

Disambiguating automatic semantic annotation based on a thesaurus structure

Véronique MALAISÉ¹, Luit GAZENDAM², Hennie BRUGMAN³

¹ Vrije Universiteit, Amsterdam

² Telematica Institute, Enschedé, Netherlands

³ Max Planck Institute for Psycholinguistics, Nijmegen, Netherlands

vmalaise@few.vu.nl, Luit.Gazendam@telin.nl,

Hennie.Brugman@mpi.nl

Résumé. La relation *voir/employé pour* d'un thesaurus est souvent plus complexe que la (para-)synonymie recommandée par l'ISO-2788, standard décrivant le contenu de ces vocabulaires contrôlés. Le fait qu'un non descripteur puisse renvoyer à plusieurs descripteurs (seuls les descripteurs sont pertinents dans le cadre de l'indexation contrôlée) fait que cette relation est complexe à utiliser dans un contexte d'annotation automatique : elle génère des cas d'ambiguïté. Dans ce papier, nous présentons CARROT, un algorithme que nous avons mis au point pour classer les résultats de notre chaîne de traitements pour l'Extraction d'Information, et son utilisation dans le cadre de la sélection du descripteur pertinent lorsque plusieurs choix sont possibles. Cette sélection s'adresse à des documentalistes, dans le but de simplifier et d'accélérer leur travail, et se base sur la structure de leur thesaurus. Nous arrivons à un succès de 95 % dans nos suggestions ; nous discutons ces résultats et présentons des perspectives à cette expérimentation.

Abstract. The *use/use for* relationship a thesaurus is usually more complex than the (para-) synonymy recommended in the ISO-2788 standard describing the content of these controlled vocabularies. The fact that a non preferred term can refer to multiple preferred terms (only the latter are relevant in controlled indexing) makes this relationship difficult to use in automatic annotation applications : it generates ambiguity cases. In this paper, we present the CARROT algorithm, meant to rank the output of our Information Extraction pipeline, and how this algorithm can be used to select the relevant preferred term out of different possibilities. This selection is meant to provide suggestions of keywords to human annotators, in order to ease and speed up their daily process and is based on the structure of their thesaurus. We achieve a 95 % success, and discuss these results along with perspectives for this experiment.

Mots-clés : désambiguïssation sémantique, algorithme de classement, annotation automatique.

Keywords: word sense disambiguation, ranking algorithm, automatic annotation.

1 Introduction

Thesauri are controlled vocabularies, often used for indexing and retrieving documents from collections. The standard thesauri contain two types of elements, preferred and non preferred terms, related with a link called *use/use for*. This link is considered as (para-)synonymy in the ISO-2788 standard (ISO, 1986) and can thus be useful for (semi-) automatic indexing applications : it enables a program to index a document with a preferred term (which is the type of thesaurus based controlled annotation we are interested in) either if the document contains an occurrence of the preferred term or if it contains occurrences of the corresponding non preferred term. In reality, this *use/use for* relationship is often more complex, and can generate ambiguity problems when used “as is” in an automatic application. We present in this paper the solution that we have developed in our project for selecting the relevant preferred term, given an occurrence of an ambiguous non preferred term in a text. This selection algorithm is based on the thesaurus’s structure. The thesaurus we used in this experiment is the GTAA, which is employed for indexing and retrieving TV programs at the Netherlands Institute for Sound and Vision, the Dutch national TV archives. Our project, CHOICE¹, is collaborating with this Institute and focuses on easing and speeding up the work of cataloguers by providing them with a ranked set of keywords referring to their thesaurus’ entries as indexing suggestions. We will present our project’s goal and the specificity of this use case in the following section (section 2), followed by a description of thesauri in general and the GTAA itself (section 3). In this section, we will show the different semantics of the *use/use for* relationships and the problem of having multiple links between preferred and non preferred terms. We then present our annotation pipeline (section 3.4), including the algorithm that we elaborated to rank the extracted keywords, and that we propose here for selecting the relevant preferred term out of multiple possibilities (section 3.5). Section 5 shows our experiment to evaluate this algorithm in this Word Sense Disambiguation context. We achieved a 95 % of success, but are still facing minor and more important problems. We discuss them and conclude with perspectives for this experiment in section 6.

