

Comparing corpora and lexical ambiguity

Patrick Ruch
Medical Informatics Division
Geneva University Hospital
Switzerland
ruch@dim.hcuge.ch

Arnaud Gaudinat
LATL
University of Geneva
Switzerland
gaudinat@latl.unige.ch

Abstract

In this paper we compare two types of corpus, focusing on the lexical ambiguity of each of them. The first corpus consists mainly of newspaper articles and literature excerpts, while the second belongs to the medical domain. To conduct the study, we have used two different disambiguation tools. However, first of all, we must verify the performance of each system in its respective application domain. We then use these systems in order to assess and compare both the general ambiguity rate and the particularities of each domain. Quantitative results show that medical documents are lexically less ambiguous than unrestricted documents. Our conclusions show the importance of the application area in the design of NLP tools.

Introduction and background

Although some large-scale evaluations carried out on unrestricted texts (Hersh 1998a, Spark-Jones 1999), and even on medical documents (Hersh 1998b), conclude in a quite critical way about using NLP tools for information retrieval, we believe that such tools are likely to solve some lexical ambiguity issues. We also believe that some special settings -particular to the application area- must be taken into account while developing such NLP tools.

Let us recall two major problems while retrieving documents with NLP engines (Salton, 1988):

1-Expansion: the user is generally as interested in retrieving documents with exactly the same words, as in retrieving documents with semantically related words (synonyms, generics,

specifics...). Thus, a query based on the word *liver*, should be able to retrieve documents containing words such as *hepatic*. This expansion process is usually thesaurus-based. The thesaurus can be built manually or automatically (as, for example, in Nazarenko, 1997).

2-Disambiguation: a search based on tokens may retrieve irrelevant documents since tokens are often lexically ambiguous. Thus, *face* can refer to a body part, as a noun, or an action, as a verb.

Finally, this latter problem may be split into two sub problems. The disambiguation task can be based on parts-of-speech (POS) or word-sense (WS) information, but the chronological relation is still a discussion within the community. Although, the target of our work (Ruch and al., 1999, Bouillon and al., 2000) is a fine-grained semantic disambiguation of medical texts for IR purposes, we believe that the POS disambiguation is an important preliminary step. Therefore this paper focuses on POS tagging, and compares morpho-syntactic lexical ambiguities (MSLA) in medical texts to MSLA in unrestricted corpora.

Although the results of the study conform to preliminary naive expectations, the method is quite original¹. Most of the comparative studies, dedicated to corpora, have addressed the problem by applying metrics on words entities or word pieces (as in studies working with n-

¹ We do not claim to be pioneer in the domain, as others authors (Biber 1998, Folch and al., 2000) are exploring similar metrics. However, it is interesting to notice that for these authors the adaptation of the NLP tools has rarely been questioned in a technical point-of view, and in order to feed back the design of NLP systems.

gram strings), or on special sets of words (the indexing terms, see Salton, 1988) as in the space-vector model (see Kilgarriff, 1996, for a survey of these methods), whereas the present paper attempts to compare corpora at a morpho-syntactic (MS) level.

1 Validating each tagger into its respective domain

In order to conduct the comparative study, we used two different morphological analysers; each one has a specific lexicon tailored for its application field. The first system is specialised for tagging medical texts (Ruch and al., 2000), while the second is a general parser (based on FIPS, cf. Wehrli, 1992).

For comparing lexical ambiguities on a minimal common base, the output of each morphological analyser is first mapped into its respective tagset (more than 300 fine-grained tags for FIPSTAG, and about 80 for the morpheme-based medical tagger). The tagsets are then converted into a subset of the medical tagger. Finally, about 50 different items constitute this minimal common tagset (MCT), which will serve for comparing both corpora.

We collected two different sets of documents to be tagged at a lexical level via the predefined MCT: this step provides a set of tags to every token. This set of tags may come from the lexicon or from the POS guesser. As we are using guessers, the empty set (or the tag for unknown tokens) is forbidden. However, first of all, it is necessary to verify the lexical coverage of each system for each corpus, as we need to be sure that the lexical ambiguities provided by each system are necessary *and* sufficient.

