

***Désambiguïstation sémantique des expressions
polysémiques en traitement automatique du langage*****LIMAME-BEN SALAH Dalila****UR Sémantique, Syntaxe & Pragmatique (Sousse/ Tunisie)****Centre Lucien TESNIERE (Besançon/France)****Résumé :**

Les expressions figées constituent un problème de poids pour le traitement automatique, du fait de leur caractère polylexical, figé, non compositionnel et polysémique. Nous nous proposons de présenter « une » méthode permettant de contrer ces différents obstacles.

Le système proposé est un système d'analyse conduisant à la reconnaissance des expressions figées polysémiques par le biais de filtres. Certes, certains indices (morphologique et syntaxique) permettent de définir si une expression est libre ou figée, mais parfois ces indices ne suffisent pas à déterminer et à identifier l'expression potentielle. La catégorisation de ces unités polysémiques et leur actualisation en contexte ont l'avantage de réduire la polysémie laissant place uniquement à une seule des différentes acceptions (réelle ou virtuelle).

1. Introduction

Les unités minimales porteuses de sens, autrement dit les mots se répartissent en deux groupes : les mots polysémiques¹ et les mots monosémiques². Les expressions figées nommées également expression idiomatique, séquence, syntème, locution, mots composés, lexie complexe, phraséologie, praxème, proverbe (í) ont une double caractéristique : l'une est de nature formelle, combinatoire, associative, l'autre est de nature sémantique. La première suppose qu'une lexie 1 se combine avec une lexie 2 pour produire un sens inhérent et un sens afférent. Les langues naturelles renferment un nombre considérable d'expressions plus ou moins figées : *tête de Turc*, *perdre la boussole*, *prendre le taureau par les cornes*, *casser sa pipe*. Le traitement automatique d'expressions figées est problématique tant au niveau de la reconnaissance qu'au niveau de la génération (en traduction automatique)³. En définissant séparément les composants d'une expression *prendre*, *taureau*, *cornes*, le calcul du sens ne pourra être que littéral, compositionnel. Or, force est de remarquer que le sens d'une expression n'est *aucunement* le résultat du sens de ses composants. La représentation du sens lexical est une donnée importante en TAL car elle réduit voire annule les ambiguïtés et ceci par un calcul

¹ Ils constituent la majorité et permettent de limiter le lexique d'une langue naturelle.

² Ils représentent l'exception, ainsi les dénominations propres aux domaines de spécialité (vocabulaires spécialisés) tels que le vocabulaire professionnel, scientifique ou technique possède un sens univoque et une valeur référentielle unique.

³ Les systèmes de traduction automatique connaissent des carences notamment en ce qui concerne le cadre sémantique.

Désambiguïisation sémantique des expressions polysémiques en traitement automatique du langage sémantique en contexte. Il est indéniable que le problème de la polysémie reste au cœur de l'étude des expressions.

2. Le figement

2.1. R.Martin (1997)

La locution est un syntagme figé, situé au-delà du mot et en deçà de la phrase figée

La phrase figée, pour lui, englobe les proverbes ou la phrase idiomatique en d'autres termes les phrases toutes faites spécifiques à une certaine situation du type *le jeu n'en vaut pas la chandelle*. Les propriétés caractéristiques du figement se résume à des restrictions sélectionnelles, la non-compositionnalité et à sa valeur intensionnelle⁴. La restriction sélectionnelle revient à parler de restriction combinatoire, certaines unités lexicales acceptent des associations variables. Ainsi le verbe *prendre* accepte un nombre considérable de combinaison (*prendre au collet, prendre de la terre avec une pelle, prendre sur son compte, prendre un verre*), en revanche le verbe *hocher* admet un champ combinatoire fort limité (*hocher la tête*). Plus le champ combinatoire est restreint plus on entre dans la sphère des expressions figées. Il est à signaler que certaines unités n'apparaissent qu'au sein même d'une expression (*convoler en justes noces, se mettre martel en tête*).

