

From the texts to the contexts they contain: a chain of linguistic treatments

Ahmed Amrani* Jérôme Azé† Thomas Heitz† Yves Kodratoff†
Mathieu Roche†

Abstract

The text-mining system we are building deals with the specific problem of identifying the instances of relevant concepts present in the texts. Therefore, our system relies on interaction between a field expert and the various linguistic modules we use, often adapted from existing ones, such as Brill's tagger or CMU's Link parser. We have developed learning procedures adapted to various steps of the linguistic treatment, mainly for grammatical tagging, terminology, and concept learning.

Our interaction with the expert differs from classical supervised learning, in that the expert is not simply a resource who is only able to provide examples, and unable to provide the formalized knowledge underlying these examples. We are developing specific programming languages which enable the field expert to intervene directly in some of the linguistic tasks.

Our approach is thus devoted to helping one expert in one field to detect the concepts relevant for his/her field, using a large amount of texts. Our approach is made off two steps. The first one is an automatic approach that find relevant and novel sentences in the texts. The second one is based on the expert's knowledge and find more specific relevant sentences.

Working on 50 different domains without an expert has been a challenge in itself, and explains our relatively poor results for the first task.

1 Introduction

In this paper, we present our approach and experiments relative to the Novelty Track of TREC-2004. We only answer the first two tasks of the Novelty Track. The data available for this Track was made of two sets of newspaper with domains of **Opinion** and **Event**. For each of these domains, the set of papers dealt with 25 different topics, each containing at least 25 documents relevant to the topic. We have thus to deal with 50 different topics.

The first task was made of two sub-tasks:

- determine the relevant sentences in the documents of each topic
- find the new sentences among the relevant sentences previously determined

The second task we participated into asked to find the novel sentences given all the relevant sentences for each of the 50 topics.

In order to answer these two tasks, we used the text mining system that our team is developing. This system is a complete chain of treatment of texts starting with an initial corpus and ending with ontologies built from the corpus.

As this chain of treatment has been used for answering the two tasks we participated in, we will first present this chain and we will detail thereafter our specific approach for the two tasks.

2 The global chain of text mining

The whole chain we deal with is as follow: **text retrieval** → **standardization** → **tagging** → **terminology extraction** → **coreferences resolution** → **concept recognition**

*ESIEA Recherche, 9 rue Vésale - 75005 Paris - France - amrani@esiea.fr

†LRI - Université Paris-Sud - 91405 Orsay Cedex - France - {aze,heitz,yk,roche}@lri.fr

2.1 Text retrieval

The first step of any chain of treatment of texts is to build a homogeneous corpus relative to a given domain. In the TREC competition, we don't have to deal with this step as the corpus was given.

The corpus is supposed to be relevant for the domain studied. One particular point for TREC Novelty is that we know those irrelevant documents have been inserted in the data relative to some topics. The topics that contains irrelevant documents can easily be identified since the set of documents must contain more than 25 documents.

Although if we can identify these topics, we cannot identify the irrelevant documents, thus we decided to ignore this default.

2.2 Standardization

Standardization is also often called text cleaning. This is obviously specific to each domain, and only the expert is able to decide what is to be cleaned or not. For the TREC Novelty, we first built a set of abbreviations. This set is useful to develop abbreviations find in texts in order to reduce the ambiguity of texts. For instance, if a text contains the words 'United States', 'United States of America' and the abbreviations 'USA' and 'U.S.', we must develop all the abbreviations to see that they refer to the same concept: 'United States of America'.

We used external resources to build a lexicon such like WordNet 2.0[10], the Moby Part of Speech II and the Webster dictionary 1913. The lexicon is about 300,000 words and has been used to do a first grammatical tagging.

2.3 Tagging of texts

The next step concerned the grammatical tagging of each word of the texts. We have choose a rule based tagger, the Brill's tagger [2], to realize the tagging. Although the Brill's tagger has learned its tagging rules from a WallStreet newspapers corpus, most of the texts need a serious revisiting of the rules and lexicon provided by Brill's basic version.

In order to do that, we can first modify by hand the tags of uncorrectly tagged words and train Brill's tagger on the modified texts. This approach can be very time expansive because we can not ensure that the correct tagging will be achieve by Brill's tagger.

The approach we used is based on two external ressources, a tagging language that allows a human expert to write basic tagging rules but also complex tagging rules that can not be learn by Brill's tagger. The other resource is a software, ETIQ [1], that allows the expert to see the context of each word and then determine if the tag associated to a particular word is good or not.

