



HAL
open science

Extraction de connaissances basée sur l'analyse formelle de concepts en vue de l'assistance aux débats en ligne

Imen Ben Sassi, Hani Guenoune, Alexandre Bazin, Marianne Huchard,
Mathieu Lafourcade, Jean Sallantin

► To cite this version:

Imen Ben Sassi, Hani Guenoune, Alexandre Bazin, Marianne Huchard, Mathieu Lafourcade, et al.. Extraction de connaissances basée sur l'analyse formelle de concepts en vue de l'assistance aux débats en ligne. TextMine @EGC 2024, Jan 2024, Dijon, France. lirmm-04459556

HAL Id: lirmm-04459556

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-04459556v1>

Submitted on 15 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Extraction de connaissances basée sur l’analyse formelle de concepts en vue de l’assistance aux débats en ligne

Imen Ben Sassi, Hani Guenoune, Alexandre Bazin, Marianne Huchard,
Mathieu Lafourcade, Jean Sallantin

LIRMM, Université de Montpellier, CNRS, Montpellier, France
imen.ben-sassi@lirmm.fr, prenom.nom@lirmm.fr

Résumé. Nous présentons un processus automatisé d’accompagnement de débats visant à extraire des associations entre termes à partir des listes de termes-clés des arguments, listes co-construites par les utilisateurs et notre système d’indexation. L’indexation encourage les utilisateurs à compléter ou corriger la liste des termes-clés, agissant comme un outil incitatif à l’élaboration de points de vue plus structurés. L’algorithme est basé sur l’analyse formelle de concepts et s’appuie sur la base de connaissances JeuxDeMots (JDM). La procédure fait intervenir plusieurs modules menant à une étape d’extraction de connaissances sous forme d’implications destinées à être intégrées dans JDM. Cette approche coopérative permet à la base de connaissances de s’enrichir à mesure que les débats sont analysés, améliorant les termes-clés suggérés par la plate-forme.

1 Introduction

Le projet AREN-DIA (ARGumentation Et Numérique - Didactique & Intelligence Artificielle)¹ vise à sensibiliser les élèves à la pratique du débat dans le cadre de leur éducation à la citoyenneté. La plate-forme de débats en ligne ainsi conçue est également disponible à la société civile, assurant une éthique aux débats. La plate-forme offre la possibilité d’engager des débats structurés à partir d’un texte, renouvelant l’approche traditionnelle des échanges argumentatifs. Elle présente la particularité d’intégrer une technologie collaborative de Traitement Automatique du Langage en vue d’augmenter l’efficacité du processus de débat. Dans cette perspective, AREN-DIA se déploie selon un axe didactique et un axe IA. Les résultats des études menées au sein des lycées révèlent que l’utilisation judicieuse de ce logiciel, insérée dans un dispositif didactique approprié, conduit à un essor marqué des compétences argumentatives chez les élèves (Bächtold et al., 2023).

L’axe IA et ses enjeux font l’objet de cet article. Nous nous intéressons à la manière de concevoir un mécanisme de renforcement incitant les utilisateurs à participer à l’amélioration du système d’IA produisant une représentation structurée des propos d’un débat. AREN² est une application web qui offre un espace de débat, réunissant un ensemble d’utilisateurs. Le débat porte sur un texte support publié en amont sur la plate-forme. Les utilisateurs interviennent à travers des commentaires exprimant une opinion, une argumentation ou un avis

1. Ce projet est financé par l’Agence Nationale de la Recherche : ANR-22-FRAN-0001.

2. La plate-forme est accessible via le lien suivant : <https://portail-aren.lirmm.fr/aren2023/>

sur un segment du texte support ou un commentaire préalablement publié, créant ainsi des embranchements dans l'arbre général du débat. Outre l'intervention des débattants, une procédure automatique vient compléter chaque commentaire en suggérant des termes-clés synthétisant les propos tenus. Cette opération d'indexation représente le point de départ de l'analyse et de l'accompagnement du débat par la machine. Elle est soumise à une complétion par les utilisateurs, qui seront invités à valider, invalider ou compléter ces termes-clés par ceux qu'ils estiment manquants. Afin de lever l'ambiguïté sémantique résultant de la polysémie des termes proposés, nous avons recours à une étape d'enrichissement sémantique des termes pour les préparer à l'opération d'extraction de connaissances basée sur l'analyse formelle de concepts (AFC). Ces connaissances, sous forme d'implications, seront utilisées pour mettre à jour les relations dans la base de connaissances exploitée lors de ces processus, JeuxDeMots (JDM).

L'article s'organise comme suit. Nous détaillons les étapes du fonctionnement général d'AREN dans la section 2. Nous nous pencherons également sur l'algorithme d'accompagnement du débat dans la section 3, qui consiste en la production de termes-clés et d'une analyse AFC pour produire des associations de termes pertinentes. Nous définissons ensuite, dans la section 4, les différentes métriques utilisées pour évaluer l'utilité de l'augmentation sémantique des termes d'indexation des propos du débat. Nous comparons, dans la section 5, les résultats obtenus avec l'AFC avant et après l'enrichissement sémantique des termes d'indexation.