Une unité lexicale est figée sémantiquement lorsque le sens est non compositionnel, R. Martin définit la non compositionnalité à travers deux processus : *l'enrichissement sémantique* et *la démotivation étymologique*. Le premier renvoie à un plus sémantique, à l'au-delà des mots. Ils sont empreints d'une culture

⁴ renvoie aux tests transformationnels

Désambiguïstation sémantique des expressions polysémiques en traitement automatique du langage

(avoir le feu vert, avoir un cœur d'artichaut), d'un passé (le talon d'Achille). La métaphorisation est le procédé le plus usité qui permet la création d'expressions. Ainsi, le sens métaphorique (non compositionnel ou afférent) peut parfois concurrencer le sens compositionnel.

D'autre part, une expression telle que *porter le chapeau*⁵ reproduit un fait historique qui n'est plus d'actualité.

2.2. Critères du figement

Le degré de figement des expressions s'établit par le biais de différents critères. D'une part, par la polylexicalité, une expression est constituée de plusieurs unités lexicales, par le blocage des propriétés transformationnelles et distributionnelles (paradigmatique)..

Ainsi, l'expression *mettre les pieds dans le plat*

Paul a mis les pieds₀ dans le plat₁

Pronominalisation Paul les₀ a mis dans le plat*

Paul a mis les pieds dedans*

Passivation Les pieds ont été mis dans le plat par Paul*

Extraction les pieds, Paul a mis dans le plat ??

Dans le plat, Paul a mis les pieds ?

Relativisation les pieds que Paul a mis dans le plat*

Clivage Ce sont les pieds que Paul a mis dans le plat

Insertion Paul a mis les [petits] pieds dans le [grand] plat ?

Nominalisation la mise des pieds dans le plat*

Le blocage des propriétés transformationnelles renforce le caractère figé de l'expression.

⁵ à l'époque médiévale, le condamné portait un chapeau grotesque.

3. La polysémie

La notion même de polysémie est à définir étant donné que certains linguistes (tel que D. LePesant, 2003) considèrent comme étant polysémiques des constructions possédant au moins deux sens, le sens figuré et le sens propre ; pour d'autres (tel que S. Méjri, 2003) la polysémie des expressions figées ne doit pas être perçue de la même façon que pour les expressions monolexicales. Ainsi, S. Méjri parle de dédoublement sémantique : une signification littérale et une signification globale et non de polysémie. De plus, il attribue le terme de polylexicalité aux *séquences figées* et le terme de polysémie aux *unités unilexicales*. D'autre part, parler de la polysémie d'unité monolexicale n'est pas parler d'unité polylexicale, en d'autres termes nous pouvons nous poser la question de savoir si le sens de *mettre les pieds dans le plat* renvoie à la polysémie du verbe *mettre* ou des noms *pieds* ou *plat*. La polysémie des unités monolexicales⁶ rapproche des paradigmes sémantiques d'une seule et même unité ce qui est différent de la polysémie des expressions figées polylexicales. Cette dernière suppose un déséquilibre entre signifié et signifiant, le poids du signifiant est plus conséquent que celui du signifié.

Deux interprétations sont à concevoir : une première interprétation, celle de l'expression littérale, compositionnelle et une seconde interprétation, celle de l'expression prise dans sa totalité l'expression polylexicale. Seulement, il est également possible de rencontrer des expressions possédant deux voire trois

⁶ Unicité du signifiant et multiplicité du signifié, ceci implique une variété considérable des significations.

Désambiguïisation sémantique des expressions polysémiques en traitement automatique du langage

significations globales. C'est le cas de l'expression *avoir du coffre* où l'unité lexicale *coffre* est polysémique et qui signifie avoir du souffle et, avoir du courage ou encore de l'expression *mettre les pieds dans le plat* possède certes une signification compositionnelle mais à ses côtés nous rencontrons deux autres significations⁷ : aborder une question délicate avec une franchise brutale et, commettre une gaffe⁸.

La polysémie est loin d'être un fait marginal, bien au contraire elle est inhérente à toute langue naturelle. La désambiguïisation sémantique d'une expression polysémique⁹ s'effectue par le contexte (phrastique ou textuel) et l'isotopie sémantique. Le contexte immédiat d'une unité lexicale permet de déceler le sens adéquat de ses différentes occurrences. La récurrence de sèmes communs permet de mettre en relation la signification potentielle de l'expression avec la signification des éléments qui la précèdent ou qui la suivent.