The corpus of the unrestricted texts consists of 16003 tokens: about one third of newspaper articles (*Le Monde*), one third of literature excerpts (provided by the InaLF, <http://www.inalf.fr>), and a smaller third being mainly texts for children. Approximately a quarter (3987 tokens) of this corpus is used for evaluating FIPSTAG tagging results (the tool together with some explanations can be found at <http://latl.unige.ch>). In parallel, we chose three types of medical texts to make up the medical corpus: it represents 16024 tokens, with 3 equal

thirds: discharge summaries, surgical reports, and laboratory or test results (in this case, tables were removed). Again, a regularly distributed quarter (4016) of this corpus is used for assessing the medical tagger.

The test samples used for assessing the results of each tagger are annotated manually before measuring the performances, but in both cases we sometimes had to modify the word segmentation of the test samples. This is particularly true for FIPSTAG, which handles several acceptable but unusual collocations (which gather more than one 'word'), as for example *en avion* (in Eng. *by plane*), which is considered as one lexical item, tagged as an adverb. For the lexical tagger we had to modify the 'word' segmentation in the other direction (for tagging items smaller than 'words'), as morphemes were also tagged. Table 1 gives the results for FIPSTAG, and table 2 gives the results for the medical tagger. In the case of the medical tagger, together with the error rate and the success rate, we provide results of the residual ambiguity rate: the basic idea is that the system does not attempt to solve what it is not likely to solve well (cf. Ruch and al. 2000, a similar idea can be found in Silberztein 1997).

1 Correct tag	3959 (99.3%)
1 Incorrect tag	28 (0.7%)

Tab. 1: Evaluation of FIPSTAG

1 Correct tag	3962 (98.5%)
1 Incorrect tag	12 (0.4%)
2 or more tags, at least 1 is correct	39 (1.0%)
2 or more tags, 0 correct	3 (0.1%)

Tab. 2: Evaluation of the medical tagger

A comparison of the tagging scores (99.3 vs. 98.5) confirms that both systems behave in an equivalent way in their respective application area².

² Out of curiosity, we ran each tagger on a small sample of the other domain. The tests were made without any adaptation. FIPSTAG made 27 errors in a medical sample of 849 tokens, i.e. an error rate of 3.2%. The medical tagger made 18 errors in a general sample of 747 tokens, which means an error rate of 2.4%. In the case of the medical tagger, 41 tokens

2 Morphological analysers, lexicons and guessers

Lexical ambiguities have two origins: the lexicon, and the guessing stages for unknown tokens. However, all the ambiguities considered in this study are strictly lexical, and so translation phenomena (Tesnière 1959, and Paroubek 1997) are not considered here.

2.1 Medical lexicon

The medical lexicon is tailored to biomedical texts, thus, with about 20000 lexemes, it covers exhaustively ICD-10. The biomedical language is not only a 'big' sub language, as its morphology is also more complex. This high level of composition (at least compared to regular French or English languages) concerns about 10% of tokens within clinical patient records; therefore the lexicon contains also about 2000 affixes. For example, the token *iléojéjunostomie* is absent from the lexicon, however, this type of token may be recognized via its compounds (see Lovis and al., 1997, for the so-called morphosemantemes): *iléo*, *jéjuno*, and *stomie*.

2.2 Morphological analysis and medical morphology

The morphological analysis associates every surface form with a list of morpho-syntactic features. When the surface form is not found in the lexicon, it follows a two-step guessing process: the first level (oracle1) is a more complex morphological analyzer, based on the morphosemantemes, while the second level guesser (oracle2) attempts to provides a set of MS features looking at the longest ending (as described in Chanod and Tapanainen, 1995).

The importance of these two levels is not clear for POS tagging, but becomes manifest when dealing with sense tagging. Let us consider three examples of tokens absent from the lexicon: *allomorphiques*, *allomorphiquement* (equivalent to *allomorphic* and *allomorphically* in Eng.

remained ambiguous after disambiguation, the residual ambiguity is therefore about 5.5%. In this sample, and before disambiguation, the number of ambiguous tokens was 150, which means an ambiguity rate of 20%. Thus, even using the same lexicon, the ambiguity rate seem higher for general corpora than for domain-specific ones.

language) and *allocation*. In the first case, the prefix *allo* and the suffix *morphiques* are listed in the morphosemantemes database (MDB). In the second case, *morphiquement* is not listed within the MDB, but *ment* can be found in it, In these two cases, therefore, oracle1 is able to provide both the MS and the WS information associated. The latter example cannot be split into any morphemes, as *cution* is absent from the MDB. Thus, oracle1 is unable to recognize it, and finally oracle2 will be applied and will provide some MS features regarding exclusively the endings. The major role given to oracle1 and the semantic features it provides is obvious for IR purposes.