2 Fonctionnement de la plate-forme AREN

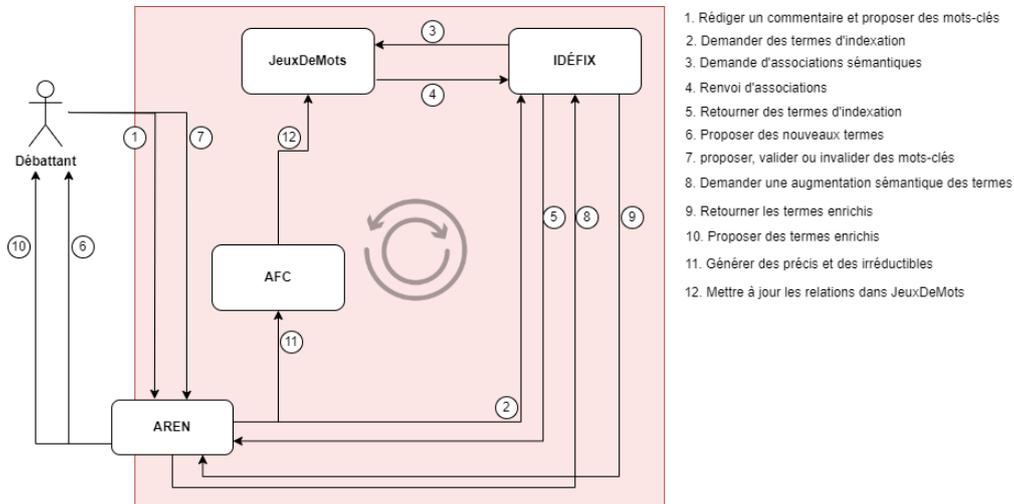


FIG. 1 – Fonctionnement de la plate-forme AREN sous forme de cycles entre les interactions utilisateurs, le calculateur de termes-clés IDÉFIX, l'analyse formelle de concepts (AFC), et la base de connaissance JeuxDeMots.

L'application constitue un espace de débat, faisant intervenir un ensemble d'utilisateurs. Un débat porte sur un texte support publié en amont sur la plate-forme, il est conjointement formé par le contenu du texte ainsi qu'un ensemble de commentaires créés par les utilisateurs

et exprimant une opinion ciblée, une argumentation ou un avis sur un segment du texte support ou un commentaire préalablement publié.

Chaque intervention utilisateur comporte une *position* (d'accord ou pas d'accord), une *reformulation*, une *argumentation* et des *mot-clés*. La partie du texte que l'utilisateur souhaite commenter est choisie en sélectionnant, dans le texte de départ, le segment correspondant. La possibilité de sélectionner un segment aussi bien dans le texte support que dans un commentaire préalablement publié, permet la création d'embranchements dans l'arbre général du débat (c.f. Figure 1). Les commentaires sont constitués d'un ensemble d'informations construisant le propos de l'utilisateur, parmi ces informations se définit, entre autres, la position que prend le débattant (d'accord, pas d'accord) vis-à-vis de la sélection (le segment auquel il réagit). Les champs de texte libres, de *reformulation* et d'*argumentation*, sont prévus afin de consolider puis définir, l'avis du débattant.

3 Algorithme d'accompagnement des débats

Nous présentons dans cet article une IA d'accompagnement de débats, KeepTalk (Knowledge Extraction for Enhanced online Public Talks and Argumentative Learning Know-how), dont un des objectifs est d'extraire des associations nouvelles entre termes à partir des listes de termes-clés des arguments d'un débat. L'approche est organisée en plusieurs modules aboutissant à une étape d'extraction de connaissances (c.f. Section 3.3) alimentée par l'analyse formelle de concepts. Après augmentation lexicale (c.f. Section 3.2), cette étape permet de créer des implications entre termes (par exemple, si A est présent alors B et C sont aussi présents). Les implications produites sont destinées à être introduites dans la base de connaissances. Par exemple, si nous avons $A \rightarrow B, C$, alors dans le réseau lexical JDM nous ajouterons : A *r_associated* B et A *r_associated* C. La procédure de description thématique (c.f. Section 3.1) sur laquelle s'assoit l'algorithme s'inscrit dans une démarche collaborative, itérative et incrémentale. Les ensembles de termes indexant chaque commentaire sont co-construits d'un côté, par la procédure automatisée (*IDÉFIX*), et de l'autre, par une *supervision* et *complétion* par les utilisateurs des termes extraits par *IDÉFIX*. Cette supervision est permise en donnant à l'utilisateur la possibilité de *proposer*, *valider* ou *invalidier* des termes de l'ensemble proposé par l'IA accompagnant le débat. Ce retour est pris en compte lors des itérations de descriptions thématiques ultérieures, menant à une indexation de meilleure qualité. L'objectif étant d'assurer une amélioration de la base de connaissances à mesure que des débats sont analysés, avec en retour une amélioration des termes-clés suggérés par la plate-forme (via *IDÉFIX*) pour les arguments d'un débat. En outre, ce mécanisme est pensé de manière à inciter les utilisateurs à proposer des termes-clés complétant les propos du débat. Plus précisément, le calcul automatique de termes-clés pour un argument est un moyen de donner envie aux utilisateurs, et en particulier à l'auteur de l'argument, de compléter voire de corriger la liste des termes-clés proposés. Un mauvais terme-clé sera en général considéré par l'utilisateur comme une tache/erreur insupportable devant être nettoyée/corrigée.

