
Un Système d'Aide à la Recherche d'Information en ligne basé sur les Ontologies (SA-RI-Onto)

Rania Soussi* — Nesrine Ben Mustapha *— Hajer baazaoui zghal*
— Marie-aude Aufaure**, ***

* Laboratoire RIADI ENSI Campus Universitaire de la Manouba 2010

{nesrine.benmustapha, , hajer.baazaouizghal }@riadi.rnu.tn

soussi_rania@yahoo.fr

** SUPELEC Plateau du Moulon 3, rue Joliot Curie 91 192 Gif sur Yvette Cedex

Marie-Aude.Aufaure@supelec.fr

*** INRIA Paris-Rocquencourt Domaine de Voluceau 78 153 Le Chesnay Cedex

Marie-Aude.Aufaure@inria.fr

RÉSUMÉ. La croissance très importante des informations disponibles sur Internet nécessite des outils de recherche de plus en plus performants permettant de discerner efficacement les informations intéressantes parmi des centaines voire des milliers de documents. Seulement, la qualité des résultats fournis par les moteurs de recherche traditionnels n'est pas toujours pertinente surtout quand il s'agit de composer plus d'une requête. Ceci est dû aux ambiguïtés linguistiques et aux concepts abstraits qui ne sont pas bien traités. L'utilisation de la sémantique et plus précisément des ontologies présente des atouts importants. L'objectif de cet article est de montrer l'apport des ontologies dans la recherche d'information en ligne. Ainsi, un système d'aide à la recherche d'information en ligne basé sur les ontologies est proposé. Ce système est composé de deux ontologies : une ontologie de domaine et une ontologie de services ainsi que WordNet, pour représenter les concepts ainsi que les services de domaine. La contribution de ces ontologies pour améliorer la recherche d'information en ligne est montrée par la proposition et l'expérimentation d'un système de recherche d'information en ligne basé sur les ontologies.

MOTS-CLÉS : recherche d'information en ligne, construction d'ontologies

ABSTRACT. The huge number of available documents on the Web makes finding relevant ones challenging. Thus, searching for information becomes more and more complex because of the growing volume of data and of its lack of structure. The quality of results that traditional full-text search engines provide is still not optimal for many types of user queries. The ambiguities of natural languages and abstract concepts are handled inadequately by full-text search engines. Ontologies provide a solution to these problems. An architecture composed of several ontologies is proposed to represent concepts as well as services of domain. In this paper, we expose the contribution of these ontologies in information retrieval and we propose a new retrieval system based on ontologies. Experimentation in the domain of the tourism is presented, and the gotten results are compared to other systems.

KEYWORDS: on-line Information retrieval, ontologies construction

2 Revue. Volume X – n° x/année

1. Introduction

Depuis l'avènement de l'Internet dans les années 90, le nombre croissant de documents disponibles rend indispensable l'utilisation d'outils adaptés permettant d'assister et d'aider les utilisateurs au cours de leurs recherches d'information en ligne. Par ailleurs, les techniques généralement employées par les moteurs de recherche reposent sur des méthodes statistiques et des traitements syntaxiques ne permettent pas de traiter la sémantique contenue dans la requête de l'utilisateur ainsi que dans les documents. D'autres problèmes, tels que celui lié au vocabulaire et à l'hétérogénéité des données, le choix des mots clés et celui du filtrage peuvent être rencontrés. Ces problèmes reposent sur la capacité des mots des langages humains à générer la polysémie et la synonymie. Aider un utilisateur à trouver l'information qu'il cherche dans ce contexte devient donc une tâche de plus en plus difficile. D'où le recours à la prise en compte du niveau sémantique qui permet de pallier ces difficultés de la recherche d'information et l'apparition du Web sémantique (Berners-Lee et al., 2001). Ce dernier évoque la notion d'ontologie dans la nouvelle architecture de base et la considère comme étant une de ses couches fondamentales. En effet, les ontologies qui sont une spécification explicite, formelle d'une conceptualisation partagée (Gruber, 1993) cherchent à donner un sens aux données du Web ce qui facilite de manière considérable l'accès à des documents pertinents. Des approches ont été proposées pour permettre l'extraction de la sémantique et mieux répondre aux requêtes émises. Par contre, la plupart de ces techniques ont été conçues pour s'appliquer au Web en entier et non pas à un domaine particulier. Une piste intéressante consiste à utiliser une ontologie non seulement pour représenter un domaine spécifique mais aussi pour le filtrage et la classification des documents retournés à l'utilisateur par sous thèmes. L'architecture ontologique pour la recherche d'information utilisée comme base comprend trois ontologies : une ontologie de domaine, une ontologie de structure des sites Web et une ontologie des services de domaine (Baazaoui et al., 2007).