The final stage transforms some of the lexical features returned by the morphological analysis in a tag-like representation to be processed later by the tagger.

2.3 FIPSTAG tagger and lexicon

The FIPSTAG lexicon is a general French lexicon, therefore it contains most well-formed French words. The overall structure of the lexicon is more or less stable, but the content is regularly updated in order to improve the coverage. Currently, the coverage is about 200000 words with around 30000 lexical items. The lexicon is designed for deep parsing, so that, together with classical morpho-syntactic features, we can also find sub categorization of verbs, semantic features, and some very specific grammatical classes.

As the system is claimed to be general, it is supposed to master efficiently any unknown words: the lexical modules supply, in an equiprobable way, all the possible lexical categories (i.e. nouns, verbs, adjectives, and adverbs), as other categories are considered to be exhaustively listed in the lexicon. Consequently, the guesser does not rely on any morphological information, and only syntactic principles are applied to choose the relevant features.

3. Results and comparison of ambiguities

	medical corpus	general corpus
ambiguities	2532 (15.8%)	4657 (29.1%)

Table 3: ambiguity rates according to the corpus

Amb. class	Si.	Fm.	Fg.	Ex. or BR
proc/v[ms]	0	0	1	lui
nc[ms]/v[n]	0	0	1.3	être
d[fs]/nc[fs]	0	0	2.3	une
v[12]/v[s03]	0.2	1.3	7	semble,
sp/v[12]/v[s0]	0.2	0.2	1	entre, contre
prop[03]/cccs	0.2	0.3	1.7	s'
nc[ms]/v[12] /v[s03]	0.3	0.4	1.3	contrôle, groupe
r/v[12]/v[mp]	0.8	1	1.3	plus
d[ms]/nc[ms]	0.8	1.6	2	son
d[bp]/proc	0.8	5.5	7	les
d[ms]/proc	0.9	7.1	8.3	le
cccs/nc[ms]/r	1	1	1	bien
nc[ms]/v[s03]	1	1	1	fait
proc/prop[12]	1	1.7	1.6	nous
cccs/r	1	2.1	2.2	que
nc[ms]/r	1.1	4.9	4.6	pas
nc[ms]/v[s03]	1.2	5.3	4.5	est
nc[fs]/v[12] /v[s03]	1.3	2.6	2	sorte, mesure, demande
proc/sp/cccs	1.6	7.5	4.6	en
d[bs]/proc	1.9	13.8	7.3	l'
d[fs]/proc	2.1	14.1	6.8	la
a/nc	4.2	1.7	0.4	patient
a/nc/v ³	5.0	1.5	0.3	patiente

Tab. 4: Similarity measure for the most frequent classes of ambiguity.

Note (tab. 4):

Column 1 gives the ambiguity class. Column 2 provides the ratio of similarity (maximum similarity

³ This class has only one representative within the medical corpus, the word *patient* (feminine *patiente*): An equivalent within the general corpus is *politique* (in Eng. it means both *political* and *politics*), but the former (0.5% of tokens) is ten times more frequent than the latter (0.05%). The frequency of the word *politique* is consistent with the frequency lists distributed by Jean Véronis (<http://www.up.univ-mrs.fr/~veronis>), which were calculated on a one million words corpus from *Le Monde Diplomatique* (1987-1997). It should be noted that this result questions the concepts of 'unrestricted corpora' and 'representativeness' (Biber, 1994), as in fact it often refers to a mix of politics and newspaper topics!

= 1, minimal similarity = 0 and 5) between the frequency of the considered ambiguity in medical (Fm.) and general texts (Fg.). Columns 3 and 4 (resp. Fm. et Fg.) indicate the frequency of the ambiguity respectively in the medical texts and in the general texts. Column 5 provides some examples or the best representative (BR) of the ambiguity class, i.e. when one lexeme represents at least 80% of the class.

List of abbreviations for the syntactic categories: *proc*, clitic pronoun; *v*, verb; *nc*, common noun; *d*, determiner; *sp*, preposition; *prop*, personal pronoun; *cccs*, conjunction; *q*, numeral. List of abbreviations for the morpho-syntactic features and sub categorizations: *ms*, masculine singular; *n*, verbal infinitive form; *fs*, feminine singular; *bs*, masculine or feminine singular; *12*, first and second person singular or plural; *s03*, third person singular; *p03*, plural third person.