4. Système et sous-systèmes

Dans ce qui suit nous présentons un système de reconnaissance (par contraintes) permettant de définir l'état d'une expression (libre ou figée) mais également de définir son sens en contexte.

⁷ Le sens dit figuré ou non compositionnel n'est pas forcément une métaphore du sens propre mais peut être également une métonymie. Autrement dit, le sens non compositionnel constitue une extension sémantique.

⁸ Définitions du *Petit Robert*, 2002

⁹ mais également d'une unité monolexicale.

4.1. Filtres ou contraintes syntaxiques

En utilisant des critères formels (transformationnels et distributionnels), nous avons pu créer plusieurs matrices. Une expression libre se prête aux diverses contraintes transformationnelles¹⁰ et distributionnelles¹¹. L'auditeur d'une expression établit une suite d'hypothèses par laquelle le sens sera défini en relation avec une série de contextes de vérité. En énonçant une expression figée, le locuteur s'exprime sans établir de correspondance entre le contenu sémantique de l'expression énoncée et son énoncé. Le problème qui se pose lors de l'interprétation d'une expression figée est celui du rendu sémantique, d'une interprétation fidèle¹². Les structures syntaxiques représentent un domaine construit et analysable par conséquent nous avons élaboré une catégorisation des expressions selon diverses structures syntaxiques

¹⁰ Les contraintes transformationnelles sont l'extraction, la pronominalisation, la relativisation, la substitution, l'insertion d'éléments, la négation, l'inachèvement (ce critère reste une donnée intéressante pour la traitement automatique, notamment concernant les proverbes ou tournures (*au petit bonheur* *au petit bonheur la chance*)

¹¹ les contraintes distributionnelles sont le paradigme, l'aspect verbal, la négation, la catégorie du sujet selon qu'il s'agisse d'un animé ou d'un inanimé (*Paul vole de ses propres ailes* *cet oiseau vole de ses propres ailes*), l'actualisation (en d'autres termes la détermination : l'absence d'article renvoie à un état archaïque, à une non actualisation ; la présence de l'article est un indicateur soulignant l'état de figement *mettre les points sur les i* *mettre les points sur le i**)

¹² le caractère de ressemblance, *il pleut des cordes*, la lexie *cordes* évoque toutes les caractéristiques aspectuelles de l'objet corde ; cette expression exprime bien plus que l'objet en soi (épaisseur, horizontal). Il s'agit d'une économie langagière.

Désambiguïstation sémantique des expressions polysémiques en traitement automatique du langage

possibles. Les expressions sont réparties en trois catégories :

Catégorie I regroupe les expressions sans contrainte (littérale ou idiomatique) : *tourner la page* admet toute transformation.

Catégorie II regroupe les expressions possédant des contraintes pour les expressions idiomatiques allant de 1 à 5 contraintes. Ainsi, *perdre la main* admet toutes les transformations hormis la relativisation qui n'est envisageable que dans le cas d'une expression littérale. Il est des cas tel que *tuer le ver*, cette expression n'admet que deux transformations (l'affirmation et l'insertion), les autres transformations sont valides pour l'expression littérale uniquement.

Catégorie III se scinde en deux groupes :

Groupe 1 présente des expressions ayant des contraintes au niveau du sens littéral.

Groupe 2 présente des expressions ayant des contraintes au niveau du sens littéral et du sens idiomatique.

4.2. Filtres ou contraintes distributionnelles

Cette nouvelle série de contraintes (sur le sujet et sur le déterminant) va permettre ou bien de confirmer les premières hypothèses ou de les infirmer.

4.2.1. Sujet

Le sujet constitue un élément de base dans la reconnaissance d'une expression idiomatique. Une expression telle que *avalier des coulevres* ne peut prétendre à une interprétation idiomatique avec un sujet

Désambiguïstation sémantique des expressions polysémiques en traitement automatique du langage

ANH¹³ (*ce chien a avalé des couleuvres*) mais avec un sujet AH ou NANH.

Chercher la petite bête possède un sens littéral avec un sujet ANH.

