraphes qui absorbent progressivement leurs titres, nous détectons les sections (paragraphes+sous-hiérarchie de titres) les mieux discriminées par les lemmes du document. Cette mesure n'est pas une probabilité. Pour calculer la probabilité qu'un lemme soit représentatif d'une section (paragraphe associé à un ou plusieurs de ses titres), l'intervalle entre les mesures *tf.idf* la plus basse et la plus haute, est corrélé à un intervalle de probabilité entre 0,25 et 0,75.

Le tableau suivant (Figure 1) présente :

- la liste des lemmes, leurs nombres d'occurrences dans les trois sections et le nombre de sections dans lesquelles ces lemmes apparaissent ;
- les mesures *tf.idf* appliquées au contenu des paragraphes (feuilles, avec titre immédiat, les deux titres les plus proches) et les probabilités corrélées (Px).

Liste des lemmes	Nbr d'occurrences dans les sections	Nbr de sections ds lesquelles le lemme apparaît (S)	tf.idf niveau 0 N/25 .log(3/S) (valeur T0)	tf.idf niveau 1 N/29 .log(3/S) (valeur T1)	tf.idf niveau 2 N/37 .log(3/S) (valeur T2)	P0 (avec T0)	P1 (avec T1)	P2 (avec T2)
Toxicité	2 - 3 - 3	2 - 2 - 2	0,014	<b>0,018</b>	0,014	0,25	<b>0,393</b>	
<b>champignon</b>	1 - 2 - 4*	1 - 1 - 3	0,019	<b>0,033</b>	0	0,354	<b>0,75</b>	
apparition	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
<b>symptôme</b>	2	1	<b>0,038</b>	0,033	0,026	<b>0,75</b>		
deux	1	1	0,019	0,016	0,013	0,35		
syndrome pr.	1	1	<b>0,019</b>	0,016	0,013	4		
syndrome ph.	2 - 2 - 4*	2 - 2 - 3	0,014	0,012	0	<b>0,354</b>		
syndrome or.	1 - 1 - 3*	1 - 1 - 3	<b>0,019</b>	0,016	0	<b>0,25</b>		
lamelle	0 - 0 - 2*	0 - 0 - 2			<b>0,0095</b>	<b>0,354</b>		<b>0,433</b>
<b>amanite</b>	1 - 2 - 2	1	0,019	<b>0,033</b>	0,026		<b>0,75</b>	
reconnaiss.	1	1	<b>0,019</b>	0,016	0,013	0,354		
identification	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
anneau	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
volve	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
espèce	2	2	<b>0,014</b>	0,012	0,0095	<b>0,354</b>		
amanita ph.	1	1	<b>0,019</b>	0,016	0,013	<b>0,25</b>		
amanita ve.	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
amanita vi.	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
<b>cortinaire</b>	1 - 2 - 2	1	0,019	<b>0,033</b>	0,026	<b>0,354</b>	<b>0,75</b>	
présence fr.	1	1	<b>0,019</b>	0,016	0,013	0,354		
cortine	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
genre	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		
variété	1	1	<b>0,019</b>	0,016	0,013	<b>0,354</b>		

Moyenne des probabilités de pertinence : 0.422 ; (\*) ces lemmes sont distribués sur les deux sous-sections

Les valeurs en gras sont les valeurs optimales retenues.

**Table 1.** Par lemme, nombres d'occurrences et mesures *tf.idf* en absorbant les titres les plus proches

Les lemmes les plus discriminants, et probablement pertinents par rapport aux sections car d'une probabilité supérieure à 0,5, sont : **champignon**, **symptôme**, **amanite** et **cortinaire**. Ils seront donc à la base de l'indexation par les associations

sémantiques initiales car seules les associations comprenant au moins un de ces lemmes seront retenues par un premier filtre.

### 5.3 Première indexation des paragraphes sur les termes lemmatisés

Voici pour la première phrase de chaque paragraphe les lemmes probablement pertinents et les associations sémantiques candidates qui les comprennent :

- La toxicité des **champignons** ingérés est d'autant plus grave que l'apparition des **symptômes** est tardive : *champignon – toxicité ; champignon – apparition ; champignon – symptôme ; symptôme – toxicité ; symptôme – apparition*

- La reconnaissance des **amanites** passe par l'identification d'un anneau et d'une volve : *amanite – reconnaissance ; amanite – identification ; amanite – anneau ; amanite – volve*

- Les **cortinaires** sont caractérisés par la présence fréquente d'une cortine : *cortinaire – présence-fréquente ; cortinaire - cortine*

### 5.4 Calcul de la probabilité de pertinence d'une association sémantique par rapport à l'ontologie et aux requêtes des utilisateurs

Une probabilité de pertinence va être calculée pour chaque association sémantique en se basant sur une **ontologie support** (à défaut un thésaurus), et sur un **premier contingent de requêtes** effectuées par les experts.

Ces deux outils sont complémentaires, l'ontologie permettant de spécifier la description des champs sémantiques notamment sur les relations de synonymie, de spécialisation et d'agrégation, et les requêtes des utilisateurs-experts ciblant leurs usages fonctionnels.

#### 5.4.1 Une ontologie pour une première probabilité de pertinence de l'association

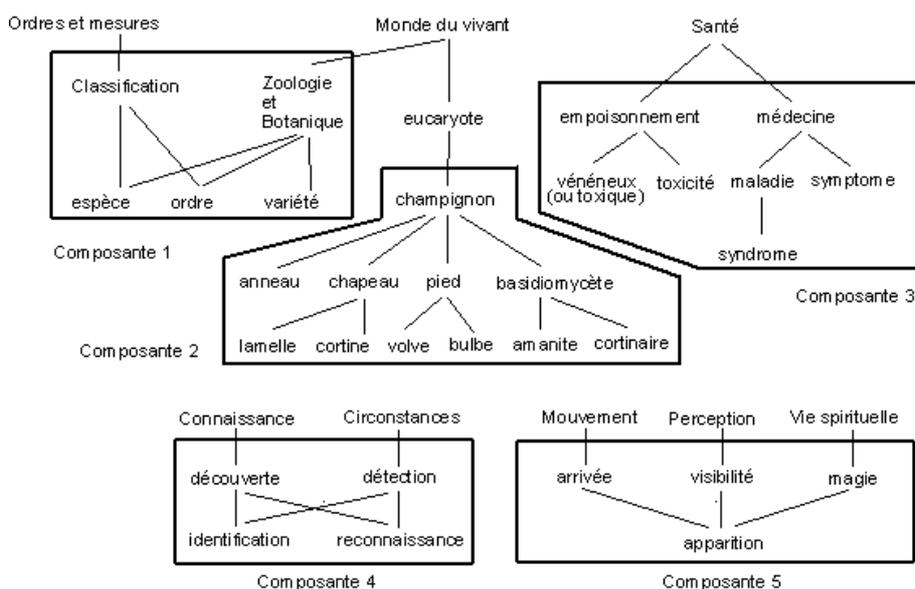
Dans notre processus, les associations sémantiques sélectionnées sont d'abord pondérées suivant un calcul de probabilité étayée par une ontologie support. Si les ontologies sont devenues indispensables à la recherche d'information en permettant de proposer aux auteurs des méta-données pour annoter les documents, d'améliorer les requêtes des utilisateurs par des mécanismes d'extension et des cheminements de recherche, nous les utilisons dans un premier temps pour valider la proximité sémantique des lemmes associés. Notre ontologie, comme souvent dans le domaine de la fouille de texte, se réduit à un « simple » thésaurus explicitant trois relations :