3.1 Indexation thématique

Les divers arguments des participants au débat sont contenus dans des textes bruts et non-structurés. L'indexation thématique des commentaires a pour objectif d'associer ces données

textuelles à une représentation structurée permettant de synthétiser les propos par des ensembles de termes-clés, référencés dans des bases de connaissances et pouvant servir de point d'entrée à une procédure automatisée. Les termes extraits peuvent désigner des concepts évoqués dans le texte ou des unités lexicales dont la saillance au sein du commentaire est jugée importante. Cette étape d'extraction de mots-clés s'appuie sur des connaissances externes issues du réseau lexico-sémantique *JDM* (Lafourcade et Le Brun, 2023), et est réalisée par le service *IDÉFIX*³.

JDM est un réseau lexico-sémantique sous forme de graphe orienté. Les nœuds du graphe représentent les termes, tandis que les arcs désignent des relations typées, pondérées et potentiellement annotées entre les termes. Le graphe représente la polysémie des mots en explicitant des raffinements sémantiques hiérarchisés, où un sens spécifique est affilié au sens général du terme. Basé sur une série de notions, principes et outils originaux (ex. la notion de raffinement, la palette des types de relations sémantiques - les éléments d'information, des liens sémantiques entre un type de relation et son inverse (*r_isa* et *r_hypo*, par exemple), l'outil contributif *Diko*, etc.), le réseau *JDM* est conçu pour une utilisation humaine et comme support de connaissances pour des processus d'intelligence artificielle (analyse sémantique de texte, raisonnement, assistance à la prise de décision, résumé automatique, etc.). Un système de pondération (arcs pondérés, éventuellement négatifs) et de valuation symbolique (annotation en méta-informations, par exemple : rare, pertinent, non pertinent, etc.) a été mis en œuvre pour faciliter des heuristiques de parcours du graphe ainsi que son exploitation.

Le réseau *JDM* peut être utilisé avec des algorithmes classiques liés aux bases de connaissances, mais également sous forme de réseau neuronal (approches hybrides, algorithmes de propagation et de rétro-propagation, etc.). Parmi ces algorithmes, *IDÉFIX* est une sur-couche du réseau *JDM* fondée sur des réseaux de neurones permettant de sélectionner des concepts pertinents pour un texte fourni en entrée. Cette sélection se fait de manière abductive et locale au commentaire, par imitation des exemples déjà appris des interactions précédentes avec l'utilisateur (validation, invalidation et proposition de termes-clés).

3.2 Enrichissement sémantique

Afin d'assurer une représentativité des propos des utilisateurs, nous procédons à l'enrichissement des ensembles de mots-clés produits à l'étape précédente. Nous cherchons, en premier lieu, à assurer une couverture sémantique suffisante en nous occupant des éventuels phénomènes d'ambiguïté lexicale⁴ engendrés par la polysémie des termes-clés (c.f. Figure 2). Ceci revient à *séparer les termes semblables en apparence, mais dont les sens sont différents*, en identifiant les raffinements sémantiques adéquats dans le réseau *JDM*.

L'enrichissement des termes de description par leurs termes synonymes ou hyperonymes pertinents, permet, à l'inverse de la désambiguïsation, de *regrouper les termes différents en apparence, dont les sens sont (quasi-)semblables*. L'ajout d'un synonyme dans les termes-clés d'un propos peut être réalisé sans ou avec restriction, c'est-à-dire dans ce dernier cas, si ce synonyme existe déjà comme mot-clé d'un autre propos du débat, ceci afin d'éviter une dérive liée à des cas de synonymie foisonnante.