2 The CHOICE project

Charting the Information Landscape Employing Context Information, the CHOICE project deals with the suggestion of metadata from textual resources to annotate video documents. In the context of the Dutch TV archives, the cataloguers check a set of textual documents, on top of watching the program itself, to make their descriptions. One of the goals of our project is to build on existing Information Extraction platforms, extend and tune them to our specific needs in order to cope with the particularities of this specific use case and provide the cataloguers with a relevant set of keywords as indexing suggestions. Our Information Extraction is based on the content of the thesaurus that they are currently using at Sound and Vision, enriched and transformed by us. We present this thesaurus in the following section, and the specificity of our task in the section describing our ranking algorithm.

¹<http://www.nwo.nl/CATCH/CHOICE>

3 The GTAA thesaurus

3.1 A thesaurus according to the ISO 2788 standard

A thesaurus is *The vocabulary of a controlled indexing language, formally organized so that the a priori relationships between concepts (for example as “broader” and “narrower”) are made explicit.*²

Although this definition mentions *concepts*, a thesaurus contains terms (preferred and non preferred terms), organized according to 5 relationships : broader term (BT), narrower term (NT), related term (RT), use (US³) and use for (UF). A preferred term is *A term used consistently when indexing to represent a given concept [...]*⁴, whereas a non preferred term should not be used for indexing, but is only useful at search time to point different words possibly expressing the same idea towards the one that has been chosen to represent it in the thesaurus. In practice, as we detail in the following section, this relationship can encode different kinds of links, as suggested by these examples : hurricane UF cyclone, insurgent UF guerilla's, organ UF church organ, oven UF magnetrons, octopus UF calamary.

The other relationships should stand only between preferred terms. BT relates a term with a more generic one, supposed to index a larger set of documents. For example *Means of transportation* is a BT of *Bus*. NT is the relationship between a term and a more specific one, that should be used to index a subset of the documents indexed by the more generic one (*Bus* and *School bus*). RT is a non hierarchical relationship between two terms in the same domain, as *Bus* and *Driver*, for example.

3.2 The GTAA

The GTAA thesaurus, a Dutch acronym for “Common Thesaurus for Audiovisual Archives”, is the controlled vocabulary used for the Sound and Vision documentation process. It contains approximately 160.000 terms. They are divided in 6 disjoint facets : Keywords (about 3800 preferred terms), Locations (about 14.000), Person Names (about 97.000), Organization-Group-Band Names (about 27.000), Maker Names (about 18.000) and Genres (113 preferred terms). The thesaurus mainly uses constructs as presented in the ISO 2788 standard and commonly used in companies or institutions : amongst others, use, use for, broader term, narrower term, related term. Terms from all facets of the GTAA may have related terms and use for relationships, but only Keywords and Genres can also have broader term/narrower term relations, organizing them into a set of hierarchies. Additionally, Keyword terms are thematically classified in 88 subcategories of 16 top Categories (Nature, Society,...). Although the data model that is used for the thesaurus allows links between terms across facets, no instances of these links currently exist. This experiment concerns only automatic indexing with terms from the Keyword facet.

²(ISO, 1986), section 3-Definitions.

³In general this relationship is encoded USE, but this acronym is the one used in the GTAA.

⁴(ISO, 1986), section 3-Definitions.

3.3 Different semantics and non uniqueness of the Use relationship

In the GTAA, there are 1377 US relationships, *i.e.* 1377 times a non preferred term is associated with a preferred term. Some of these non preferred terms are associated with multiple different preferred terms and some of the preferred term are associated with multiple non preferred terms. In the first case, the non preferred terms is polysemic or has different domains of application, each meaning or domain having an explicit preferred term : for example, the non-preferred term *minority*⁵ has two preferred terms, *ethnic minority* and *religious minority*. In the latter case (one preferred term associated with different non preferred terms), different notions were grouped under one common and single preferred term. This is either done for easing the thesaurus' use (the fewer terms there are, the easier it is to find the most appropriate one when indexing), or because the distinction was not relevant for indexing the TV programs of Sound and Vision : for example, the preferred term *diplomats* groups two non preferred terms, *ambassadors* and *consuls*. When having a close look at the nature of the US, UF relationship we see four different types :