- la relation de synonymie (« *lépiote élevée* » synonyme de *coulemelle*) ;
- la relation d'hyponymie (*lépiote* hyperonyme de *coulemelle*) ;
- la relation d'association de sens (*amanite* en relation avec *toxicité*) (bien que certaines amanites soient excellentes)

qui ne sont d'ailleurs pas forcément explicitées.

Cela-dit, il est rare de pouvoir disposer d'ontologies adaptées aux documents manipulés : nous ne disposons souvent que de thésaurus linguistiques trop généralistes qui ne couvrent que partiellement les champs sémantiques des corpus spécialisés (ici la phytotoxicité), ou au contraire de taxonomies trop focalisées...

Considérons les fragments de l'ontologie (*Figure 2*) dont nous disposons, et qui organisent le sens des lemmes de notre exemple. Ces fragments au nombre de cinq sont en fait des composantes du graphe sous-jacent à l'ontologie :



(\*) est appelée ici **composante** les sous-graphes denses comprenant des nœuds fortement connectés.

**Figure 2.** Fragments de l'ontologie relevant des champs sémantiques couverts par l'exemple

Nous remarquerons que cette ontologie possède une superstructure qui met en relation toutes les composantes, et que si elle recense en très grande majorité des substantifs, certains adjectifs sont également présents (ici *vénéneux/toxique*). Par ailleurs, même si les ontologies peuvent être manipulées comme des graphes complets, il est naïf de pouvoir mesurer une proximité sémantique entre deux lemmes de l'ontologie en calculant simplement le nombre de pas du chemin le plus court reliant le premier lemme d'une association au second de celle-ci **en dehors d'une composante**, et cela à cause de deux raisons. D'une part, le thésaurus n'est jamais homogène (certains champs sémantiques sont toujours plus détaillés que d'autres), notamment quand il est progressivement complété par du vocabulaire spécialisé, et d'autre part les liens de superstructure sont souvent factices.

Par rapport à notre exemple, notre procédure a donc isolé cinq composantes, soit cinq sous-graphes de nœuds fortement inter-connectés :

- la composante *c1* des termes de classification : *espèce*, *genre* et *variété* ;
- la composante *c2* du terme *champignon* et de ses agrégats (*lamelle*, *anneau*, *volve*...) ou de ses spécialisations (*amanite* et *cortinaire*) ;
- la composante *c3* centrée autour du terme *toxicité* ;
- la composante *c4* centrée autour d'*identification* et de *reconnaissance* ;

- et enfin la composante *c5* centrée autour du terme *apparition*.  
 (Utiliser les connexions entre ces composantes est dépourvu de signification, même le lien entre les deux composantes *monde du vivant* et *santé* est peu prégnant).

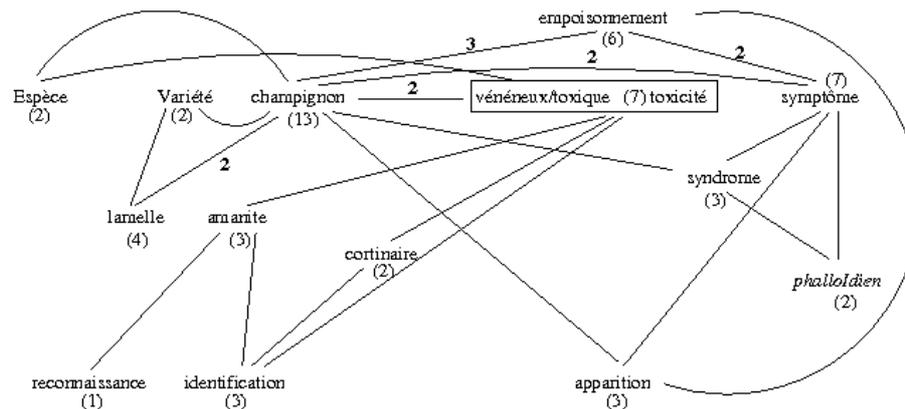
Par rapport à notre document, nous pouvons déjà inventorier les lemmes qui sont potentiellement synonymes (en pouvant être employés l'un pour l'autre **dans ce contexte**). Ce sont ceux, qui au sein d'une composante ont les mêmes pères : *espèce* – *variété*, *lamelle* – *cortine*, *amanite* – *cortinaire* ... *identification* – *reconnaissance*.

#### 5.4.2 Prise en compte des requêtes des utilisateurs

Les associations sémantiques associées aux paragraphes sont filtrées par rapport aux termes lemmatisés composant les requêtes des utilisateurs (*Figure 3*), dont voici les treize premières comprenant au moins un lemme « probablement » pertinent, et ayant sélectionné le document d'illustration :

amanite – identification	espèce – champignon – toxique
amanite – toxicité	identification – cortinaire – toxique
champignon – empoisonnement	symptôme – empoisonnement – champignon
champignon – syndrome	symptôme – syndrome-phalloïdien
reconnaissance – amanite	variété – champignon – lamelle
champignon – lamelle – toxique	apparition – symptôme – empoisonnement – champignon

En effet, à partir du moment où tous les lemmes probablement pertinents ont été utilisés dans les requêtes des utilisateurs, celles-ci participent aux calculs des probabilités de pertinences des associations sémantiques. Le premier prétraitement consiste à fixer les termes synonymes dans le contexte de notre document : ce sont ceux qui pré-identifiés peuvent être interchangeables dans un nombre significatif de requêtes. Sur les six couples de termes pré-identifiés, les termes *identification* et *reconnaissance* peuvent être substitués dans deux requêtes : ils sont donc confondus. Une analyse morpho-syntaxique préalable à l'exploitation de l'ontologie a aussi permis de confondre les noeuds *véneux-toxique* et *toxicité*.