3. L'outil *IDÉFIX* est accessible via le lien : https://www.jeuxdemots.org/intern_extract.php

4. Ambiguïté traitée via le service *Belléophon* : https://www.jeuxdemots.org/intern_desamb.php

Commentaire : <i>la monnaie locale est un outil financier.</i>
Indexation : outil conceptuel; être utile; outil>moyen d'action; MLC; économie locale; moyen d'action; monnaie locale; crise commerciale; monnaie locale complémentaire et citoyenne; outil; financier; Sol-violette; économie; monnaie locale complémentaire; monnaie; outil financier; local
Désambiguïsation de l'indexation : outil>moyen d'action; monnaie>argent; économie>activité économique; financier>finance; MLC>monnaie locale complémentaire; monnaie>unité monétaire; local>propre à un lieu
Augmentation sémantique - synonymes : régional (depuis local>propre à un lieu); sous>argent (depuis monnaie>argent)

FIG. 2 – Exemple de désambiguïsation et d'augmentation sémantique d'indexation d'un propos d'un débat sur les monnaies locales. L'ajout du synonyme régional n'est autorisé que parce qu'il est présent ailleurs dans le débat (dans le cas de l'augmentation avec restriction).

3.3 Extraction de connaissances

L'extraction de connaissances à partir de l'indexation des commentaires utilise l'AFC, un cadre mathématique basé sur la théorie des treillis permettant la représentation de l'information contenue dans des données sous des formes algébriques ou logiques (Ganter et Wille, 2012).

3.3.1 Contexte formel et fermeture de Galois

L'AFC part de données sous la forme d'un *contexte formel*; un triplet $(\mathcal{O}, \mathcal{A}, \mathcal{R})$ où $\mathcal{R} \subseteq \mathcal{O} \times \mathcal{A}$ est une relation binaire entre des *objets* O et les *attributs* A qui les décrivent. Cette relation peut être représentée sous la forme d'un tableau de croix (c.f. Figure 3).

Dans AREN, les objets sont les commentaires du débat et les attributs sont les mots-clés (ou les termes) proposés par les débattants ou ajoutés lors de la phase d'indexation. Un terme est en relation avec un commentaire s'il l'indexe. Par exemple, dans l'exemple de la Figure 3, $(c_4, \text{monnaie} > \text{argent}) \in \mathcal{R}$ signifie que l'objet $c_4 = \text{la loi donne une existence légale aux monnaies locales}$ est indexé par le terme $\text{monnaie} > \text{argent}$.

	<i>organisation</i>	<i>monnaie > argent</i>	<i>cours légal</i>	<i>monnaie complémentaire</i>	<i>monnaie locale</i>
c_1	×		×		
c_2			×	×	×
c_3		×			×
c_4		×	×	×	×

FIG. 3 – Exemple de contexte formel avec une relation binaire entre quatre commentaires (c_i): $c_1 = \text{« la monnaie est une manière de faire et d'organiser la société »}$; $c_2 = \text{« L'acceptation dans le cadre de la loi rend la monnaie locale légale »}$; $c_3 = \text{« les monnaies locales nous font nous questionner sur un outil que nous banalisons la monnaie »}$; $c_4 = \text{« la loi donne une existence légale aux monnaies locales »}$ et cinq termes (t_j): « organisation »; « monnaie>argent »; « cours légal »; « monnaie complémentaire »; « monnaie locale ».

Un contexte formel donne lieu à deux *opérateurs de dérivation*, tous deux notés \cdot' et définis tels que

$$\cdot' : 2^{\mathcal{A}} \mapsto 2^{\mathcal{O}}$$

Extraction de connaissances pour l'accompagnement de débats en ligne

$$\begin{aligned}
 A' &= \{o \in \mathcal{O} \mid \forall a \in A, (o, a) \in \mathcal{R}\} \\
 \cdot' &: 2^{\mathcal{O}} \mapsto 2^{\mathcal{A}} \\
 O' &= \{a \in \mathcal{A} \mid \forall o \in \mathcal{O}, (o, a) \in \mathcal{R}\}
 \end{aligned}$$

Les compositions \cdot'' de ces opérateurs forment des opérateurs de fermeture. Par exemple, dans la Figure 3, la fermeture de $\text{monnaie} > \text{argent}$ est $\{\text{monnaie} > \text{argent}, \text{monnaie locale}\}$.

3.3.2 Génération des irréductibles

Un contexte formel est dit *clarifié* s'il n'a pas deux objets ayant exactement la même description ou deux attributs décrivant exactement les mêmes objets. Dans un contexte clarifié, un attribut a est dit *irréductible* si l'ensemble $\{a\}'$ des objets qu'il décrit n'est pas égal à l'intersection des ensembles d'objets décrits par d'autres attributs (Liquiere, 2021), c'est-à-dire qu'il n'existe pas d'ensemble d'attributs X tel que $\{a\}' = \bigcap_{x \in X} \{x\}'$. Dans l'exemple de la Figure 3, seul l'attribut *monnaie complémentaire* n'est pas irréductible puisque $\{\text{monnaie complémentaire}\}' = \{\text{cours légal}\}' \cap \{\text{monnaie locale}\}'$. Le reste des termes, à savoir *organisation*, *monnaie > argent*, *cours légal* et *monnaie locale*, sont tous des irréductibles.