- Synonyms : To cleanse US To clean
- Meronym : Sabbath US Jewish religion
- Hyponym : Scanner US Hardware
- Semantically related : Geiger counter US radioactivity

83 non preferred terms are associated with more than one preferred term in the thesaurus, ranging from 2 to 3 different preferred terms. This non unique association can be a source of problems when using the thesaurus' content as a basis for automatic indexing. If we select the wrong preferred term, we might for example suggest *petrol* (*aardolie*) as an indexing term for a document about food, because the non preferred term *oil* (*oliën*) has both *petrol* and *vegetable oil* (*plantaardige oliën*) as its preferred term. We will present in the next section our semi-automatic annotation pipeline, the ranking algorithm applied to the term extraction, CARROT, and its usefulness for selecting the right preferred term out of 2 to 3 different possibilities.

3.4 Semi-automatic annotation pipeline

3.4.1 The pipeline

As stated in section 2, the goal of the semi automatic annotation pipeline is to suggest appropriate indexing terms to cataloguers, with the goal of easing their job and increasing their productivity. From discussion with the cataloguers it followed that they like a focussed and limited set of keywords : focussed because they only experience a suggestion as supportive if it closely matches the main topic of the document, limited because actual work process of cataloguers only allows for a limited number of terms to be attached to a document. Another reason for that requirement is that the inspection of the suggested terms should improve the work process, so the inspection time and the mental processing of the suggestions need to be bounded in order not to generate additional burdens.

The pipeline consists of tree parts : a term detector, a term collector and a term ranker. As input to our pipeline we use our selected corpus and the GTAA. The output of the pipeline is a ranked

⁵ All the terms we mention in this paper are translated from Dutch to English out of consideration for our readers. We tried to select examples which have the same ambiguity in their semantics in the English translation

list of GTAA preferred terms.

3.4.2 The input : GTAA in a RDF-OWL representation

As input we use an RDF-OWL representation of the GTAA, based on the SKOS Working Draft (see (van Assem *et al.*, 2006) and (Miles & Brickley, 2005)). The SKOS representation of a thesaurus is “concept based” : instead of terms, the entities are nodes with identifiers (ID), to which labels are attached, a `prefLabel` to represent the preferred term, and one or more `altLabel(s)` to represent the non preferred term(s). As the GTAA entries are in plural form, we also extended this model to add the information of the singular form corresponding to the original thesaurus terms. This model has drawbacks, and has an obvious conceptual bias, but it helps gathering pragmatically different strings corresponding to the same annotation ID. These strings are called “textual representations of the concept” in the GATE pipeline, and we decided to keep this terminology here.

3.4.3 The term detector : GATE with the Apolda plug-in

The term detector scans a text and looks for all possible textual representations of concepts. The detector is built with the Apolda plug-in in GATE architecture (Maynard *et al.*, 2003). After tokenization, the Apolda plug-in makes a simple string matching. It annotates a piece of text with the ID of the “concept” corresponding to the longest matching textual representation. If for a piece of text multiple concepts have the same longest matching textual representation, which can be the case for a non preferred term with multiple preferred terms, the plug in generates all possible annotations. This means that the string `minority` will receive two annotations : `Keyword_ethnical_minority` and `Keyword_religious_minority`. The string `religious minority` however will only receive the latter. The term detector is not case sensitive.

3.4.4 The term collector

The outcome of the term detector is an annotated text. In this text, multiple annotations can correspond to the same “concept”. The term collector collects all the annotation ID’s, computes their number of occurrences and writes the output into one file.

3.5 The term ranker and WSD algorithm : CARROT

The file with ID’s and number of occurrences computed at the previous step is fed into the Cluster And Rank Related to Ontology and Thesauri algorithm (CARROT algorithm) (Gazendam *et al.*, 2006).