**Figure 3.** Graphe des requêtes des utilisateurs (l'arité des nœuds étant spécifiée)

En partant du nœud ayant la plus forte arité (*champignon*) sont progressivement englobés tous les nœuds connexes ayant une arité supérieure ou égale à l'arité moyenne (de 3,9). Est ainsi créé un sur-nœud (Figure 4) qui représente la principale composante sémantique fonctionnelle liée au document et qui :

- confond les deux composantes *c2* et *c3* (*monde du vivant* et de la *santé*) ;
- confère une forte probabilité de pertinence (0,75) à toutes les associations formées par tout couple de lemmes pris dans ce sous-graphe ;
- crée une transitivité entre des lemmes joignables par son intermédiaire.

Tous les nœuds dont l'arité est supérieure à l'arité moyenne sont absorbés. Autrement, l'identification de sur-nœuds (composantes sémantiques fonctionnelles) serait réitérée sur les nœuds restants d'arité supérieure à la moyenne.

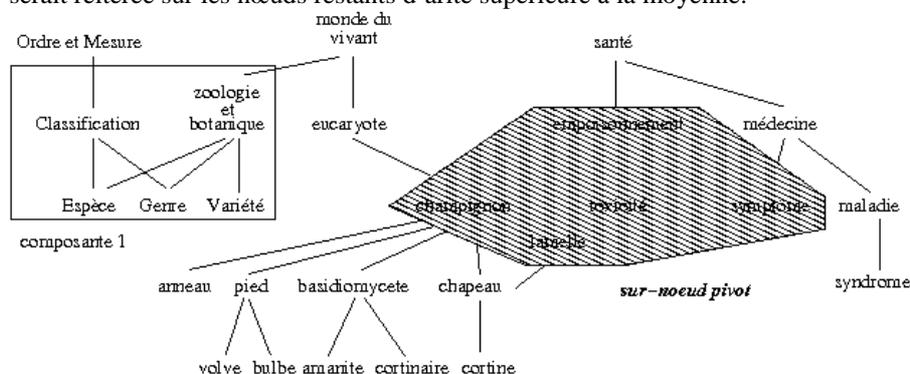


Figure 4. Détection des composantes fonctionnelles (ici une seule) liées à l'exploitation du corpus

Nous présentons maintenant les éléments participant au calcul des probabilités de pertinence associées aux associations sémantiques :

- le nombre de pas minimal qui relie leurs termes dans une composante ;
- la probabilité de proximité sémantique entre deux termes d'une association sémantique corrélée à la distance intra-composante ;
- une maximisation à 0,75 si les deux lemmes sont associés dans la même composante sémantique fonctionnelle (post-exploitation du corpus) ;

(*) 1- (distance/(distance max. +1)) / facteur de minoration	Distance intra-composante entre parenthèses (distance max.)	Probabilité de pertinence par rapport à l'ontologie (*)	Probabilité de pertinence : ontologie et requêtes
champignon – toxicité/symptôme	-	0,2	0,75
champignon – apparition	-	0,2	0,25
symptôme – toxicité	4 (5)	0,33	0,75
symptôme – apparition	-	0,2	0,2
amanite – reconnaissance/ident.	-	0,2	0,2
amanite – anneau	3 (4)	0,4	0,4
amanite – volve	4 (4)	0,2	0,2
cortinaire – cortine	-(cortine inconnu)	0,5	0,5
cortinaire – présence-fréquente	-	0,2	0,2
reconnaissance – identification	1 (2)	0,66	(synonymes)
anneau – volve	3 (4)	0,4	0,4
syndrome-phalloïdien – espèce	-	0,2	0,2
champignon – lamelle	2 (4)	0,6	0,75

Figure 5. Probabilités de pertinence des associations sémantiques

Cette probabilité peut être par la suite minorée car elle est vulnérable à l'insertion de vocabulaire de structure entre lemmes peu usité dans les documents (cette minoration dépend de l'homogénéité de l'ontologie). La probabilité d'association de deux lemmes inter-composantes est fixée à la plus basse des probabilités de pertinence calculée entre deux lemmes d'une même composante (ici 0,2 pour l'association entre *amanite* et *volve*), tandis que la probabilité d'association d'un lemme connu avec un lemme inconnu est fixée à 0,5.

### 5.5 Association « a priori » des lemmes par rapport aux associations

Nous calculons maintenant la probabilité de pertinence « a priori » des lemmes par rapport aux associations. Ce *prior* sera utilisé dans le théorème de Bayes pour finalement filtrer les associations sémantiques par les titres des sections dans lesquels elles se trouvent. La probabilité qu'un lemme présent dans un titre soit pertinent sachant que ce lemme est associé à une ou plusieurs associations sémantiques se trouvant dans la section qu'il introduit, est calculée comme suit.

Une probabilité de 0,5 correspondrait à ce que le terme ait le même poids dans le document que dans les requêtes. Pour étager ces probabilités entre 0,25 et 0,75, une première étape consiste à rechercher le lemme le plus sur-représenté dans les requêtes par rapport au document, et celui qui l'est le moins.

Dans notre exemple, le lemme le plus sur-représenté est *champignon*. Par rapport aux seules requêtes, sa sur/sous-représentation est de  $R = (7/31 / 4/37) = (0,226 / 0,108) = 2,1$ . Par contre le lemme *cortinaire* étant présent deux fois dans le document et une seule fois dans les requêtes, sa représentation est  $R = (1/31 / 2/37) = (0,032 / 0,054) = 0,6$ , puis à une probabilité de  $P = 0,25 + ((R-0,6)/3)$ , les probabilités étant étagées de 0,75 à 0,25.

Le tableau suivant présente les lemmes des titres, leur fréquence dans les requêtes et le document, et leur probabilité de pertinence par rapport aux requêtes :

lemmes présents dans les titres :	Fréquence du lemme dans les requêtes et dans le document	Probabilité de pertinence suivant les associations définies par les requêtes
Amanite < basid. < champignon	3/31 = 0,097 et 2/37 = 0,054	0,65
Champignon	7/31 = 0,226 et 4/37 = 0,108	0,75
Cortinaire < basid. < champignon	1/31 = 0,032 et 2/37 = 0,054	0,25
Lamelle	2/31 = 0,065 et 2/37 = 0,054	0,45
Syndrome et syndrome phalloïdien	2/31 = 0,065 et 4/37 = 0,027	0,25
Toxicité < empoisonnement	4/31 = 0,129 et 3/37 = 0,081	0,58

Figure 6. Fréquence et probabilité de pertinence des lemmes présents dans les titres

### 5.6 Calcul final de la probabilité de pertinence d'une association par rapport à une section ( $P(\text{association}|\text{lemme du titre})$ )

Pour contextualiser chaque association sémantique, nous appliquons le théorème de Bayes qui aura pour but de valider, via des probabilités conditionnelles, la pertinence de l'association sémantique par rapport aux lemmes des titres introduisant la section dans laquelle cette association a été retrouvée.