3.3.3 Extraction des implications

Nous cherchons à extraire des régularités dans la cooccurrence des mots-clés dans l'indexation des commentaires. L'AFC offre différentes possibilités de représentation de ces régularités : implications, règles d'association, treillis de concepts ou relations causales (Bazin et al., 2022). Une implication est une règle constituée d'une paire d'ensembles d'attributs A et B , habituellement notée $A \rightarrow B$. Une implication est dite *valide* dans un contexte formel donné si et seulement si tous les objets décrits par les attributs de A sont aussi décrits par les attributs de B , c'est-à-dire $B \subseteq A''$. Ainsi, dans l'exemple de la Figure 3, les deux implications $\{\text{cours légal}, \text{monnaie locale}\} \rightarrow \{\text{monnaie complémentaire}\}$ et $\{\text{organisation}\} \rightarrow \{\text{cours légal}\}$ sont valides tandis que $\{\text{cours légal}\} \rightarrow \{\text{organisation}\}$ ne l'est pas. Afin de réduire le nombre de règles à présenter aux débattants, notre attention se focalise spécifiquement sur les implications de la forme $\{a\} \rightarrow B$ telles que a est un terme irréductible.

3.4 Enrichissement de la base de connaissance

Les implications obtenues avec l'AFC sont utilisées pour mettre à jour les relations dans la base de connaissances exploitée lors de ce processus. Donc, depuis une implication de la forme $\{a\} \rightarrow \{b, c, d, e, \dots\}$ nous ajoutons dans la base de connaissances des relations $a \rightarrow \mathbf{x}$ avec $\mathbf{x} \in \{b, c, d, e, \dots\}$.

Dans l'exemple de la Figure 3, la mise à jour de la base de connaissances JDM se fait par l'ajout de l'association des termes « *cours légal* » et « *organisation* » et celle de « *monnaie locale* » et « *monnaie > argent* ». Ces modifications améliorent globalement la composante associative des calculs ultérieurs des indexations des propos.

4 Mesures d'évaluation des règles

Afin d'étudier l'impact de l'augmentation sémantique sur la qualité des règles, nous avons utilisé diverses métriques, notamment le support, la nouveauté et la surprise (fondée sur la co-occurrence ou le voisinage des termes).

4.1 Support

Le support peut-être perçu comme un indicateur de « confiance statistique » d'une règle. Le support d'un ensemble d'attributs ou termes T est le nombre d'objets (ou de commentaires) décrits par T divisé par le nombre total d'objets. Il peut être défini par l'Equation 1.

$$Supp(r) = p(T_r^p \ T_r^c) / |C| \quad (1)$$

où T_r^p et T_r^c sont respectivement les termes de la prémisse et de la conclusion de la règle r et C sont les commentaires.

4.2 Nouveauté

La nouveauté est une métrique qui a été utilisée dans les domaines de découverte de sous-groupes et de découverte de clauses (Wrobel, 1997). Une règle est considérée nouvelle si sa prémisse et sa conclusion ne sont pas statistiquement indépendantes (Lavrac et al., 1999).

La nouveauté d'une règle est définie par l'Equation 2.

$$Nov(r) = p(T_r^p \ T_r^c) - p(T_r^p) \ p(T_r^c) \quad (2)$$

où r est une règle (implication), T_r^p et T_r^c sont respectivement les termes de la prémisse et de la conclusion de la règle r .

4.3 Surprise

Bien que la pertinence peut être facilement évaluée à l'aide du support, la mesure de la surprise (ou de l'inattendu) des règles est une tâche complexe qui nécessite souvent des études coûteuses à mener, impliquant des utilisateurs (ou des ressources externes, dans notre cas). Une règle nouvelle peut être rétrospectivement surprenante ou non, dans le sens où la connaissance disponible ne permet pas de l'expliquer rapidement/facilement.

Dans ce travail, nous ajustons deux définitions de la mesure de surprise utilisées dans le domaine de recommandation (Kaminskas et Bridge, 2014), l'une basée sur le degré d'association sémantique entre les termes indexant les propos du débat et l'autre basée sur les termes associés aux termes d'indexation. Les deux mesures produisent un score qui indique le niveau de surprise que le terme cible a apporté à la règle.

4.3.1 Surprise basée sur la co-occurrence des termes

L'information mutuelle spécifique (Point-wise mutual information notée PMI) indique à quel point deux termes sont statistiquement dépendants, en fonction du nombre de propos indexés par les deux termes et chaque terme séparément (c.f. Equation 3). Les valeurs de PMI

varient entre -1 et 1 , où -1 signifie que les deux termes ne sont jamais utilisés ensemble pour indexer un propos, 0 signifie l'indépendance des termes et 1 signifie une co-occurrence systématique des termes.

$$PMI(i, j) = \log_2 \frac{p(i, j)}{p(i)p(j)} / -\log_2 p(i, j) \quad (3)$$

où $p(i)$ et $p(j)$ représentent respectivement les probabilités qu'un propos soit indexé par les termes i et j , tandis que $p(i, j)$ est la probabilité qu'un propos soit indexé par les deux termes i et j .

