

Etude comparative de modèles d'extraction non supervisée de mots-clés pour la recommandation d'emploi

Bissan Audeh*, Maia Sutter**, Christine Largeron**

* Inasoft, 2507 avenue de l'Europe, 69140, Rillieux-La-Pape, France
bissan.audeh@inasoft.fr,

** Laboratoire Hubert Curien UMR 5516, F-42023, Saint-Etienne, France
maia.sutter@etu.univ-st-etienne.fr, chistine.largeron@@univ-st-etienne.fr

1 Introduction

La recherche d'emploi peut être facilitée en augmentant la visibilité des postes qui correspondent le mieux à un candidat. Dans un tel contexte, il est essentiel d'attribuer des mots-clés aux offres d'emploi et aux CV pour réaliser l'indexation et la mise en correspondance mais aussi effectuer des analyses intéressantes pour la prise de décision quant au marché de l'emploi. Dans cet article, nous cherchons à évaluer dans quelle mesure les mots-clés extraits avec des approches non supervisées peuvent représenter du contenu textuel même sans apprentissage et sans connaissance externe dans un contexte de recommandation d'emploi. Nous avons sélectionné six méthodes non supervisées d'extraction de mots-clés que nous avons évaluées sur de vrais offres d'emploi et CV anonymisés. Pour cela nous avons élaboré un protocole d'évaluation qui inclut la construction d'un gold standard.

2 Cadre méthodologique

Nous avons évalué quatre modèles parmi les plus courants dans l'état de l'art : deux modèles statistiques TF-IDF et RAKE (Rose et al. (2010)), TextRank (Mihalcea et Tarau (2004)) un modèle basé sur le score PageRank et KeyBERT (Grootendorst (2020)) qui utilise les plongements de texte. En plus, nous avons évalué les extensions KeyBERT+ et TF-IDF+ des modèles KeyBERT et TF-IDF respectivement, où nous remplaçons la sélection de mots-clés candidats des méthodes par une sélection basée sur le modèle de langue français de Spacy¹ qui permet d'identifier des termes composés. Nous avons adapté toutes ces méthodes à notre corpus par des pré-traitements et post-traitements. Les données disponibles pour ce projet étaient constituées de 818 CV et 858 offres d'emploi mais seul un sous-ensemble a été annoté en vérifiant manuellement 12316 mots-clés; ce qui a permis de construire un jeu d'évaluation contenant 57 CV et 29 offres d'emploi. Le texte des CV a déjà été extrait et anonymisé en supprimant les données personnelles, et les offres d'emploi étaient au format HTML. Bien que la taille de ce jeu soit limitée, l'expérimentation a néanmoins permis de comparer les performances des algorithmes.

1. https://spacy.io/models/fr#fr_core_news_lg