When possible this tagset follows the MULTEXT (Ide and Véronis 1994) morpho-syntactic description, modified within the GRACE action. But we must notice that the original MULTEXT description and the GRACE version (Paroubek and al. 1998, Rajman and al. 1997) for the French language have not been foreseen for annotating morphemes.

Previously, while attempting to assess the performance of our tools, only a sample of the ad hoc corpus we built up was used, whereas the following studies on the ambiguities will be carried out on the whole corpus. Like in the validation task, the lexical ambiguities are based on the morphological analysis of each tagger, expressed in the MCT. First of all, table 3 gives the general ambiguity rate in each corpora: it clearly states that the total ambiguity rate in general corpora is about twice as big as in medical texts.

A more precise table (tab. 4) provides at least two remarkable results. First, it shows that in the general corpus, less than a dozen words are responsible for half of the global ambiguity rate. These results must be compared to (Chanod and Tapainen, 1995), who situate this number around 16, while about six words generate the same level of ambiguity in the medical corpus! This table also shows that the distribution of the ambiguity type is also domain dependant. Thus, the ambiguity *d[fs]-[bs]/proc* is twice more frequent in medical texts, and the ambiguity represented by the tokens *patient/patiente* (masculine and feminine form of *patient*; which

may be a noun, an adjective, or some form of verb) is five times more frequent. On the contrary, some classes of ambiguity are simply absent or very rare in the medical domain (as for example v[12]/v[s03], or nc[ms]/v[n]).

Finally, in table 5, we give the distribution of the most frequent syntactic categories according to the corpus. In this table, a particularly interesting result concerns the imbalance between categories of noun phrases (determiner, noun, adjective...) and categories of verb phrases (verb, adverb...); the former being much more frequent in medical texts, whereas the latter are more frequent in general texts. Here we verify a well-known stylistic manner: medical reports are often written in a telegraphic style, where the verb is frequently implicit. As a corollary, nominalization phenomena are very frequent. Simple or complex numeral tokens (date, time, expressions with digits and measure symbols) are also much more frequent.

General		Medical	
r	505	276	v[n]
v[n]	721	301	v[12]; v[s03]; v[p03]
cccs	765	550	q
v[12]; v[s03]; v[p03]	837	587	cccs
sp	1356	1283	a
d	1659	1529	f
nc	1707	1784	d
f	2179	2138	sp
-	-	3472	nc

Tab. 5: Distribution of the most frequent morpho-syntactic categories according to the domain.

Note (tab. 5) : f refers to the punctuations.

4. Discussion and conclusion

We have showed that the lexical ambiguity in medical texts (considered as a paradigm of any particular domain) is different to the one in general texts, both at a purely quantitative level, and at a deeper qualitative level. Another result concerns the difference in the distribution of the POS categories. All these particularities must be added to others: lexical, morphological, spelling and grammar errors. This last point has been

rarely studied, but errors in documents, which are not intended for publication, may be quite impressive (the spelling error rate in our medical corpus was about 2%, i.e. up to one error every five sentences!). Finally, our conclusion is of two types: First, concerning the study, we showed that the use and comparison of taggers tailored for different corpora, supports a measure of the difference between these corpora; second, at a more methodological level, if it seems that the syntax may be *ceteris paribus*- regarded as a domain-independent field (at least at a computational level, cf. Wehrli 1995), we argued that natural language processing applications require domain-adaptable tools. Therefore, the use of NLP tools by other research fields must be very carefully related to the design of these tools. We suggest that adaptability should be explored in at least three directions⁴:

1. Systems must allow lexical items to be added (custom lexicon) and removed from the lexicon; therefore access to the main lexicon must be available – at least negatively.
2. Systems must be optionally applied with a specialised morphological analyser module.
3. MS description (tagset) should be parametrable, and this should include the ability to provide a mapping table.

Acknowledgements

This study was supported by the FNRS (Swiss National Science Foundation).

References

- Douglas Biber (1994) *Representativeness in Corpus Design*. Zampolli, Antonio, Nicoletta Calzolari and Martha Palmer (Eds.). 377-407.
- Douglas Biber, Susan Conrad, Randi Reppen (1998) *Corpus Linguistics: Investigating Language Structure and Use*. Cambridge University Press.
- Pierrette Bouillon, Patrick Ruch, Robert Baud, Gilbert Robert (2000) *Indexing by statistical tagging*. In proceedings of the 7th JADT'2000. Vol. 1, pp. 35-42. Lausanne. Switzerland.

