

# Etude de La Pertinence lors de La Sélection de Collections dans les Systèmes Distribués

Kheira Mechach\*, Lougmiri Zekri\*, Mustapha Kamel Abdi\*,\*\*

\*Département d'informatique  
Université d'Oran1 BP 1524 El-M'naouer Maraval, Oran, Algérie  
mechach.kheira@gmail.com, lougmiri@gmail.com  
\*\*abdi.mustapha@univ-oran.dz

## 1 Introduction à Sélection basée sur le Degré de Pertinence

Les bibliothèques numériques sont actuellement très répandues. Elles renferment des quantités d'informations énormes et nécessitent des mécanismes efficaces d'indexation et de manipulation. Les moteurs de recherche du type général ne peuvent pas les indexer car ils exigent que l'information qu'ils manipulent soit composée d'entités indépendantes. Dans le besoin de traiter rapidement et efficacement les requêtes, des méthodes basées des approches différentes ont été inventées. On rencontre alors, des méthodes se basant sur les réseaux bayésiens comme CORI Callan et al. (1995), d'autres méthodes qui se basent sur les statistiques TF\*IDF. Il existe aussi des méthodes qui se basent sur le modèle de langage et la pseudo-pertinence. Ces méthodes utilisent des résultats déjà obtenus pour de réponses futures. Puisque le modèle centralisé souffre du problème de passage à l'échelle, certaines méthodes ont été mises pour tourner sur les systèmes pair-à-pair. La méthode CORI a été une source d'inspiration et a été utilisée comme moyen de classification dans beaucoup de travaux. Cette méthode fonctionne sur un système bayésien pour localiser des réponses probables aux utilisateurs. La fonction de score donnée dépend de certains paramètres obtenus à partir d'expérimentations sur des datasets. Ce paramétrage fait que CORI est devenue instable. Ces paramètres doivent être réajustés pour chaque nouvelle collection. Afin de réduire le nombre de collections interrogées, Abbaci et al. (2002) présente la méthode CS. Celle-ci définit ndoc le nombre de documents à retourner et tient compte uniquement des deux premiers termes lors de l'évaluation des requêtes longues. Bien que l'objectif de réduction de flux est atteint, CS produit des faux positifs et faux négatifs importants à cause des restrictions imposées.

Soit un système distribué où un serveur appelé courtier est lié à un ensemble de serveurs. Le courtier détient un index Terme/Serveur qui indique pour chaque terme  $t_i$  la liste des serveurs qui le manipulent. Chaque serveur  $S_i$  est responsable d'une collection de documents  $c_i$  et manipule un index Terme/Documents. Cet index définit pour chaque terme  $t_i$  la liste des documents où il figure. Par cette définition, le courtier sélectionne de façon déterministe le sous-ensemble de serveurs pertinents. Ces index permettent de réduire la charge du système. Un document est jugé pertinent s'il partage au moins un terme avec la requête. Plus un document partage de termes avec la requête plus son degré de pertinence s'élève, induisant ainsi que le score d'une