Dans cet article, nous présentons notre système de recherche d'information en ligne en utilisant les relations entre l'ontologie de domaine et l'ontologie de services issues du processus de recherche d'information et du processus de classification des documents. L'expérimentation menée dans le domaine du tourisme est présentée à la fin de cet article ainsi que la comparaison des résultats obtenus à ceux d'autres systèmes.

2. Ontologie et recherche d'information

Les systèmes de recherche d'information devraient fournir à l'utilisateur un accès facile à l'information à laquelle il s'intéresse. Cependant, les systèmes traditionnels ont des difficultés à fournir un résultat pertinent. Ainsi, certains des travaux actuels dans la RI tentent d'améliorer le procédé de récupération avec l'aide des ontologies. Les expériences effectuées en utilisant WordNet (Miller, 1995) avec une stratégie

intelligente de recherche ont prouvé qu'un gain significatif est possible sur le TREC de l'ensemble des données (Baziz., 2004) (Voorhees, 1994). Elles peuvent aider l'utilisateur à détecter ses besoins et trouver les mots-clés appropriés qui utilisent les concepts existants dans l'ontologie et leur description. La connaissance qui représente les ontologies peut être utilisée à différents niveaux dans le processus de RI. Elle peut aider à l'indexation des documents, alors appelée indexation sémantique. Les ontologies peuvent également aider à la formulation du besoin de l'utilisateur et à l'accès aux documents. Enfin l'ontologie peut être utilisée dans le modèle lui-même pour réaliser l'appariement entre le besoin et les documents. Afin de profiter des apports des ontologies dans la recherche d'information, nous proposons un système d'aide à la recherche d'information en ligne basé sur les ontologies. Ce système repose sur deux ontologies : ontologie de domaine et ontologie de services ainsi que WordNet, pour représenter les concepts ainsi que les services de domaine. Ainsi, un ensemble de services disponibles à une ontologie du domaine spécifique (tourisme, médecine, héritage culturel, etc.) sont associés. L'ontologie des services est en rapport avec les tâches, telles que, dans le domaine du tourisme, réservation d'hôtel, réservation d'une voiture, etc. Dans ce qui suit, la construction des ontologies, l'architecture de SARIOnto et l'apport de ces ontologies sont détaillés.

3.1. Construction des ontologies

L'ontologie de domaine construite semi-automatiquement en utilisant notre technique d'apprentissage d'ontologies développée dans OntoCoSemWeb (Baazaoui et al., 2007). Cette ontologie est la base de notre système de recherche d'information en ligne, et contient la plupart des concepts du domaine et leurs propriétés. Les tâches reliées au domaine, aussi appelées services sont modulées. En effet, selon le grand dictionnaire terminologique¹, une activité est un ensemble des tâches élémentaires ou des travaux exécutés par un individu ou un groupe et qui conduisent à la réalisation de biens ou de services. Chaque domaine est caractérisé par une liste de services, activités et tâches. Ces services sont liés aux concepts existants dans l'ontologie du domaine. Un concept contient un ensemble de propriétés. Par exemple dans notre ontologie de domaine, le concept « hotel » a les propriétés « Name », « Star Number » et « Address ». Chaque concept peut avoir des sous_concepts, comme « room » et « suite » pour le concept « hotel ». Chaque concept de l'ontologie du domaine est le sujet d'un ou plusieurs services, activités ou tâches. Par exemple, le service «Lodging_hotel » est associé au concept « hotel ». Cette relation donne la possibilité d'améliorer la recherche et aider les utilisateurs à mieux exprimer leurs besoins (figure 1). La construction de cette ontologie des services, ainsi que la mise en correspondance avec l'ontologie de domaine, sont pour le moment réalisées manuellement.

¹ <http://www.granddictionnaire.com/>

4 Revue. Volume X – n° x/année

3.2. Architecture de SARIOnto

Le système est constitué de trois modules principaux : un module de traitement de la requête, un module pour la recherche et le traitement des documents et le dernier pour leur classification. Le module de traitement de la requête permet d'enrichir la requête initiale de l'utilisateur en se basant sur les concepts et relations de l'ontologie de domaine et sur WordNet. La requête enrichie est fournie à un moteur de recherche. Le résultat obtenu est ensuite traité par le module de traitement des documents. Ces documents seront classifiés par services en utilisant le module de classification par service. Ce module guide l'utilisateur à construire une deuxième requête en utilisant les services choisis. Ce système est basé sur :

- L'usage d'une ontologie du domaine combiné avec une ontologie des services ;
- La classification par service des résultats d'une requête utilisée pour améliorer une recherche basée sur les services correspondants; ce qui permet d'assurer un usage plus facile de l'information;
- L'adaptation du modèle vectoriel (Salton, 1983) dans lequel nous substituons les termes par les concepts.