Contrairement à de nombreux travaux de fouilles de textes utilisant les réseaux bayésiens que nous allons rappeler rapidement, nous utilisons simplement le théorème de Bayes pour discriminer les associations sémantiques suivant la structure

du document. En 1992, le système de recherche d'information INQUERY [Callan et al., 1992] introduit l'utilisation des réseaux bayésiens dans la RI. Ceux-ci auront pour but de calculer la probabilité qu'une requête soit satisfaite par un document. Mais ce modèle utilisé dans INQUERY fait encore un traitement « plat » des documents, c'est-à-dire que leur structure n'est pas prise en compte et que tous les mots sont traités de la même manière quelque soit l'endroit où ils se trouvent. Une extension de INQUERY proposée par Myaeng [Myaeng et al., 1998] prend en compte la structure des documents en plus de leurs contenus, celle-ci étant représentée par un arbre. Chaque feuille de cet arbre est prolongée par les nœuds termes contenus dans l'élément de structure représenté par cette feuille. D'autres travaux ont appliqué les réseaux bayésiens à des corpus de documents XML [Zargayouna, 2004], mais ces modèles reposent sur des documents XML qui doivent être structurés uniformément.

La probabilité finale de pertinence d'une association par rapport aux lemmes de son titre est la suivante (lemme ayant la signifiant lemme du titre) :

$$P(\text{association}|\text{lemme}) = P(\text{lemme} | \text{association}) \cdot P(\text{association}) / P(\text{lemme})$$

$$(P(\text{lemme} | \text{association}) = P(\text{lemme} | \text{associations}) * \text{moyenne } P(\text{lemme}) \text{ des associations})$$

Voici le tableau récapitulatif des différents résultats des étapes de notre processus :

Association sémantique candidate	Lemme filtre du titre	P ( lemme   ass)	P (asso.)	P (lemme)	P (asso.   lemme)
<b>champ. – toxicité</b>	toxicité	0,58 * 0,422 = 0,245	0,75	0,393	<b>0,468</b>
champ. – apparition	toxicité	0,58 * 0,422 = 0,245	0,2	0,393	0,125
champ. – symptôme	toxicité	0,58 * 0,422 = 0,245	0,75	0,393	0,468
<b>champ. – toxicité</b>	champignon	0,75 * 0,422 = 0,317	0,75	0,75	<b>0,317</b>
champ. – apparition	champignon	0,75 * 0,422 = 0,317	0,2	0,75	0,085
champ. – symptôme	champignon	0,75 * 0,422 = 0,317	0,75	0,75	0,317
<b>symptôme – toxicité</b>	toxicité	0,58 * 0,422 = 0,245	0,75	0,393	<b>0,468</b>
symptôme – apparition	toxicité	0,58 * 0,422 = 0,245	0,2	0,393	0,125
symptôme – toxicité	champignon	0,75 * 0,422 = 0,317	0,75	0,75	0,317
symptôme – apparition	champignon	0,75 * 0,422 = 0,317	0,2	0,75	0,085
amanite – recon./identif.	amanite	0,65 * 0,422 = 0,274	0,2	0,75	0,073
amanite – anneau	amanite	0,65 * 0,422 = 0,274	0,4	0,75	0,147
amanite – volve	amanite	0,65 * 0,422 = 0,274	0,2	0,75	0,073
<b>cortinaire – cortine</b>	cortinaire	0,25 * 0,422 = 0,106	0,5	0,75	<b>0,071</b>
cortinaire - présence	cortinaire	0,25 * 0,422 = 0,106	0,2	0,75	0,028
reconnaiss. – identific.	amanite	0,65 * 0,422 = 0,274	- : synon.	0,75	-
<b>anneau – volve</b>	amanite	0,65 * 0,422 = 0,274	0,4	0,75	<b>0,147</b>
<b>syndrome-ph. – espèce</b>	amanite	0,25 * 0,422 = 0,106	0,2	0,75	<b>0,028</b>
<b>champignon – lamelle</b>	toxicité	0,58 * 0,422 = 0,245	0,75	0,393	<b>0,468</b>
<b>champignon – lamelle</b>	champignon	0,75 * 0,422 = 0,317	0,75	0,75	<b>0,317</b>

Figure 7. Récapitulatif des probabilités de pertinence des associations sémantiques

## 5.7 Bilan

Dressons un comparatif des calculs de probabilité effectués sur l'exemple (où les associations n'emploient que les lemmes des premières phrases de chaque paragraphe), et sur le corpus complet lié à la toxicité des champignons. Ce qui est indexé par les experts (lemmes ou associations) est en gras dans les tableaux :

Classement des lemmes par comptage de fréquence	Classement des lemmes par mesures <i>tf.idf</i> maximisées	Classement des associations par rapport à l'ontologie et aux requêtes	Classement des associations suivant le théorème de Bayes
<b>champignon</b> 3 <b>syndrome ph.</b> 3 <b>toxicité</b> 3 symptôme 2 <b>amanite</b> 2 <b>cortinaire</b> 2 syndrome or. 2 amanite ph. 1 amanita ve. 1 amanita vi. <b>anneau</b> 1 apparition 1 <b>volve</b> 1 <b>cortine</b> 1 espèce 1 genre 1 identification 1 <b>lamelle</b> 1 présence fréqu. 1 reconnaissance 1 syndrome pr. 1 variété 1	0.75 : <b>champignon</b> symptôme <b>amanite</b> <b>cortinaire</b> 0,435 : ----- <b>lamelle</b> 0,393 : <b>toxicité</b> 0,354 : ----- apparition syndrome prin. syndrome orol. reconnaissance identification <b>anneau</b> <b>volve</b> amanita ... présence <b>cortine</b> 0,25 : ----- <b>syndrome ph.</b> Espèce	0.75 : <b>champignon – toxicité</b> <b>champignon</b> – <b>symptôme</b> <b>symptôme – toxicité</b> <b>champignon – lamelle</b>  0.5?: <b>cortinaire–cortine</b>  0.4 : <b>anneau – volve</b> amanite – anneau  0.2 : amanite – volve champignon – apparition symptôme – apparition cortinaire – présence-fr. <b>syndrome-ph.–espèce</b>  0 : reconn. – identif.	0.468 (sachant toxicité) : <b>champignon – toxicité</b> <b>champignon – symptôme</b> <b>symptôme – toxicité</b> <b>champignon – lamelle</b> 0.317 (sachant champignon) : <i>mêmes que précédemment</i> 0,147 (sachant amanite) : <b>anneau – volve</b> amanite – anneau 0,125 (sachant toxicité) : champignon – apparition symptôme – apparition 0,085 (sachant champignon) : champignon – apparition symptôme - apparition 0,073 (sachant amanite) : amanite – reconn./identific./volve 0,071 : <b>cortinaire–cortine</b> (cortin.) 0,028 : cortinaire–présence (cortin.) <b>syndrome-ph.–espèce</b> (amanite) 0 : reconnaissance–identification