Pour *mordre à l'hameron*, l'interprétation littérale se déduit de la catégorie du sujet ainsi un sujet ANH (tel que *poisson*) renvoie à la littéralité en revanche un sujet AH renvoie à l'idiomaticité.

4.2.2. Déterminant

Le déterminant est l'élément d'actualisation qui se produit sous la forme de l'article (défini ou indéfini) ou de la catégorie (article ou adjectif (numéral, possessif ou démonstratif)).

Article défini vs article indéfini

Ainsi, *dérouler le tapis rouge* peut prétendre à une interprétation idiomatique ou littérale alors que *dérouler un tapis rouge* ne peut admettre qu'une interprétation littérale.

Article singulier vs article pluriel

Tirer les ficelles admet les deux interprétations alors que *tirer la ficelle* admet une seule interprétation possible (littérale)

Absence d'article vs présence d'article

Acheter chat en poche est une structure archaïque, idiomatique, en revanche *acheter un chat* est une structure libre.

¹³ ANH : animé non humain ; AH : animé humain ; NANH : non animé non humain

4.2.3. Sémantisme

Certaines expressions outrepassent ces premiers filtres : *retourner sa veste* et *vider son sac* ne possèdent pas de contraintes que ce soit sur le plan syntaxique ou sur le plan distributionnel. Par conséquent, une analyse plus approfondie doit être menée pour ces expressions restées dans le flou.

Le cadre phrastique en général se révèle insuffisant pour la reconnaissance automatique de l'expression idiomatique puisque bien souvent les traducteurs automatiques peinent face à ces expressions et donnent une traduction mot à mot ce qui fausse l'interprétation. Un contexte plus large est plus opportun ; certains éléments textuels procurent des indices permettant de lever toute ambiguïté sémantique.

5. Désambiguïsateur

Notre système permet d'identifier une expression telle que *prendre une veste* possédant un sens littéral et un sens idiomatique et de lever l'ambiguïté, qui plus est il permet également la reconnaissance d'expressions polysémiques et l'interprétation sémantique en contexte. Prenons l'exemple *marcher sur les pieds*. L'analyse morphologique repère les mots clefs inclus dans une expression par le biais d'une base de données lexicales. Cette dernière nous informe sur les lexies (temps, nombre, partie du discours, etc), leur reconnaissance se fait au moyen de leurs racines et de leurs flexions : *pieds* représente un mot clef dans l'expression idiomatique dans sa forme pluriel uniquement ; l'unité prise au singulier ne pose aucun problème d'analyse dans le cadre syntagmatique énoncé et constitue ainsi une expression

Désambiguïstation sémantique des expressions polysémiques en traitement automatique du langage

libre. Notre base de données lexicales prend la forme :
<ped> <ped, nom, masculin, pluriel, +mot clef>

Le mot clef est mis en relation avec les mots collocatifs pouvant entrer dans le cadre d'une expression polysémique.

Casser sa pipe mot clef : <pipe>

Éléments collocatifs : verbe *casser* possessif : *sa, leur*

Si le système rencontre le mot *pipe*, il poursuivra sa recherche au niveau des mots collocatifs *casser* et les possessifs. Dans le cas où le système ne trouve aucun élément collocatif, il déduira une lecture littérale. Certaines expressions potentielles élimine d'emblée l'interprétation idiomatique dans la mesure où le mot clef ne possède pas les critères morphologiques requis. *Plomb* ne peut prendre la forme pluriel dans aucunes de nos expressions.

L'étape suivante s'intéresse à l'analyse syntaxique en d'autres termes à l'étiquetage du texte en faisant ainsi usage d'un analyseur morphologique, le système Labelgram (Cardey & Greenfield, 2003). A chaque mot correspond une fonction, la suite de fonctions constitue un patron syntaxique qui permettra la mise en relation avec la base de données syntaxiques. Chaque expression rassemble ses propres structures syntaxiques selon les possibilités transformationnelles.