Dans la suite, nous décrivons le fonctionnement de SARIOnto.

3.3. Apport des ontologies dans SARIOnto

Dans cette partie nous détaillons le fonctionnement des modules de SARIOnto en expliquant le rôle des ontologies dans chacun d'eux.

3.3.1. Enrichissement de la requête

Afin d'obtenir une requête plus pertinente qui met en relief le domaine de recherche, la requête émise par l'utilisateur est enrichie sur la base des concepts et des relations de l'ontologie du domaine et les relations sémantiques de WordNet pour trouver les synonymes, hypéronymes, hyponymes des mots clés de la requête. SARIOnto, permet à l'utilisateur de s'appuyer sur l'ontologie de domaine pour choisir les concepts qu'il veut ajouter à sa requête. La requête est traitée selon les étapes suivantes :

Analyse morphologique : Elle permet la reconnaissance des différentes formes de mots à partir d'un lexique (dictionnaire, thésaurus). La lemmatisation permet la transformation d'une forme conjuguée ou fléchie à sa forme canonique ou lemme. Par exemple, des règles de dérivation permettant de retrouver la forme de base du dérivé (exemple : constitutionnel → constitution). Dans notre approche la lemmatisation de la requête est assurée par TreeTagger (Schmid, 1994).

Analyse sémantique : Après la lemmatisation de la requête, les lemmes des termes sont filtrés. Les mots vides sont éliminés, nous ne traitons que les verbes (verbe non vide : différent de find, do, get...) et les noms pour détecter par la suite les concepts correspondants existant dans l'ontologie de domaine. Pour poursuivre

l'analyse sémantique, nous utilisons WordNet, qui est organisé en un ensemble de synonymes et qui fournit les différents sens associés aux mots de la langue anglaise (par exemple, le mot *room* est associé à 5 sens). Il s'est donc avéré intéressant de permettre aux utilisateurs de sélectionner les sens des mots clés. Puis, l'algorithme suivant est appliqué aux termes lemmatisés :

- Si le terme est un verbe, nous conservons sa forme lemmatisée. Par exemple : le verbe « *traveled* » a pour lemme « *travel* », c'est ce lemme qui est sauvegardé.

- Si le terme est un nom : ses synonymes, hyponymes et hypéronymes sont récupérés à partir de WordNet, en utilisant le sens choisi par l'utilisateur. Deux sous cas se présentent, (1) si le terme ou l'un de ses composants existe dans l'ontologie alors le terme sera remplacé par ce concept, ses sous concepts et les propriétés qui lui sont relatifs, (2) si ni le terme ni ses composants ne sont présents dans l'ontologie alors ce terme est éliminé de la requête.

Ainsi, nous obtenons une nouvelle requête construite avec les concepts et les relations extraites à partir de l'ontologie de domaine. Par exemple, la requête initiale envoyée par l'utilisateur est: "find hotel ". Cette requête est lemmatisée avec TreeTagger. La forme lemmatisée rendue est: find VV find, hotel NN hotel(terme, type grammatical, et lemme). Le mot "hotel" a un sens unique dans Wordnet qui est: «un bâtiment où les voyageurs peuvent payer pour se loger et repas et autres services», et "find" est un verbe vide alors il est enlevé de la requête. L'analyse sémantique retourne une requête enrichie qui est "Hotel Suite Room Address Name Price Star".

3.3.2. *Processus de recherche*

Nous utilisons l'ontologie pour généraliser des mots à des concepts comme il sera expliqué dans les étapes ultérieures. Après son enrichissement, la requête est soumise à un moteur de recherche (google). Google retourne un résultat trié suivant le score de l'algorithme Page Rank qui mesure la popularité d'une page Web. Nous obtenons une première liste de documents qui sont rafraîchis dynamiquement par Google. Les documents résultants sont traités comme suit:

Analyse sémantique : Chaque document, se trouvant dans la liste (qui n'est pas dans le format pdf ou doc), est téléchargé à partir de son URL sous la forme d'un document HTML puis analysé avec DOM (Stenback et al., 2003). En récupérant le texte et en l'analysant morphologiquement avec TreeTagger, nous pouvons extraire les formes de base des mots pour pallier aux problèmes des variations morphologiques et détecter les concepts qui existent dans l'ontologie de domaine et dans la requête.

