

Apport des traitements morpho-syntaxiques pour l'alignement des définitions par une classification SVM

Laura Dioşan^{*,**}, Alexandrina Rogozan^{*}, Jean-Pierre Pécuchet^{*}

^{*}LITIS (EA 4108) - INSA Rouen, France

^{**}Babeş Bolyai University, Computer Science Department, Cluj Napoca, Romania
lauras@cs.ubbcluj.ro, arogozan@insa-rouen.fr, pecuchet@insa-rouen.fr

Résumé. Cet article propose une méthode d'alignement automatique de définitions destinée à améliorer la fusion entre des terminologies spécialisées et un vocabulaire médical généraliste par un classifieur de type SVM (Support Vecteur Machine) et une représentation compacte et pertinente d'un couple de définitions par concaténation d'un ensemble de mesures de similarité, afin de tenir compte de leur complémentarité, auquel nous ajoutons les longueurs de chacune des définitions. Trois niveaux syntaxiques ont été investigués. Le modèle fondé sur un apprentissage à partir des groupes nominaux de type *Noms-Adjectifs* aboutit aux meilleures performances.

Les systèmes de recherche d'informations reposent sur une terminologie spécifique d'un domaine d'application que seuls les experts possèdent. En effet, les utilisateurs naïfs utilisent un langage généraliste pour formuler leurs requêtes. Pour qu'un système de recherche puisse répondre efficacement aux requêtes de ces derniers, il devrait pouvoir tirer parti des liens sémantiques entre des concepts véhiculés dans le langage généraliste et dans le langage spécialisé. Une des tâches du projet *VODEL* est de réaliser un alignement automatique de définitions, c'est-à-dire de mettre en correspondance des définitions associées à un même concept, mais ayant des vedettes différentes. Le cadre choisi étant celui du domaine médical, les ressources terminologiques de spécialité sont tirées du thésaurus *MeSH* et du dictionnaire *VIDAL*, alors que le vocabulaire généraliste est représenté par des définitions appartenant à l'encyclopédie *Wikipédia* et au réseau sémantique LDI de *Memodata*¹.

Aligner deux définitions revient à résoudre efficacement un problème de classification binaire supervisée. Notre modèle d'alignement passe par deux étapes : premièrement, une représentation compacte des définitions et deuxièmement, une classification supervisée de couples de définitions. Chaque définition a été représentée par un sac des mots, après un traitement linguistique (segmentation, lemmatisation et étiquetage morpho-syntaxique) permettant de filtrer les mots vides et de ne garder que les noms (*N*), les noms et les adjectifs (*NA*), et respectivement les noms, les adjectifs et les verbes (*NAV*). Nous proposons une représentation compacte et pertinente d'un couple de définitions par concaténation d'un ensemble de mesures de similarité classiques (Matching, Dice, Jaccard, Overlap, Cosine), afin de tenir compte de leur complémentarité, auquel nous ajoutons les longueurs de chacune des définitions. Nous proposons un alignement des terminologies par un classifieur de type SVM (Séparateur à Vaste

¹Le corpus de définitions a été réalisé dans le cadre du projet *VODEL* par G. Lortal, I. Bou Salem et M. Wang.