Sur la base de la PMI, la mesure de surprise d'un terme i pour la règle r est définie comme la valeur moyenne de PMI des termes dans la règle (c.f. Equation 4).

$$Surprise_{co-occ}^{avg}(i, r) = 1 - \frac{1}{|T_r|} \sum_{j \in T_r} PMI(i, j) \quad (4)$$

où i est un terme, r est une règle (implication) et T_r sont les termes de la règle r .

La surprise basée sur la co-occurrence permet de tenir compte du contexte local du débat et des rapprochements de termes que celui-ci peut engendrer. Toutefois, l'indépendance statistique n'implique pas une similarité sémantique faible. En effet, deux contributeurs peuvent respectivement préférer utiliser le terme « *vélo* » et « *bicyclette* ». Ces deux termes sont alors, dans le débat, en co-occurrence nulle ou faible, alors qu'ils sont sémantiquement très proches. La surprise basée sur le contenu sémantique des termes (leur voisinage, c.f. Section 4.3.2) permet de tenir compte de ce type de phénomènes.

4.3.2 Surprise basée sur le voisinage des termes

Notre deuxième mesure de surprise est basée sur la distance appliquée aux termes associés aux termes cibles. Le voisinage d'un terme t dans la base de connaissances JDM est l'ensemble des termes auquel t est relié par la relation d'association d'idées. Nous avons utilisé le complément de la métrique de similarité de Jaccard pour comparer les termes (c.f. Equation 5).

Par exemple, sont associés au terme « *monnaie* » de façon non-exhaustive les termes : « *argent, pièce, billet, euro, devise* », le terme « *fric* » aura comme termes associées : « *argent, pièce, billet, euro, thune* ». La distance entre ces deux termes est de $1 - 4/6 = 1/3$.

$$dist(i, j) = 1 - \frac{A_i \cap A_j}{A_i \cup A_j} \quad (5)$$

où A_i et A_j sont respectivement les ensembles de termes associés aux termes i et j . Dans le cas où le terme A est polysémique, on considère sa désambiguïsation lexicale pour extraire les termes qui sont associés au contexte des règles.

Pour mesurer la surprise d'un terme, nous calculons la distance moyenne entre le terme cible i et les autres termes T_r de la règle r comme indiqué dans l'équation 6.

$$Surprise_{vois}^{avg}(i, r) = \frac{1}{|T_r|} \sum_{j \in T_r} dist(i, j) \quad (6)$$

où i est un terme, r une règle (implication) et T_r sont les termes de la règle r .

5 Résultats et discussions

Afin d'aider à l'interprétation des résultats de l'algorithme, nous commençons dans cette section par présenter les données ayant servi à cette évaluation. Nous cherchons ensuite à mettre en évidence la pertinence de chaque module employé, ceci en mettant en place des configurations contrastives de l'algorithme rendant possible la comparaison des résultats permis par chaque sous-module.

5.1 Jeux de données et configurations

Nous procédons à l'évaluation de notre approche à l'aide des données issues d'un débat sur la plate-forme AREN concernant les monnaies locales⁵ intitulé « Les monnaies locales sont-elles un outil pour sauver l'économie locale et dans quelles conditions ? ». Les principales caractéristiques de notre jeu de données sont présentées dans le Tableau 1.

Débatants	Arguments	Mots-clés	Période
8	48	464	Mars 2020 – Mai 2023

TAB. 1 – Statistiques du débat sur les monnaies locales.

Chaque argument d'un débattant est associé à un texte initial du débat et décrit par une reformulation, une phrase qui reflète sa compréhension du texte argumenté (« *La monnaie locale est un outil financier* » : Figure 2), et une opinion (83.33% des arguments sont « *plutôt d'accord* » et 16.67% ne sont « *plutôt pas d'accord* »). En total, 464 mots-clés distincts ont été utilisés pour indexer les reformulations dont 125 termes uniques sont proposés par les utilisateurs et 339 par IDÉFIX. En moyenne, chaque débattant a utilisé 5.39 termes par argument.

Nous comparons les résultats de trois variantes de notre approche pour mesurer l'effet de l'augmentation sémantique sur la qualité des résultats de l'AFC. Les détails de nos méthodes sont énumérés ci-dessous :

- KT : Les implications sont calculées à partir du contexte d'extraction initial, défini par la relation binaire entre les reformulations des débattants et les termes-clés qui les indexent.
- KT^\dagger : Le contexte d'extraction est enrichi par les synonymes des termes qui définissent les attributs pour générer les implications. On s'intéresse aux termes synonymes qui sont déjà utilisés lors de l'indexation initiale (R^\bullet) et aussi ceux qui ne le sont pas (R°), donc, de nouveaux termes n'apparaissant pas dans le débat.
- KT^\ddagger : Identique à la configuration précédente avec des hyperonymes au lieu de synonymes.

