

Spécialité : Informatique Laboratoire : LINA

Équipe : TALN

Extraction automatique de termes-clés à partir de sujets

Directrice : Béatrice DAILLE Adrien BOUGOUIN E-mail : adrien.bougouin@univ-nantes.fr

Encadrant : Florian BOUDIN

Problématique

- Les termes-clés sont des unités textuelles capables de synthétiser le contenu d'un document
- Quelle est la nature linguistique des termes-clés ?
- Comment distinguer les termes-clés des autres unités textuelles de même nature linguistique ?

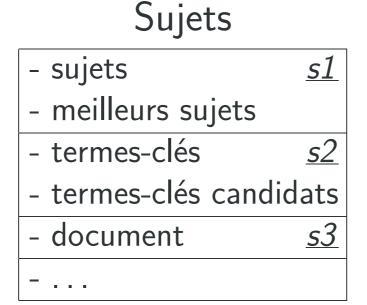
Solution proposée

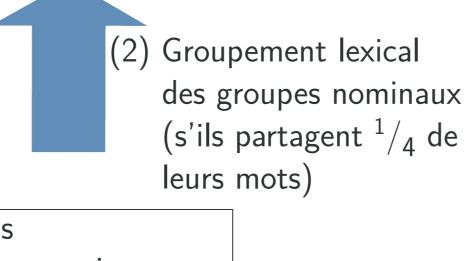
- 1. Extraction des groupes nominaux du document
- 2. Groupement des groupes nominaux en sujets (concepts)
- 3. Ordonnancement par importance des sujets
- 4. Sélection d'un terme-clé pour chaque sujet

Modélisation du lien sémantique entre chaque sujet

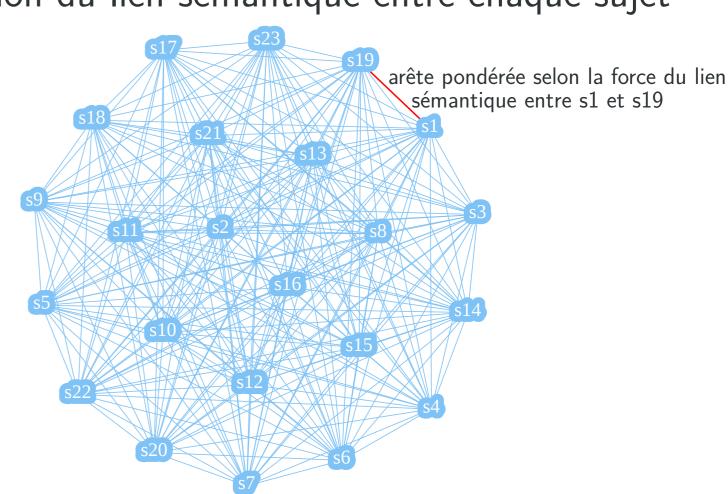


(1) Extraction des plus longues séquences de noms, noms propres et adjectifs





- sujetsmeilleurs sujets
- termes-clés
- termes-clés candidats
- document
- Groupes nominaux



- (3) Ordonnancement des sujets, selon TextRank [1]
- (4) Pour les k meilleurs sujets : extraction du groupe nominal apparaissant en premier dans le document



Termes-clés

Résultats

- Comparaison de TopicRank avec TextRank [1]
- Application à 100 articles journalistiques (WikiNews)
- Extraction de 10 termes-clés par document

Méthode	Précision	Rappel	F-mesure
TextRank	9,3	8,3	8,6
TopicRank	35,0	37,5	35,6

Perspectives

- Usage de connaissances linguistiques pour le groupement en sujets
- Exploration de diverses stratégies de selection du terme-clé le plus représentatif d'un sujet

Références

[1] Rada Mihalcea and Paul Tarau.

TextRank: Bringing Order Into Texts.

In Dekang Lin and Dekai Wu, editors, *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 404–411, Barcelona, Spain, July 2004. Association for Computational Linguistics.



Ce travail a bénéficié d'une aide de l'Agence Nationale de la Recherche portant la référence (ANR-12-CORD-0029).