

Analyse comparée et interprétation des résultats de trois classifications de textes littéraires

Baptiste Bohet¹, Nicole Vincent²

¹Université Sorbonne Nouvelle – baptiste.bohet@sorbonne-nouvelle.fr

²Université Paris Cité – nicole.vincent@u-paris.fr

Abstract

This article presents a comparative study of the automatic classification results of different methods. The aim is to categorize complex literary texts according to specific problems: dating the text, identifying the author's gender, determining the language in which the text was written (i.e. whether it is translated or not) and, in the case of translation, recognizing the original language (among five: English, Italian, Spanish, German and Japanese). In addition to showing that these methods are effective, with the majority of results showing recognition rates in excess of 90%, this study proposes hypotheses for interpreting the results.

Keywords: statistical analysis, textual data, text-mining, NLP, LLM, explainability, BERT, TF-IDF, deep learning, Artificial Neural Network (ANN).

Résumé

Cet article présente une étude comparative des résultats de différentes méthodes de classification automatique. L'objectif est la catégorisation de textes littéraires, donc complexes, en fonction de problématiques spécifiques : la datation du texte, l'identification du sexe de l'auteur, la détermination de la langue de rédaction, c'est-à-dire si le texte est traduit ou non, et enfin, dans le cas d'une traduction, la reconnaissance de la langue originale (parmi cinq : anglais, italien, espagnol, allemand, japonais). En plus de montrer que ces méthodes sont efficaces, la majorité des résultats font apparaître des taux de reconnaissance supérieurs à 90 %, cette étude propose des hypothèses d'interprétation des résultats.

Mots clés : analyses statistiques, données textuelles, fouille de textes, TALN, grand modèle de langue (LLM), explicabilité, BERT, TF-IDF, apprentissage profond, Réseaux de Neurones Artificiels (RNA).