5.2 Résultats

Nous commençons par proposer une vue quantitative des résultats des différentes configurations. Nous rapportons, dans le Tableau 2, le nombre d'*attributs*, *irréductibles* et *implications* dans les configurations se limitant, ou pas, aux termes-clés du débat.

Quand l'ajout de termes n'est pas restreint à ceux du débat, nous observons une augmentation du nombre d'attributs. Inversement, si on se restreint aux termes du débat, le nombre

5. <https://portail-aren.lirmm.fr/aren2023/debates/6>

Extraction de connaissances pour l'accompagnement de débats en ligne

d'attributs est constant. Dans tous les cas, l'utilisation de la désambiguïsation lexicale réduit le nombre d'objets produits (irréductibles et implications), ceci est conforme à l'intuition car la désambiguïsation réduit l'éparpillement lexical. Par ailleurs, l'ajout d'hyperonymes est plus productif que l'ajout de synonymes, car il est possible de trouver au moins un hyperonyme pour la quasi-totalité des attributs (qui sont des termes), ceci est beaucoup moins vrai pour les synonymes.

	KT		KT^\dagger		KT^\ddagger	
			DL°	DL^\bullet	DL°	DL^\bullet
<i>Avec restriction aux termes-clés du débat</i>						
Attributs	464	464	464	464	464	464
Irréductibles	70	73	68	83	79	79
Implications	43	54	46	75	71	71
<i>Sans restriction aux termes-clés du débat</i>						
Attributs	464	3831	2240	1125	866	866
Irréductibles	70	161	85	121	103	103
Implications	43	114	50	106	88	88

TAB. 2 – Résultats de KT (KeepTalk) avec et sans restriction aux termes-clés du débat : Le nombre d'objets demeure constant pour toutes les configurations et est égal à 48. DL°/DL^\bullet désignent l'utilisation ou non de la tâche de désambiguïsation lexicale.

La première expérimentation rapporte la proportion des relations d'association qui sont considérées comme correctes/pertinentes. Ces associations sont générées à partir des règles produites (implications). Cette étape consiste en une évaluation menée manuellement par 4 intervenants adoptant le rôle « d'experts ».

Sans augmentation	60.11 %	
	DL°	DL^\bullet
Augmentation avec restriction (R^\bullet)	63.12 % (1)	72.07 % (2)
Augmentation sans restriction (R°)	42.60 % (4)	76.77 % (3)

TAB. 3 – Pourcentage des bonnes associations selon une évaluation manuelle menée par 4 experts.

Nous cherchons à travers le Tableau 3 à classer les configurations de notre système en termes de « qualité », du point de vue d'utilisateurs humains. Le cas 1 signifie que dans une configuration se restreignant aux termes du débat et sans procédure de désambiguïsation, 63.12% des relations d'association à ajouter à la base de connaissance sont jugés correctes par les experts. Pour la meilleure configuration (cas 3), où nous procédons à une désambiguïsation sans se limiter aux termes-clés du débat, nous obtenons 76.77% de bonnes associations.

Dans le Tableau 4, nous constatons que la nouveauté est globalement très faible, ce qui indique que l'on trouve peu d'associations n'existant pas dans la base de connaissances. Ceci est positif du point de vue de la complétude de la base. On constate par ailleurs que la surprise est globalement très haute, ce qui veut dire qu'une information nouvelle n'aurait pas pu être

	KT	KT^{\dagger}		KT^{\ddagger}	
		DL°	DL^{\bullet}	DL°	DL^{\bullet}
Avec restriction aux termes-clés du débat : R°					
Support	0.0667	0.1300	0.0836	0.1658	0.1390
Nouveauté	0.0546	0.0681	0.0564	0.0733	0.0690
Surprise par co-occurrence	0.5662	0.7995	0.5304	0.1270	0.1294
Surprise par voisinage	0.9488	0.9623	0.9706	0.9674	0.9718
Score agrégé	0.2103	0.2872	0.2215	0.1965	0.1863
Score agrégé syn+hyper		0.4838		0.4079	
Sans restriction aux termes-clés du débat : R°					
Support	0.0667	0.1352	0.0978	0.1667	0.1359
Nouveauté	0.0546	0.0780	0.0615	0.0799	0.0676
Surprise par co-occurrence	0.5662	0.7167	0.4534	0.3481	0.1957
Surprise par voisinage	0.9488	0.9502	0.9602	0.8997	0.9497
Score agrégé	0.2103	0.2911	0.2262	0.2541	0.2032
Score agrégé syn+hyper		0.5452		0.4294	

TAB. 4 – Comparaison des résultats de l’analyse formelle de concepts : KT avec le contexte initial; KT^{\dagger} avec le contexte augmenté avec les synonymes; KT^{\ddagger} avec le contexte augmenté avec les hyperonymes. L’augmentation est faite avec et sans restriction aux termes du débat.