For the Novelty Track, our expert have produced a set of 450 rules in order to improve the quality of the tagging. Some of these rules are specific to the corpus studied, but many of them are complex rules that are valid for any english text.

The language we developed was customized to be more powerful than the Brill's tagger rules. In the Brill's tagger, the rules context varies from -3 to +3 around the word to be tagged. The analyze of the bad tagged words reveals that this context is too poor to express the correct tagging rules, so we introduce in our language the possibility to use the whole sentence as context.

2.4 Terminology

Once a good tagging is obtained, we can go to the step of finding terms often used in the texts. The tagging is a really important step as we used grammatical patterns to extract collocations (bags of words) and we need to have a good tagging if we want to obtain significant collocations. We used simple patterns to extract collocations. These patterns are Adjective-Noun, Noun-Noun, Noun-Preposition-Noun, etc. Significant collocations are collocations that represent an occurrence of a concept in a text. Collocations can be order using statistical criteria [3, 13, 5, 11]. We used the software EXIT [12] built in our team. For the TREC Novelty Track, we extracted terms linked to Opinions or Event. These terms have been associated to concepts using ACT [8, 11] (see section 2.6).

We also extracted collocations from the description of topics. Then, using the software FASTR [7], we looked for variations of terms in the texts associated to each topic. For example, in the description of the

topic N88, the collocation “Vieques Island” have been extracted. This collocation appears only 7 times in the texts associated to this topic. But it appears 20 times in another form, which is a variation of “Vieques Island”: “Island of Vieques”.

Terms and their variants occurring both in the topic’s description and in the associated texts have been used to compute a score for each sentence of the texts.

2.5 Co-references resolution

One of the main difficulty we found in the TREC texts is the coreference problem. In other words, how can we automatically determined the person who is referred by the word “he” or “she” in a given sentence ?

The resolution of the coreferences on the names is performed topic by topic using predefined lexical patterns, list of possible first names and list of possible social roles. For each person, we look for the functions which can be associated with this person (using the social roles list). We also look for the possible pseudonyms and the gender for each individual.

This knowledge is used to replace, as much as possible, all the references of the type: he, she, his, her and I by the referred person.

For example, in the topic N51 which title is “General Pinochet Arrested”, many possible co-references have been detected in the sentences presented in table 1. These co-references are of two types: “he/his” or a social role (Foreign_Minister).

Table 2 shows how the different co-references have been solved.

sentence number	text
14-NYT19981017.0086	British authorities refused to say where Pinochet is being held, nor did they set a date for when he would be questioned.
15-NYT19981017.0086	ever since he led a violent coup to overthrow Salvador-Allende-Gossens, the elected socialist president in 1973, Pinochet has been a political icon throughout Latin-America , representing the excesses of a long period of military rule.
16-NYT19981017.0086	an estimated 3000 Chileans were shot in the-streets or ”disappeared ” during his rule, and a senior member of his regime was imprisoned under United-States pressure for the murder of former Foreign-Minister Orlando-Letelier in Washington in 1976.

Table 1: Some sentences from NYT19981017.0086 text.

sentence number	text
14-NYT19981017.0086	British authorities refused to say where augusto-pinochet is being held, nor did they set a date for when augusto-pinochet would be questioned.
15-NYT19981017.0086	ever since augusto-pinochet led a violent coup to overthrow salvador-gossens , the elected socialist president in 1973, augusto-pinochet has been a political icon throughout Latin-America , representing the excesses of a long period of military rule.
16-NYT19981017.0086	an estimated 3000 Chileans were shot in the-streets or ”disappeared ” during augusto-pinochet rule, and a senior member of augusto-pinochet regime was imprisoned under United-States pressure for the murder of former Foreign-Minister(<i>orlando-letelier</i>) orlando-letelier in Washington in 1976.

Table 2: Coreferences solved on sentences 14, 15 and 16 from NYT19981017.0086 text.

These ontologies informations were gather to achieve the co-reference resolution (see Tab 3).