Filtrage avec le modèle vectoriel : Pour récupérer les documents les plus pertinents, nous avons utilisé le modèle vectoriel de Salton (Salton, 1983). Mais, pour que ce modèle soit plus adapté à notre approche, nous substituons les termes par les concepts. Plus précisément, dans la formule de similarité nous calculons les poids des concepts au lieu de celui des termes. Un concept est représenté par

6 Revue. Volume X – n° x/année

l'ensemble de ses synonymes, hypéronymes et hyponymes. Dans ce modèle chaque document est représenté par un vecteur : Soit t_i , un concept de la requête Q , $D_j = (d_{1j}, d_{2j}, d_{3j}, \dots, d_{Nj})$, avec d_{ij} : poids du concept t_i dans le document D_j , et N est le nombre de concepts qui forment la requête. Et chaque requête est représentée par un vecteur : $Q = (q_1, q_2, q_3, \dots, q_N)$, q_i : poids du concept t_i dans la requête Q .

Dans notre approche, les poids des termes dans la requête valent toujours initialement 1 (la requête est reformulée avec les concepts et relations de l'ontologie, et sans répétition du même terme). Plus explicitement, si notre requête reformulée est par exemple « hotel room », alors $Q = (1,1)$. Pour les documents, nous utilisons la formule de poids TF*IDF normalisée (Buckley et al., 1996) pour donner une chance égale à tous les documents et ne pas favoriser les plus longs. La mesure de similarité entre la requête et les documents sera calculée avec la formule du cosinus (Salton, 1983). Après avoir calculé la similarité pour tous les documents, nous obtenons une deuxième liste triée par ordre croissant suivant leurs similarités par rapport à la requête. Seuls les documents qui ont une similarité non nulle sont affichés à l'utilisateur.

3.3.3. Classification des documents

La classification permet d'affecter un document à une des classes ou sous-classes d'un plan de classement ou table de classification, c'est-à-dire dans un des domaines de la connaissance. Un résultat classifié par catégorie permet de faciliter la récupération de l'information et même de détecter d'autres besoins. Lorsque le résultat est traité et affiché à l'utilisateur, il peut faire l'objet d'une classification par service fourni par le module de classification. La requête reformulée par le module de traitement de la requête contient un ensemble de concepts qui appartiennent à l'ontologie du domaine. Ces derniers sont en liaison avec un ensemble de services de l'ontologie de services, plus précisément chaque concept est l'objet d'un ou plusieurs services. Par exemple, au concept « restaurant » est associé au service « Restauration ». Cette liaison est construite manuellement au cours de la construction de l'ontologie de services. En utilisant cette liaison, nous pouvons extraire, à partir des concepts ajoutés aux requêtes, les services qui lui sont relatifs dans l'ontologie des services. Le modèle vectoriel est de nouveau utilisé, un service est représenté par un vecteur : $Servi = (c_1, c_2, \dots, c_N)$ avec N est le nombre de concepts relatifs à un service. Pour chaque document retenu, nous calculons sa similarité avec les services en utilisant la formule du cosinus, puis il est affecté au service avec lequel il est le plus similaire. Ainsi les URL sont affichées par service. Dans notre exemple précédant, la requête initiale était "find hotel". Les services, activités et tâches détectées à partir des concepts dans l'ontologie du service sont:

Le service: Lodging_Hotel, Les activités: Reservation, Hotel_Search,
 Modify_Reservation, Disponibilité_Verification, Les tâches:
 Verify_RoomNumber, Verify_RoomType, Search_HotelName,
 Search_HotelAdress.

Chaque document ayant une valeur de similarité supérieure à 0, est attribué à un service, une activité et une tâche.

3.3.4. *Reformulation de la requête*

La classification effectuée pour le résultat permet à l'utilisateur de cerner les catégories des tâches, activités et services qui concernent sa requête. L'utilisateur peut alors choisir les services correspondants à ce qu'il recherche. Le module de traitement de la requête capture les services choisis et utilise la liaison entre l'ontologie de services et celle de domaine pour formuler une deuxième requête avec les nouveaux concepts et les relations détectés. Par exemple, en choisissant le service « Restauration », nous trouverons forcément le concept « restaurant » dans la nouvelle requête. La nouvelle requête est envoyée au moteur de recherche et le système répète le processus déjà décrit pour afficher un nouveau résultat plus raffiné. Dans notre exemple, si l'utilisateur choisit la tâche "Search_HotelName", la nouvelle requête sera "Hotel Name" parce que le concept "Hotel" et la propriété "Has_Name" sont liés à cette tâche. Une expérimentation de l'ensemble des propositions a été menée et est décrite dans la section suivante.