<lever le pied>

<lever Ø pied> <X, liste des structures syntaxiques>

dès lors qu'une instanciation du patron syntaxique est possible, dans ce cas la séquence est syntaxiquement

Désambiguïstation sémantique des expressions polysémiques en traitement automatique du langage

valide et constitue alors une expression polysémique potentielle. Les expressions polysémiques refermant des structures syntaxiques spécifiques sont reconnues qu'elles soient littérales ou idiomatiques. Les matrices transformationnelles donnent accès à une reconnaissance (d'une certaine catégorie) d'une expression.

Danser devant le buffet >

Paul ne danse pas devant le buffet

la négation indique une expression libre, l'expression idiomatique rejette la structure négative. Toute expression potentielle ayant été identifiée subit une analyse des structures transformationnelles et une analyse distributionnelle qui vient appuyer celle-ci. En effet, une expression potentielle peut valider tous les tests transformationnels, autrement dit ne posséder aucune contrainte. L'expression sera reconnue au niveau structural comme étant une expression potentielle non catégorisée, par conséquent les matrices distributionnelles seront un moyen de catégoriser l'expression. Ainsi, *mordre à l'hameçon* fait partie de la Catégorie I, la structure affirmative de l'expression *Paul a mordu à l'hameçon* [N, Vbconj, Prép, N] est acceptée que ce soit pour le sens idiomatique ou le sens littéral. L'échec d'identification va conduire le système à poursuivre sa recherche au moyen d'une analyse distributionnelle puisque certaines expressions ont des contraintes au niveau du sujet et/ou du déterminant. Ici, le sujet AH impliqué par Paul renvoie après vérification de la matrice distributionnelle sujet à une expression validée idiomatique.

Désambiguïisation sémantique des expressions polysémiques en traitement automatique du langage

Paul tient la chandelle passe à travers tous les filtres cités précédemment, le système va analyser le contexte immédiat de l'expression en relevant les mots avoisinants figurant à droite et à gauche de l'expression les unités permettant de désambiguïiser le sens d'une expression. Ayant établi une liste d'éléments avoisinants potentiels, ces unités sont vérifiées dans cette liste comprenant des noms mais également des verbes.

Avaler des couleuvres peut se trouver à proximité d'unités telles que *nourrir, alimenter, bouche, ventre, festin*

Verbe {nourrir, alimenter, mastiquer, affamer, engloutir, mâcher, dévorer}

Nom {bouche, ventre, faim, nourriture, repas, collation, festin, dîner, déjeuner}

Si, après vérification, le système valide l'unité *festin* donc il en déduira qu'il s'agit d'une expression littérale.

Conclusion

Le système présenté a pour but d'extraire des expressions potentielles de les analyser et également de permettre la génération en vue d'une traduction automatique. Il est encore aujourd'hui difficile de délimiter les expressions figées puisque certaines sont plus ou moins figées, et de les chiffrer puisque d'un linguiste à l'autre les chiffres varient.

Références

CARDEY S. & GEENFIELD P. (2003), *Disambiguating and tagging using systemic grammar*, Proceedings of the 8th international symposium on social communication, Santiago de Cuba, pp.559-564

ENJALBERT, P. (2005), *Sémantique et traitement automatique du langage naturel*, Paris : Hermès.

GIGUET, E., VERGNE, J. (1997), From part of speech tagging to memory-based deep syntactic analysis, *Proceedings of International Workshop of Parsing Technologies*, Massachusetts : MIT

KLEIBER, G. (1994), *Problèmes de sémantique: la polysémie en question*, Paris : PUF

LIMAME D., ALSHARAF H., CARDEY S., GREENFIELD P., SKOURATOV I. (2003), Fixedness, the complexity and fragility of the phenomenon: some solutions for natural language processing, in Acts of XVII International Congress of Linguists, Prague

MARTIN, R. (1997), Sur les facteurs du figement lexical, in *Locution entre langue et usage*, Fontenay-aux-roses : ENS Editions.

MARTINS-BALTAR, M. (éd.), (1997), *La locution entre lexicale, syntaxe et pragmatique*, Paris : Klincksieck.

MEJRI, S. (éd.), (2003), *Syntaxe & Sémantique 5 : Polysémie et polylexicalité*, Caen : Presses Universitaires de Caen.

VICTORRI, B. & FUCHS, C. (1996), *La polysémie : construction dynamique du sens*, Paris : Hermès.