Figure 8. Classement des lemmes et des associations suivant les probabilités calculées

Voici les associations sémantiques faites par les experts pour chaque section :

La toxicité des champignons :

La toxicité des champignons ingérés est d'autant plus grave que l'apparition des symptômes est tardive. Ces symptômes sont le fait de deux syndromes principaux : les syndromes phalloïdien et ...

**champignon – toxicité**  
**champignon – symptôme** sachant toxicité  
**symptôme – toxicité** sachant champignon

Le syndrome phalloïdien et le syndrome orellanien de certains champignons à lamelles :

Les amanites :

La reconnaissance des amanites passe par l'identification d'un anneau et d'un volve. Les principales espèces responsables du syndrome phalloïdien sont amanita phalloïdes, ...

**champignon – lamelle** sachant toxicité  
**anneau – volve** sachant amanite  
**syndrome-phalloïdien – espèce** sachant amanite

Le syndrome phalloïdien et le syndrome orellanien de certains champignons à lamelles :

Les cortinaires

Les cortinaires sont caractérisés par la présence fréquente d'une cortine. C'est un genre qui comprend une très grande variété d'espèces dont certaines sont comestibles, d'autres d'une grande toxicité.

**champignon – lamelle** sachant toxicité  
**cortinaire – cortine**

A cette annotation comparons les taux d'adéquation de nos classements :

- par **comptage** : discrimination de 7 lemmes parmi les 9 pré-sélectionnés dont 5 sont parmi les 7 distingués par le classement :  $(5/7) * (7/9) = 55,5\%$ .
- par des **mesures *tf.idf* maximisées** : discrimination de 6 lemmes parmi les 9 pré-sélectionnés dont 5 sont parmi les 6 distingués :  $(5/6) * (7/9) = 65\%$ .
- par l'exploitation **de l'ontologie et des requêtes** : ce classement discrimine 7 associations sur 7 pré-sélectionnées dont 6 distinguées :  $(6/7) * (7/7) = 86\%$  qui est le meilleur taux brut d'adéquation.
- par le **théorème de Bayes** : discrimination de 6 associations sur 7 pré-sélectionnées dont 5 sont parmi les 6 discriminées :  $(5/6) * (6/7) = 72\%$ .

**Sur le corpus complet, les taux d'adéquation sont respectivement de 51% (comptages simples), 63% (mesures *tf.idf* maximisées), 83% (exploitation de l'ontologie et d'une indexation manuelle via les requêtes), et 70% (par Bayes).**

Nous remarquons que si le résultat donné par l'application du théorème de Bayes est donc moins bon que celui obtenu par l'établissement d'associations sémantiques prenant en compte seulement l'ontologie et l'exploitation du corpus, il permet d'enrichir l'annotation du document **en précisant le contexte de celle-ci** (grâce à la probabilité conditionnelle), et de définir des cartes sémantiques associées au corpus.

### **5.8 Evaluation « objective » des annotations générées**

Nous sommes en cours d'essais pour évaluer les co-occurrences que notre système produit, via une mesure que nous espérons objective. Notre objectif est de mettre en place un protocole d'évaluation via le calcul d'une mesure de similarité sémantique entre les termes des co-occurrences et en prenant comme référence le réseau sémantique pour la langue anglaise Wordnet [Wordnet] développé à l'université de Princeton. Wordnet est basé sur la notion de « synsets », soit d'ensembles de synonymes dénotant des concepts inter-connectés par différentes relations linguistiques. A l'instar de la quasi totalité des termes présents dans nos co-occurrences, nous n'utilisons que les substantifs présents dans ce réseau sémantique. L'utilisation de Wordnet nous oblige évidemment à mettre en place une correspondance entre le français et l'anglais (que nous ne détaillerons pas ici) pour nos corpus de langue française (nous envisageons également d'utiliser WOLF [WOLF], la version francisée, mais pour l'instant incomplète, de Wordnet).

Pour évaluer la pertinence des co-occurrences, nous avons commencé par utiliser les mesures de Resnik [Res 95], puis de Lin [Lin 98] légèrement différente de la précédente. Resnik a défini la notion de contenu informationnel (CI) en spécifiant que la similarité sémantique entre deux concepts ( $C_1$  et  $C_2$ ) est mesurée par la quantité de l'information qu'ils partagent. Le contenu informationnel est obtenu en calculant la fréquence du concept dans le corpus (par exemple Wordnet). La mesure de similarité est alors donnée par  $Sim(C_1, C_2) = \text{Max}[E(CS(C_1, C_2))] = \text{Max}[-\log_p(CS(C_1, C_2))]$  où  $CS(C_1, C_2)$  représente le concept le plus spécifique qui subsume les deux concepts dans l'ontologie. Lin [Lin 98] a formalisé une approche hybride qui combine plusieurs sources de connaissances différentes et qui représente la similarité en prenant en compte le chevauchement des concepts descendants de  $C_1$  et  $C_2$ . Nous nous sommes aussi intéressés au travail de Seco et consorts [Seco 2004] qui modifient le calcul du contenu informatif en postulant que plus un concept a de

descendants, moins il est informatif, et utilisent donc les hyponymes de ces concepts pour calculer leurs contenus informatifs.

