Prise en compte de la structure des documents pour une indexation performante

Pascal Cuxac, Nicolas Kieffer

INIST-CNRS

2, allée du parc de Brabois 54519 Vandoeuvre lès Nancy pascal.cuxac@inist.fr, nicolas.kieffer@inist.fr

1 Introduction

L'indexation d'un document a des intérêts multiples : répondre à une requête, faire une classification, comparer des textes... (Gupta et Lehal, 2017).

Nous présentons Skeeft, une méthode automatique d'indexation en langage libre, sans ressource appliquée à du texte intégral (des articles scientifiques). Ce choix permet de s'affranchir du domaine scientifique et traiter un corpus multi-thématiques sans difficulté. L'analyse du texte intégral permet l'extraction de termes pertinents et une meilleure pondération. La structure du document est souvent utilisée afin de cibler des zones d'extractions ou pondérer les termes (You et al., 2013). Notre approche est très différente car nous ne hiérarchisons pas les différentes parties mais nous les mettons en concurrence : il suffit d'identifier les parties du document sans avoir à identifier leur "rôle" (introduction, méthodologie...).

2 Méthodologie

Nous assimilons un document à une classification dont chaque partie est une classe composée de termes. A partir de là le procédé se déroule en 4 grandes étapes :

- extraction des termes pour chaque partie identifiée (toute méthode est utilisable ici);
- application d'une sélection de variables : un terme est une variable, la partie ou il apparaît est sa classe d'appartenance (Lamirel et al., 2015);
- pondération des termes sélectionnés pour chaque partie;
- filtrage final, fusion des résultats et affichage.

Le lecteur se rapportera à l'article de Lamirel et al. (2015) pour plus de détails.

3 Résultats expérimentaux

Skeeft est testée sur des documents multi-disciplinaires du réservoir ISTEX ¹, indexés manuellement par des experts et comparée aux méthodes suivantes : TopicRank (Bougouin et

^{1.} https://api.istex.fr/documentation/