⁴ A future version of FIPSTAG should integrate some of these specifications.

- Jean-Pierre Chanod, Pasi Tapanainen (1995) *Tagging French: comparing a statistical and a constraint-based method*. In ACL, Ed., 7th Conference of the European Chapter of the Association for Computational Linguistics (EACL'95). pp. 149-156. Dublin.
- Helka Folch, Serge Heiden, Benoît Habert, Serge Fleury, Gabriel Illouz, Pierre Lafon, Julien Nioche, Sophie Prévost (2000) *TypTex: Inductive Typological Text Classification by Multivariate Statistical Analysis for NLP Systems Tuning/Evaluation*. In Proceedings of the 3rd International Conference on Language Resources and Evaluation (LREC'2000), Athenes, Greece.
- William R. Hersh, Price S., Kraemer D., Chan B., Sacherek L, Olson D (1998a) *A Large-Scale Comparison of Boolean vs. Natural Language Searching for the TREC-7 Interactive Track*. TREC 1998, pp. 429-438.
- William R. Hersh (1998b) *Information Retrieval at the MILLENIUM*. In R MASY, Ed., American Medical Informatics Association Annual Symposium (AMIA'1998, ex-SCAMC). Orlando.
- Nancy Ide and Jean Véronis (1994) *MULTEXT: Multilingual Text Tools and Corpora*. In Proceedings of the 15th International Conference on Computational Linguistics (COLING-94), Kyoto, Japan.
- Adam Kilgariff (1996) *Which words are particularly characteristic of a text? A survey of statistical approaches*. ITRI Technical report 96-08. (<http://www.itri.brighton.ac.uk/~Adam.Kilgariff/publications.html>)
- Christian Lovis, Robert Baud, Pierre-André Michel, Jean-Raoul Scherrer (1997) *Morphosemantems decomposition and semantic representation to allow fast and efficient natural language recognition of medical expressions*. In R MASY, ed., American Medical Informatics Association Annual Symposium (AMIA'1997, ex-SCAMC). Washington.
- A. Nazarenko, Pierre Zweigenbaum, Jean Bouaud (1997) *Corpus-based identification and Refinement of Semantic Classes*. In R MASY, ed., American Medical Informatics Association Annual Symposium (AMIA'1997, ex-SCAMC), pp. 585-589. Washington.
- Patrick Paroubek, Gilles Adda, Joseph Mariani, Josette Lecomte, Martin Rajman (1998) *The GRACE French Part-Of-Speech Tagging Evaluation Task*, In Proceedings of the 1st International Conference on Language Resources and Evaluation (LREC), Granada, Spain.
- Martin Rajman, Patrick Paroubek, Josette Lecomte (1996) *Format de description lexicale pour le français - partie 2: Description morpho-syntaxique, rapport GRACE GTR-3-2-1*. (<http://www.limsi.fr/TLP/grace/www/gracdoc.html>)
- Patrick Ruch, Pierrette Bouillon, Gilbert Robert, Robert Baud (2000) *Minimal Commitment and Full Lexical Disambiguation: Balancing Rules and Hidden Markov Models*. In Proceedings of the 5th CoNLL Conference (ACL-SIGNLL). Lisbon. Portugal.
- Patrick Ruch, Judith Wagner, Pierrette Bouillon, Robert Baud (1999) *Tag-like semantics for medical document indexing*. In N. M. LORENZI, ed., American Medical Informatics Association Annual Symposium (AMIA'1999, ex-SCAMC), pp. 137-141. Washington.
- Gerald Salton (1988) *Term-weighting approaches in automatic text retrieval*. McGraw.Hill. Vol. 24. New-York.
- Max Silberstein (1997) *The Lexical Analysis of Natural Languages*, In Finite-state Language Processing, Yves Shabes and Emmanuel Roche ed., MIT Press, pp. 175-203. Cambridge.
- Karen Spark-Jones (1999) *What Is The Role for NLP in Text Retrieval*. Strzalkowski, ed., Natural Language Information Retrieval, Kluwer Publishers, pp.1-25
- Lucien Tesnière (1959) *Elements de syntaxe structurale*. Klincksieck. Paris.
- Eric Wehrli (1992) *The Interactive Parsing System*, In ACL, ed., *Proceedings of COLING-92*. 870-4. Nantes. France.
- Eric Wehrli, Robin Clark (1995) *Natural Language Processing: Lexicon and Semantics*, Methods of Information in Medicine, Vol. 34, p. 68-74.