Malheureusement toutes ces mesures sont relativement inadaptées à notre système pour la raison principale que nos co-occurrences relient des termes qui n'ont souvent pas de sens proche dans Wordnet hors de leur contexte, (ni bien sûr dans notre premier thésaurus dérivé de celui de LAROUSSE). Ainsi, dans le paragraphe suivant portant sur l'exploitation de mini-cartes sémantiques lors des sessions de navigation dans le corpus, nous verrons que les trois occurrences suivantes qui sous-tendent les axes sémantiques proposés à l'utilisateur, ne font sens que par rapport au lemme contexte :

- ainsi *anxiété – poussée* n'a de sens dans notre corpus que par rapport au lemme contexte *amanita* : poussée d'anxiété lors de la consommation d'amanita...
- *administration - sédatif / traitement* : administration d'un sédatif lors d'un traitement...
- *destruction – foie / phase* : destruction du foie lors de la première phase de l'intoxication...

Nous étudions donc actuellement l'intégration de ce troisième concept (le terme décrivant le contexte sémantique de la co-occurrence) dans notre mesure.

Par ailleurs notre système étant fortement adaptatif aux requêtes expertes qui modifient le thésaurus de référence, qui lui même modifie les mesures de similarité de contrôles (en vérifiant leur proximité dans le thésaurus), nous étudions également la modélisation de cette rétroaction pour mesurer la qualité des convergences opérées.

## **6. Génération de cartes sémantiques interactives**

Outre l'évaluation de la qualité de notre indexation, notre principal axe de travail est la création dynamique, lors des sessions de recherche par les utilisateurs, de cartes sémantiques interactives.

En effet, grâce à la spécification d'un contexte (grâce à notre processus fondé sur le théorème de Bayes) à chaque annotation d'une section du document par une co-occurrence, nous sommes en mesure de définir différentes vues sur le corpus. Ces vues sont calculées dynamiquement au cours des recherches effectuées par les utilisateurs (qui peuvent être ou non des experts, mais seules les recherches effectuées par les utilisateurs ayant un rôle d'expert sont susceptibles de modifier le processus d'indexation (voir paragraphe 5.4.2)). Ces vues interactives constituent des cartes sémantiques (réduites) constitutives de la cartographie sémantique du corpus.

### **6.1 La cartographie sémantique**

La cartographie sémantique est un outil de construction et de représentation graphique servant à mettre en exergue des réseaux d'informations ou de connaissances sur un domaine donné. Ce réseau précise les relations qui existent entre ses constituants, auxquels par exemple dans les cas d'exploitation de bases

documentaires, vont être associées des listes de documents et de mots clés (traditionnellement ordonnées par importance décroissante).

Toujours dans le contexte de la gestion électronique de documents, la cartographie sémantique passe par exemple par une classification automatique de ceux-ci selon la méthode des K-Means axiales [Lelu 98], résumant le corpus en ses composantes thématiques principales, l'homogénéité des composantes (des clusters générés) étant un bon critère de qualité. Relativement peu exhibée dans les travaux effectués dans le domaine de la cartographie sémantique (nous citerons bien sûr les travaux dérivés des cartes auto-adaptatives de Kohonen [Kohonen 1999, 2001]), et ignorée dans la majorité des moteurs de recherche hormis le célèbre KartOO [KARTOO], nous développons une cartographie **interactive** lors de la navigation dans le corpus.

## 6.2 Navigation dans le corpus à l'aide de la cartographie sémantique

Reprenons les annotations de plus fortes probabilités effectuées par notre système sur le corpus d'où est tiré notre exemple d'illustration, en les récapitulant dans un tableau spécifiant les contextes sous-jacents. Ce tableau se lit ainsi :

les lemmes désignés en ligne et en colonne sont associés dans une **co-occurrence L1-L2 (ou L2-L1) sachant les contextes C1, [C2...] classés suivant un ordre de probabilité de pertinence**, (les contextes soulignés correspondant à un lemme directement présent dans la co-occurrence).

	champ.	toxicité	sympt.	lamelle	anneau	volve	amani.	appari.	syn-ph.	espèce
champ.		<u>toxicité champ.</u>	<u>toxicité champ.</u>	<u>toxicité champ.</u>				<u>toxicité champ.</u>		
toxicité	<u>toxicité champ.</u>		toxicité champ.							
sympt.	toxicité champ.	<u>toxicité champ.</u>						toxicité champ.		
lamelle	champ. toxicité									
anneau						amani. toxicité	<u>amani.</u>			
volve					amani. toxicité		<u>amani.</u>			
amani.					<u>amani.</u>	<u>amani.</u>				
appari.	toxicité champ.		toxicité champ.							
syn.-ph.										amani. syn-ph.
espèce									amani. syn-ph	
...										

(champ. = champignon, sympt. = symptôme, amani. = amanite, syn.-ph. = syndrome phalloïd.)

**Figure 9.** Tableau récapitulatif de quelques contextes de co-occurrences

Suivons maintenant une session de recherche effectuée par un utilisateur du corpus, (cette session fera apparaître d'autres lemmes que ceux précédemment vus).

La première requête effectuée est composée des termes suivants :

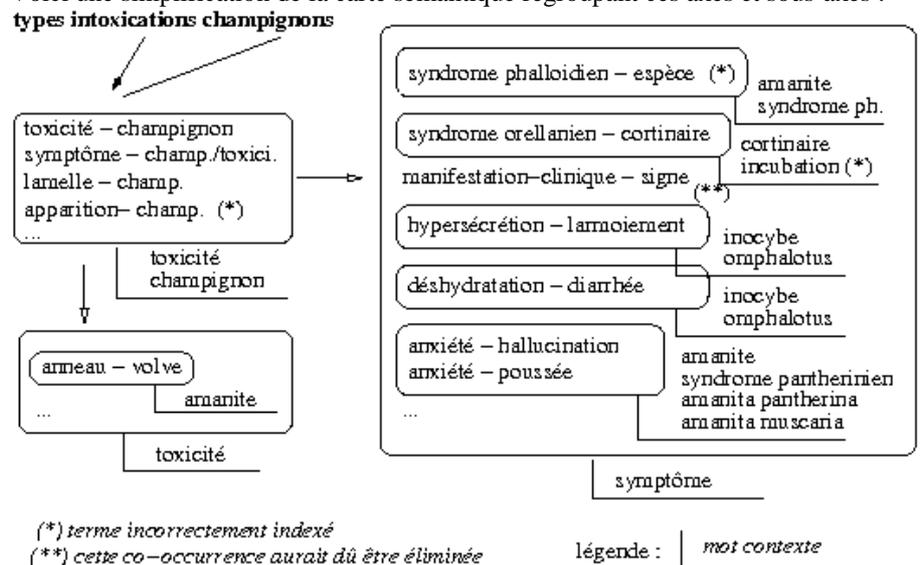
**types intoxications champignons**

Après lemmatisation de ceux-ci (et rapprochement avec les termes qui leurs sont fortement corrélés dans le thésaurus), le système sélectionne les co-occurrences contextualisées par ces termes dans la base d'indexation.