inférée, dans la quasi-totalité des cas. Ceci est un autre résultat très positif qui justifie l’utilité de notre approche d’extraction de connaissances. Nous observons un effet conjoint à l’étape de

Sans augmentation	0.1264	
	DL°	DL^{\bullet}
Augmentation avec restriction (R^{\bullet})	0.3048 (1)	0.2939 (2)
Augmentation sans restriction (R°)	0.2290 (4)	0.3264 (3)

TAB. 5 – Combinaison des résultats des métriques avec la proportion des bonnes associations (évaluation manuelle) - Il s’agit d’un score et non d’un pourcentage.

désambiguïsation DL et à la restriction R ou non aux termes-clés du débat. Ce croisement est clarifié dans le Tableau 5. La configuration la plus favorable est celle avec une augmentation avec synonymes et hyperonymes sans R et avec étape de DL (cas 3). La seconde meilleure configuration est celle sans DL et avec R (cas 1). La pire configuration, qui a de très mauvais résultats, est la combinaison de R° et DL° (cas 4). Le score du dernier cas (cas 2 : R^{\bullet} et DL^{\bullet}), bien que correct, est inférieur aux cas 1 et 3.

La désambiguïsation lexicale (DL) et la restriction (R) aux termes déjà présents dans le débat, visent le même but, contrôler le foisonnement lexical, et ne pas tomber dans le piège de polysémie. L’approche avec R permet de ne pas introduire de termes qui ne sont pas apparus dans le débat, il n’y a donc aucune chance d’introduire, par accident, un terme sans rapport. Le cas 2 est intéressant car il n’est pas intuitif : en effet, on s’attendrait à ce que l’action conjointe de D et R donne les meilleurs résultats, or ce n’est pas le cas. A priori l’effet restrictif conjoint de DL et R empêche un rapprochement efficace des propos du débat. Ne pas faire de R permet d’augmenter la richesse des associations, toutefois cette richesse doit être contrôlée par la DL .

6 Conclusion et perspectives

Les résultats obtenus sont prometteurs et soulignent l'efficacité d'effectuer conjointement une analyse basée sur l'AFC et une augmentation lexicale à partir d'une base de connaissances. En perspective, il serait important, sur la base des scores de l'évaluation manuelle, d'agréger les scores des métriques automatiques de manière à obtenir un score global qui serait représentatif de la qualité (que nous avons cherché à obtenir ici par une évaluation manuelle). Concernant l'extraction des connaissances, en perspective, le projet explorera d'autres représentations des régularités : autres implications, règles d'association et relations causales ; ce qui permettra d'ajouter dans la base de connaissances des informations sur des types de relations plus précises (autres que les associations d'idées).

Références

- Bächtold, M., G. Pallarès, K. De Checchi, et V. Munier (2023). Combining debates and reflective activities to develop students' argumentation on socioscientific issues. *Journal of Research in Science Teaching* 60(4), 761–806. 1
- Bazin, A., M. Couceiro, M.-D. Devignes, et A. Napoli (2022). Steps towards causal formal concept analysis. *International Journal of Approximate Reasoning* 142, 338–348. 6
- Ganter, B. et R. Wille (2012). *Formal concept analysis : mathematical foundations*. Springer Science & Business Media. 5
- Kaminskas, M. et D. Bridge (2014). Measuring surprise in recommender systems. In *Proceedings of the Workshop on recommender systems evaluation : dimensions and design held in conjunction with RecSys 2014, REDD '14, Silicon Valley, USA*. ACM. 7
- Lafourcade, M. et N. Le Brun (2023). Apport du jeu pour la construction de connaissances : le projet jeuxdemots. *Technologie et innovation* 8(Le jeu pour innover). 4
- Lavrac, N., P. A. Flach, et B. Zupan (1999). Rule evaluation measures : A unifying view. In *Proceedings of the 9th International Workshop on Inductive Logic Programming, ILP '99, Berlin, Heidelberg*, pp. 174–185. Springer-Verlag. 7
- Liquiere, M. (2021). Utilisation des irréductibles d'un treillis de concepts pour la sélection de motifs. Technical report, LIRMM, Université de Montpellier. 6
- Wrobel, S. (1997). An algorithm for multi-relational discovery of subgroups. In *Principles of Data Mining and Knowledge Discovery*, Berlin, Heidelberg, pp. 78–87. 7

Summary

We present an online debate analysis algorithm aiming to extract new associations between terms from co-constructed keyword lists of arguments by users and our indexing system. The calculation of keywords encourages users to supplement or correct the keyword list, serving as an incentive tool for developing more structured contributions. The algorithm is based on formal concept analysis and relies on the JeuxDeMots knowledge base. The procedure involves multiple modules leading to a knowledge extraction step in the form of implications which are to be integrated into JDM. This cooperative approach allows the knowledge base to enrich itself as debates are analyzed, improving the platform's suggested keywords.