4. Evaluation de SARIOnto

SARIOnto permet de fournir à l'utilisateur un service en ligne. Il utilise l'Api Jena pour manipuler les ontologies et Google API pour effectuer les recherches à travers le web. Dans le but d'évaluer le système proposé, nous utilisons les mesures de précisions et rappel classique. En revanche, évaluer des systèmes de recherche basés sur la sémantique est une tâche complexe vu que les mesures classiques de précision/rappel s'avèrent difficile à être calculée de manière automatique. On a choisi un protocole d'évaluation centré utilisateurs. En effet, nous avons effectué un test basé sur 15 requêtes dans le domaine du tourisme et évalué par 10 utilisateurs. Puis, nous avons comparé ces mesures avec les mesures obtenues en utilisant deux autres systèmes; à savoir le premier, Lucene est un système de recherche d'information traditionnel basé sur la recherche par les mots-clés. Le deuxième est basé sur la reformulation de la requête, mais sans faire appel au modèle vectoriel. D'après l'évaluation faite, notre système est meilleur que les deux autres systèmes en terme de précision, par exemple pour 30 documents retournés, SARIOnto offre une amélioration de 20% de précision par rapport à lucene et 10% par rapport à l'autre système.

5. Conclusion et perspectives

La qualité des résultats fournis par les moteurs de recherche traditionnels n'étant pas toujours pertinente pour plusieurs requêtes ce qui a conduit les chercheurs à trouver des solutions reposant sur les ontologies. Dans ce sens, le travail présenté par cet article propose un système de recherche d'information en ligne basé sur cette

8 Revue. Volume X – n° x/année

architecture et utilisant deux ontologies, une de domaine et l'autre de services, afin d'améliorer la pertinence des résultats présentés à l'utilisateur. L'expérimentation et l'évaluation menées montrent une amélioration du taux de précision comparé à d'autres systèmes. Notre proposition permet de détecter facilement les services d'un domaine précis, et de retourner un ensemble de documents classés par rapport à ces services. Dans nos prochains travaux, l'utilisateur aura la possibilité d'importer une ontologie de domaine pour effectuer ses recherches. Une autre perspective concerne la visualisation spatiale des résultats en fonction de leur similarité. Ce travail permet aussi de fournir des documents pré classés et filtrés pour améliorer la construction des composants ontologiques et surtout les méthodes de constructions basées sur les techniques d'apprentissage (BenMustapha et al., 2007).

6. Bibliographie

- Berners-Lee, T., J. Hendler, O. Lassila, (2001). The Semantic Web, Scientific American.
- Baazaoui Zghal, H., Aufaure, MA. and Ben Mustapha, N. (2007). Extraction of Ontologies from Web Pages: conceptual modeling and tourism, Journal of internet Technologies Volume 8 No. 4 (2007), octobre 2007, ISSN 1607-9264.
- Baziz, M., (2004). Towards a Semantic Representation of Documents by Ontology-Document Mapping. The Eleventh International Conference on Artificial Intelligence(AIMSA 2004)
- Ben Mustapha, N., Baazaoui Zghal, H. and Aufaure, MA. (2007). A Prototype for knowledge extraction from semantic web based on ontological components construction, Internet technology, Web information System and Technolmogies (WEBIST07).
- Buckley, C., Singhal A., Mitra M., Salton G. (1996). New Retrieval Approaches using SMART. Proceedings of TREC'4, pages 25-48.
- Gruber, T. (1993). Toward principles for the design of ontologies used for knowledge sharing, International Journal of Human-Computer Studies, special issue on Formal Ontology in Conceptual Analysis and Knowledge Representation. Eds, Guarino, N. & Poli, R.
- Miller, G.A. (1995). WordNet: A Lexical Database for English. Communications of the ACM, 11, 39-41.
- Salton G., MacGill M.J. (1983). Introduction to modern information retrieval, McGraw Hill International Book Company, ISBN 0-07-Y66526-5.
- Schmid, H. (1994). Probabilistic Part-of-Speech Tagging Using Decision Trees. IMS-CL, Institut Für maschinelle Sprachverarbeitung, Universität Stuttgart, Germany.
- Stenback, J., Le Hors, A., Le Hégaret, P. (2003). Document Object Model (DOM) Level 2 HTML Specification, <http://www.w3.org/TR/>, 9.
- Voorhees, E. M., Gupta, N. K. et Johnson-Laird, B. (1994). The collection fusion problem. Proceedings of the 3rd Text REtrieval Conference (TREC-3). p. 95-104.