CARROT uses the fact that terms in the Keyword facet of the GTAA are related to others via the related term, broader term and narrower term relations. We hypothesise that terms which relate to a lot of the other terms found in the text can be semantically more representative of the core topics of the TV program than terms which are found more often but without any relations to others. If one of the thesaurus relationships exists between two of the found terms we say that a relation of distance 1 exists. We also check if an intermediate term connects two terms in the

GTAA. These connections via intermediate terms are defined as relations of distance 2. We do not make any distinction in the type of relationships.

To rank the extracted keywords, we use the following rules :

- Step 1. We select the keywords with both a distance 1 and a distance 2 relation. We then order these keywords based on their number of occurrences, putting the most frequent on top of the list.
- Step 2. We select the remaining keywords with a distance 2 relation to keywords found during Step 1. We order these keywords based on their number of occurrences and add them to the list.
- Step 3. We select the remaining keywords with a relation. We order these keywords based on their number of occurrences and add them to the list.
- Step 4. We order the remaining keywords based on their number of occurrences and add them to the list.

This algorithm creates clusters of ranked terms (several terms can have the same rank, they are then simply ordered alphabetically).

Our previous experiment in (Gazendam *et al.*, 2006) showed that only the top clusters provided relevant keywords, so we intend to present the cataloguers with only these top clusters by default, with the possibility to access the whole ranked list if they wish to. In this paper we propose this CARROT algorithm as a means for selecting the right preferred term (right interpretation) for a non preferred term with multiple preferred terms (an ambiguous word). For example, the text : " *Snacks do not contain a lot of minerals.*" contains the non preferred term *minerals* and the preferred term *snacks*. *minerals* has three preferred terms : *food*, *fertilizer* and *ore*. All are considered to occur once, because their common non preferred term occurs once. These three plus *snacks* are fed into CARROT. Due to the direct relation between the terms *food* and *snacks*, *food* now ranks higher than the other two preferred terms. This means that we here interpret *minerals* as referring to *food* in this case.

As the non preferred term attributes the same number of occurrences to all its preferred terms, three scenarios are possible :

- One of the preferred terms has more direct or indirect relations to other found terms and ranks higher as a result ;
- One of the preferred terms combines a higher number of occurrences due to the fact that the preferred term appeared itself in the text or one of its other non preferred terms appeared in the text ;
- The different preferred terms rank equally high.

The output of the pipeline is the same list of annotation ID's as the input, but ranked. Therefore, our hypothesis for Word Sense Disambiguation is that the irrelevant preferred terms will not be connected to any of the other found keywords, and thus will be ranked at the bottom of the list. As a consequence, they will not be shown to the cataloguers as indexing suggestion. We present the positioning of our experiment with the state of the art in Word Sense Disambiguation in the following section, followed by the experiment itself.

4 Related Work

The task we are interested in in this paper can be related to Word Sense Disambiguation. In (Ide & Véronis, 1998), the authors describe the typical two-step process for this task :

1. Define the set of senses per lexical unit ;
2. Use either a context-based method to determine which of the senses corresponds to the occurrence of the lexical unit considered, or an external knowledge source.

Many works mention the use of a dictionary as an external knowledge for that purpose ((Veronis & Ide, 1990), for example), whereas statistically-based or machine-learning methods advertise the corpus-based contextual approach (see for example (Yarowsky, 1995)). Of course, some mixed approaches exist, as (Stevenson & Wilks, 2001). In our use case, the set of senses to take into account is the set of possible preferred terms for each ambiguous non preferred term. The method that we experiment here is using external knowledge, but instead of the lexical content of dictionary definitions, or instead of trying to map the lexical environment of the external knowledge to the corpus content, we use the thesaurus independently, and take only into account the number of occurrences of each term as a contextual information. The selection of the relevant sense, *i.e.* of the relevant preferred term, is made only based on relationships crafted by hand by cataloguing experts when building the thesaurus. Therefore it is still different from (Yarowsky, 1992), who also based his Word Sense Disambiguation algorithm on a thesaurus.

5 The experiment

5.1 Experiment : selecting the right keyword when multiple USE relations are possible

For this experiment, we annotated our documents with all the possible preferred terms related to the non preferred terms we found in the texts, along with their number of occurrences, and we will check whether the algorithm designed for ranking the IE output will help us disambiguating between the different possibilities. We will evaluate whether CARROT

1. Ranks the relevant preferred term higher ;
2. Ranks the irrelevant preferred terms low enough for them not to be part of the keywords suggested to the cataloguers.