Ainsi, si le lemme *type* n'est pas contexte dans le corpus d'illustration, *toxicité* (corrélé à *intoxication par le thésaurus*), et *champignon*, le sont pour de nombreuses co-occurrences (125 fois pour *toxicité* et 324 fois pour *champignon*). Face à cette profusion, une restriction est alors opérée : seules les co-occurrences dans lesquelles les lemmes *toxicité* ou *champignon* apparaissent directement sont filtrées. En effet, un lemme qui apparaît simultanément comme terme de la co-occurrence et contexte de celle-ci, est non seulement présent dans un ou plusieurs des titres de la section mais apparaît également un nombre significatif de fois dans le corps de celle-ci (et inversement moins dans les autres sections du document). Il a donc une forte probabilité de correspondre à un **axe thématique** de navigation dans le corpus.

A ces principaux axes de recherche, des **sous-axes** sont également proposés à l'utilisateur : ils correspondent aux autres contextes des co-occurrences filtrées par les contextes précédents (et correspondent donc aux axes thématiques secondaires).

Voici une simplification de la carte sémantique regroupant ces axes et sous-axes :



**Figure 10.** Carte sémantique liée à la requête « Type intoxication champignons »

Le tableau suivant reprend la présentation (simplifiée) des axes thématiques sur le thème : **toxicité champignon** :

Axe thématique	Sous-axe 1	Sous-axe 2	Sous-axe 3	Sous-axe 4
toxicité	amanite			
symptôme	syndrome phalloïdien	cortinaire incubation	inocybe omphalotus	amanita pantherina amanita muscaria

Si l'utilisateur sélectionne un axe et un de ses sous-axes, il pourra ainsi accéder aux paragraphes des différents documents du corpus qui recèlent des co-occurrences ayant été indexées par les contextes associés (l'interface de l'application est en cours de développement) :

Axe thématique	Sous-axe thématique	Fragment de la hiérarchie de titres d'une section d'un document [numéro] et <b>mots contextes</b>	Fragment du paragraphe du document indexé et <b>co-occurrences</b> ayant été contextualisées
toxicité	amanite	[1] - <b>Toxicité</b> des champignons -- Le syndrome phalloïdien et le syndrome orellanien ... --- Les <b>amanites</b>	La reconnaissance des amanites passe principalement par l'identification d'un <b>anneau</b> et d'une <b>volve</b> ...
symptôme	syndrome ph.	[1] - Toxicité des champignons -- Le <b>syndrome phalloïdien</b> et le syndrome orellanien ... --- Les amanites [23] - Tableaux cliniques -- <b>Symptômes</b> associés au syndrome phalloïdien	Les principales <b>espèces</b> responsables du <b>syndrome phalloïdien</b> sont amanita phalloïdes, ... On estime que l'ingestion d'une seule moitié de chapeau d'une des trois principales <b>espèces</b> d'amanites associées au <b>syndrome phalloïdien</b> . ...
symptôme	cortinaire incubation	[23] - <b>Cortinaires</b> -- <b>Symptômes</b> associés au syndrome orellanien --- Durée d' <b>incubation</b>	Ce cortinaire (cortinarius orellanus) de petite taille reste peu courant, mais le <b>syndrome orellanien</b> dont ses consommateurs sont victimes est redoutable.
symptôme	inocybe omphalotus	[14] - Syndrome sudorien ( <b>inocybe, omphalotus</b> ) -- <b>Symptômes</b> cliniques	Signes d'imprégnation cholinergique ( <b>hypersécrétions, larmolement, fourmillements</b> dans les membres) et digestifs ( <b>déshydratation, diarrhées</b> )



Et voici maintenant les paragraphes des documents accessibles via la sélection des axes et sous-axes thématiques proposés par la carte sémantique :

Axe thématique	Premier sous-axe thématique	Fragment de la hiérarchie de titres de la section et <b>mots contextes</b>	Fragment du paragraphe associé et <b>co-occurrences</b>
syndrome phalloïdien	phase	[27] - Tableaux cliniques -- Symptômes associés au <b>syndrome phalloïdien</b> --- <b>Phases</b>	Intoxication en deux phases : - troubles digestifs - <b>destruction du foie</b>
syndrome panthérinien	toxicomanie	[14] - <b>Syndrome panthérinien</b> (amanita pantherina, amanita muscaria) -- Symptômes cliniques --- <b>Toxicomanies</b>	Ces champignons utilisés par les <b>toxicomanes</b> pour leurs <b>propriétés hallucinogènes</b> peuvent être responsables d'intoxications graves...
amanita pantherina amanita muscaria	traitement	[14] - Syndrome panthérinien ( <b>amanita pantherina, amanita muscaria</b> ) -- Symptômes cliniques -- Principes du <b>traitement</b>	<b>Administration de sédatifs</b> et surveillance de la victime...

Nous sommes donc en mesure grâce aux contextes sémantiques basés sur les calculs de probabilités conditionnelles liées au théorème de Bayes, de dériver des cartes sémantiques réduites, mais focalisées sur les thématiques d'une session de recherche de documents, et dans ces documents, de fragments de documents.

Pour l'instant, notre système a un rappel qui satisfait les experts ayant testé le prototype, mais présente un bruit encore trop important pour que les cartes sémantiques puissent être exhibées de manière prégnante. L'augmentation de la précision du système devrait passer par une amélioration des filtres syntaxiques (le bruit vient notamment du fait que trop de co-occurrences sont sélectionnées, les motifs syntaxiques sous-tendant le choix des mots à associer étant trop laxistes).

En termes de temps de calcul et volumétrie, si le moteur de recherche proposé par le logiciel de gestion documentaire de la société C6 (Documentum) est performant, il ne réalise pas une indexation contextuelle contrairement à notre système, qui génère des associations sémantiques jugées pertinentes avec peu de faux positifs par les experts. Par ailleurs, notre système devrait d'une part proposer à terme l'exploitation de cartes sémantiques dynamiques et interactives d'exploration des corpus, et d'autre part a déjà été mis en oeuvre sur plusieurs corpus et donne des résultats de qualité similaire (à condition d'être étayé par un minimum de connaissances recélées dans les ontologies-support appropriées ce qui est quasiment toujours le cas dans le domaine biomédical).