5.2 Material

We constructed our corpus from a set of over 500 catalogue descriptions from Sound and Vision, related to TV programs. Each of these catalogue descriptions contains specific fields, that are described in Dublin Core : *e.g.* maker, title and keywords. One of the fields is a free text description called summary. In the Keyword field the topic of the program is described by a limited set of preferred terms from GTAA's Keyword facet. From this set of catalogue descriptions we selected all files which :

1. contain a non preferred terms which has multiple preferred terms and
2. have one of its related preferred terms appear in the keyword field

Based on these requirements we selected automatically a corpus of 121 documents, of averagely 200 words each. The second requirement is related to evaluation purpose : the preferred term that was chosen to describe the document can be seen as the correct interpretation of the non preferred term present in the description text. We base ourselves on this assumption to evaluate

the results of our ranking algorithm : the preferred term present in the Keyword field should be ranked higher than the other possible preferred terms.

5.3 Experiment

We ran our pipeline on our corpus. After completion we looked at the non preferred term, the rank of all associated preferred terms in the ranked list and compared this ranked list with the preferred term in the Keyword field of the catalogue description. We have three possible outcomes of this comparison :

1. Correct suggestion : the suggested preferred term⁶ is the preferred term in the keywords
2. Wrong suggestion : the suggested preferred term is not the preferred term in the keywords
3. Undecidable : No suggestion is made because two (or all three) preferred terms rank equally high

When evaluating the results, we also came across a set of unusable data. We discuss this point in the following section. The results are shown in table 1

correct	undecidable	wrong	unusable data	total
43	26	2	50	121

TAB. 1 – Results

5.3.1 Discussion

One of the issues that arose when evaluating our results was that we still have numerous unusable documents in our corpus : it turned out that for some documents, the non preferred term is found in the keyword field. According to the production rules of Sound and Vision a non preferred term cannot be used in the keyword field, but the set of keywords changes over time : a preferred term may be ambiguous and as a consequence be changed to a non preferred term. Because we used old descriptions in our corpus, some of these contained previously preferred terms which now became non preferred ones in their keyword field. This is the case, for example, for **murder assault** (8 occurrences) and **tent kampen** (23 occurrences). These two examples account for two thirds of the unusable data. We excluded these from our analysis.

For the remainder of the corpus, in approximately 19 out of 20 (95%) cases, the suggestions are not incorrect. We found only two cases in which we gave a wrong suggestion. Both mistakes are with the same non preferred term **clubs** which has as preferred terms **hotel**, **restaurant** and **cafe** (HRC) and **association**. This word club was used in the context of football clubs. One text was on the share issue of soccer club Ajax. The other text was on the showing of a documentary on the soccer club Ajax in a theater. The term **club** had the meaning of **association** in both cases, referring to the soccer association. However the **hotel**, **restaurant** and **cafe** was suggested. In both cases terms at distance 2 from HRC were present in the text : **theater** via the intermediate term **nightlife** and **director** via the intermediate term **enterprice**. On the other hand **associations** did not have direct or distance 2 connections to other extracted terms in the football domain as **soccer**, **supporter**, **match**, **trainer** : the distance in the thesaurus between

⁶i.e. the preferred term with the highest rank in the list.

these terms and **association** was too big. In our corpus, we have two other instances of **club** for which the matching to its preferred terms is successful once and undecidable another time. Both these texts were also in the soccer domain and having the preferred term **association**.

This could suggest that we have one “preferential preferred term” in the corpus, and that this information could be used for solving in a light way the ambiguity problem. Unfortunately, this is not always the case : the non preferred term **windmills** occurs once as **wind turbine** and once as **mill** ; in both cases the correct suggestion is made by our system. Other non preferred terms with a bigger number of occurrences also have a non regular distribution of their preferred terms.

Another remarkable feature of the results is the big number of undecidable cases. The reason why we encounter this big number of undecidable cases is manyfold :

1. Our method uses general conditional rules. These conditions are not really specific : *having any distance 1 relation satisfies a condition*. As a result, in many cases both preferred terms fit the same conditions. This can be amended by sharpening these conditions, for example by counting the number of terms at distance 1 or distance 2.
2. The texts of our corpus are relatively small, so the number of found (and related) terms is also small, and the number of occurrences too low to disambiguate between different possibilities.
3. In many cases the different preferred terms have a distance 1 relation to other extracted terms, increasing the chance of a tie. At the same time this means that the difference in meaning between the preferred terms can be subtle, giving value to the undecidability. For example, it is very difficult and maybe not relevant to distinguish between the three preferred terms related to **toxin**, namely **poison**, **venom** and **dangerous substance**, in the context of a TV program about farmers getting ill after using a toxin as a form of herbicide.

The last remark that we can make is that, due to the small number of different keywords in the different texts, very few clusters were created. As a consequence, it was hardly ever the case that the non relevant preferred terms found place low enough in the ranked list not to be proposed for indexing suggestion. Therefore, we should modify our algorithm in order to make it take into account only the preferred term with higher rank, and remove the other related preferred terms from the suggestion list.

6 Conclusion and Perspectives

We investigated whether our method and the CARROT algorithm could be used for disambiguation in an indexing setting. In cases of ambiguity, it only gives suggestions for which preferred term to choose in two cases out of three, but when it gives a suggestion, it is correct so in approximately 19 out of 20 cases. The two bad suggestions came from the same thesaurus concept, and were due to its lack of structure. Using another external resource like the Princeton University’s WordNet thesaurus could help us cope with that problem. However, the interpretation of our success rate and percentage of undecidable cases must be subject of study : it is up to the cataloguers to determine whether these numbers are fair ⁷. This is the subject of another study, that we will also conduct in the course of our project.

⁷A success of 19 out of 20 seems quite reasonable in the perspective of IR publications, but when talking about automatically securing railway crossings, the same success ratio is considered really bad.

Acknowledgements

This research was partly supported by the NWO funded CATCH program, including the CHOICE project. We want to thank our colleagues, both at the University and at Sound and Vision for their daily help, support and our fruitful collaboration.

Références

- GAZENDAM L., MALAÏSÉ V., SCHREIBER G. & BRUGMAN H. (2006). Deriving semantic annotations of an audiovisual program from contextual texts. In *Proceedings of First International workshop on Semantic Web Annotations for Multimedia (SWAMM 2006)*.
- IDE N. & VÉRONIS J. (1998). Introduction to the special issue on word sense disambiguation : The state of the art. *Computational Linguistics*, **24**(1), 1–40.
- ISO (1986). *Documentation - guidelines for the establishment and development of monolingual thesauri*. International Organization for Standardization, ISO 2788-1986 edition.
- MAYNARD D., TABLAN V. & CUNNINGHAM H. (2003). Ne recognition without training data on a language you don't speak. *ACL Workshop on Multilingual and Mixed-language Named Entity Recognition : Combining Statistical and Symbolic Models*.
- MILES A. & BRICKLEY D. (2005). Skos core guide. 2nd W3C Public Working Draft.
- STEVENSON M. & WILKS Y. (2001). The interaction of knowledge sources in word sense disambiguation. *Computational Linguistics*, **27**(3), 321–349.
- VAN ASSEM M., MALAÏSÉ V., MILES A. & SCHREIBER G. (2006). A method to convert thesauri to skos. In *Proceedings of the Third European Semantic Web Conference (ESWC'06)*.
- VERONIS J. & IDE N. M. (1990). Word sense disambiguation with very large neural networks extracted from machine readable dictionaries. In *Proceedings of the 13th conference on Computational linguistics*, p. 389–394, Morristown, NJ, USA : Association for Computational Linguistics.
- YAROWSKY D. (1992). Word-sense disambiguation using statistical models of roget's categories trained on large corpora. In *Proceedings of 14th International Conference on Computational Linguistics (COLING-92)*.
- YAROWSKY D. (1995). Unsupervised word sense disambiguation rivaling supervised methods. In *Meeting of the Association for Computational Linguistics (ACL' 95)*.