## 7. Conclusion

Le but de notre travail est de pouvoir offrir à un système documentaire cohérent, une indexation rapide (et donc purement incrémentale), générant une volumétrie d'annotation sémantique raisonnable. Nous avons donc choisi de mettre en oeuvre la

règle de Bayes pour filtrer les associations sémantiques les plus pertinentes d'une section du document, (une section étant un sous-arbre quelconque de paragraphes), par rapport aux termes du titre l'introduisant. Cette heuristique efficace en termes de calcul, donne de bons résultats qualitatifs, reconnus par les utilisateurs des corpus, et apporte deux avantages liés à la contextualisation des annotations : la création de **cartes sémantiques interactives** qui vont faciliter l'exploration du corpus, et la **révision des ontologies sous-jacentes à l'indexation**, lors des sessions de recherche effectuées par les utilisateurs-experts des corpus. Cela dit, deux problèmes peuvent être soulevés.

Une première critique tout à fait naturelle de ce travail, est l'effet délétère que pourrait causer des mots de structure par nature vides (comme « introduction », « conclusion »). En fait ce problème est rarement fondé car si ces termes ne se retrouvent que dans les titres et non dans les sections qu'ils introduisent, ils n'apparaîtront pas dans les associations sémantiques retenues pour indexer les paragraphes des documents. Au contraire un mot-concept apparaissant dans un titre va être également sur-employé dans les sections où il apparaît, et donc discriminé.

La seconde critique plus sérieuse, est l'impact de certains biais. Le premier biais est dû au fait que l'analyse des requêtes composées par les utilisateurs et qui modifie la probabilité de pertinence d'une association sémantique par rapport aux lemmes qu'elle contient, n'est faite qu'à partir du moment où un nombre critique de requêtes se trouve dans l'entrepôt de données. Dans l'exemple, cette analyse n'a été déclenchée qu'à partir du moment où tous les lemmes contenus dans les associations sémantiques filtrées par la première étape du processus ont été retrouvés dans les requêtes des utilisateurs. Cette analyse étant amorcée trop tôt, les derniers lemmes intégrés sont défavorisés. Son seuil de déclenchement doit donc être affiné.

Le second biais, corrélé au premier, est que l'exploitation d'un corpus n'est ni homogène dans le temps, ni dans le vocabulaire utilisé, certains descripteurs étant plus utilisés que d'autres, par habitude ou par notoriété. Ainsi, peu d'utilisateurs ont interrogé notre corpus avec le lemme « cortinaire », le nom de champignon le plus fréquemment utilisé comme point d'entrée dans la recherche d'information sur leur toxicité, étant « amanite », *même si les experts étaient également à la recherche d'information sur la toxicité des cortinaires*. Cette sous-représentation conduit à la dépréciation des associations fondées sur ce lemme, et doit être corrigée.

## 8. Bibliographie

- Besançon R., Rajman M., « Filtrages syntaxiques de co-occurrences pour la représentation vectorielle de documents », TALN 2002
- Bratko A., Filipic B., « Exploiting Structural Information in Semi-structured document classification », Andrej P., 13th Int. Electrotechnical and Computer Science Conf., ERK'2004
- Callan J. P., Croft W. B., and Harding S. M., « The INQUERY Retrieval System », In Tjoa A. M. and Ramos I. editors, *Database and Expert Systems Applications, Proceedings of the International Conference*, pages 78–83, Valencia, Spain. Springer-Verlag, 1992
- Denoyer L., « Apprentissage et inférence statistique dans les bases de documents structurés : application aux corpus de documents textuels », thèse de doctorat de l'Un. Paris VI, 12/2004

- Denoyer L., Gallinari P., « Bayesian network model for semi-structured document classification », *Inf. Processing Management* 40(5):807-827, 2004
- Hascoet M., Beaudouin-Lafon M., « Hypermedia exploration with interactive dynamic maps », *International Journal on Human Computer Studies*, 43:441-464, 1995
- Kohonen T., « Self-Organizing maps », Springer Series in Information Sciences, extended edition, 2001
- Kohonen T., Kaski S., Lagus K., Salojärvi J., Honkela J., Paatero V., Saarela A.. « Self organization of a massive text document collection », in: E. Oja, S. Kaski (Eds.), *Kohonen Maps*, Elsevier, 171-182., Amsterdam, 1999
- Lelu A., Halleb M., Delprat B., « Recherche d'information et cartographie dans des corpus textuels à partir des fréquences de n-grammes », *Actes des 4emes Journées Internationales d'Analyse Statistique des données Textuelles*, Nice, 1998
- Lin. D., « An Information-Theoretic Definition of similarity ». In *Proceedings of the Fifteenth International Conference on Machine Learning (ICML'98)*. Morgan-Kaufmann: Madison, WI, 1998.
- Lin X., Soergel D., Marchionini G. « A self-organizing Map for Information Retrieval », *Proceedings of the 14<sup>th</sup> annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 262-269, Chicago, 1991
- Myaeng S., Jang D., Kim M., and Zhoo Z., « A Flexible Model for Retrieval of SGML documents », *Proc 21st ACM SIGIR*, 138-140, Melbourne, ACM Press, New York, 1998
- Piwowski, Gallinari P., « A bayesian network model for page retrieval model in a hierarchical structured collection », *XML w. of the 15th ACLM SIGIR Conf.*, Tampere, Finland 2002
- Resnik Ph, «Using Information Content to Evaluate Semantic Similarity in a Taxonomy», *IJCAI 95*
- Salton G., Buckley C., « Term-weighting Approaches in Automatic Text Retrieval », *Information Processing and Management* , 24(5), pp. 513-523, 1988.
- Seco N., Veale T., Hayes J., « An Intrinsic Information Content Metric for Semantic Similarity in WordNet », *Proceedings of ECAI'2004*, Valence, Espagne, 2004
- Wilkinson R., «Effective retrieval of structured documents», *Proc. 17<sup>th</sup> ACM SIGIR*, Dubin, 1994
- Zargayouna H., «Contexte et sémantique pour une indexation de doc. semi-structurés» *CORIA'04*

Webographie :

Documentum : [www.emc.com](http://www.emc.com)

KartOO : <http://www.kartoo.com>

LUCENE : <http://fr.wikipedia.org/wiki/Lucene>

TreeTagger : <http://ims.uni-stuttgart.de/projekte/corplex/TreeTagger>

Wordnet : <http://wordnet.princeton.edu>

WOLF : <http://alpage.inria.fr/~sagot/wolf.html>