2.6 Concept recognition in texts

Once terminology is completed, we use it to recognize the occurrence of a concept in a text, by clustering the terms into classes, each term being an instance of a concept. We are building a system of concept recognition in texts called ACT [8, 11].

```
<person>
  <id>augusto-pinochet</id>
  <lastName>pinochet</lastName>
  <firstName>augusto</firstName>
  <gender>1</gender>
  <pseudo>pinochet</pseudo>
  <socialRole>General</socialRole>
  <socialRole>Mister</socialRole>
  <socialRole>President</socialRole>
  <socialRole>Senator</socialRole>
  <socialRole>dictator</socialRole>
  <socialRole>leader</socialRole>
  <socialRole>president</socialRole>
</person>
<person>
  <id>orlando-letelier</id>
  <lastName>letelier</lastName>
  <firstName>orlando</firstName>
  <gender>1</gender>
  <socialRole>Foreign-Minister</socialRole>
</person>
```

Table 3: Part of the ontologies built from topic N51.

Instance of the concept : legal-charges-crimes
(charge:Verb,for:Preposition,death:Object)
(suspect:Verb,in:Preposition,kill:Object)
(white-man:Object,accuse:Verbe)
(accuse:Verb,of:Preposition,kill:Object)
Instance of the concept : aircraft-accident
(crash:Sujet,kill:Verbe)
(jet:Sujet,crash:Verbe)
(plane:Sujet,crash:Verbe)
(kill:Verbe,in:Preposition,crash:Objet)

Table 4: Examples of instances of two concepts.

ACT allows the user to build ontologies from the texts using the terms, the words, their tags and the context of each word or term.

In the present state of ACT, we use two types of information in order to spot the presence of a concept in the texts. As in [6], we use a superficial syntactic parser in order to obtain syntactic relationships among the words. In fact, most of them can be looked upon as terms. For instance, consider the following set of syntactic relationships we actually used in order to recognize concepts (see for example table 4).

In each case, you can notice that the grammatical relationship is not as important as the co-occurrence of the words in order to define a concept. This is why we are presently developing an extension of our terminology programs to apply them to more relations including the verb-noun relationships.

2.6.1 Concept definition

The step of concept definition is entire in hands of the domain expert. It is important to stress that this entire process is strictly domain specific. The expert has to define what the interesting concepts are. One of the main difficulty of the Novelty TREC Track was that we have to deal with 50 different topics for which we don't have any expert.

As our results will show, our tools are not really dedicated to this task when used by a non-expert. But it is quite important to notice that even if we were not expert of the different fields, we have successfully used ACT to define concepts that help us in the relevant and novelty task.

2.6.2 Inductive step

The last step in this process of concept characterization, and the one which will receive the largest attention in the future, is the inductive phase during which an existing categorization is automatically completed. The present induction algorithm is based on the determination of a few thousands of groups, called **seeds**¹ containing words having the largest possible number of syntactic relations in common. These seeds grow by adding new syntactic relations to them, when this does not result in a too great reduction of the number of relations in common. The notion of the "reduction of the number" is fixed by parameters chosen by the user. We experimentally fixed the parameters so that the algorithm proposes about 50000 concepts. These are compared to the ones obtained by hand. When the induced concepts contain at least one relation in common with the concepts obtained by hand, and none of the relations in common are different from those of concepts obtained by hand, then the relations that are not in common and their associated concepts are added to the categorization.

3 Task 1 of Novelty Track

We used two different approaches to answer this task: an automatic approach and a semi-automatic approach. The automatic approach used all the steps of our chain of treatments except concepts. The semi-automatic used all the informations including concepts.

¹In other words, we use a parallel version of Diday's algorithm[4] called "nuées dynamiques", or Michalski's "AQ" algorithm [9].

Run 1		Run 2	
Information	Weights	Information	Weights
$core.f_1$ from topic	1	$core.f_1$ from topic	1
$core.f_1$ from texts	0	$core.f_1$ from texts	0.1
$core.f_2$ from topic	1	$core.f_2$ from topic	1
$core.f_2$ from texts	0	$core.f_2$ from texts	0.1
terms	1	terms	1
nouns	1	nouns	1
locate places	1	locate place	1
verbs	0.1	verbs	0.1
numbers	0.4	numbers	0.4

Table 5: Weights used for Runs 1 and 2.

3.1 Relevant Sentences Retrieval

3.1.1 Automatic approach

For this approach, we have proposed two runs using different informations. The first run (Run 1) only used informations identify in the topic’s definition: persons, noun, verb and numbers. The second run (Run 2) used these same informations and all persons identified in the sentences of the texts.

For the person name, we make difference between co-references “he, his, she, her” and “I”. We identify the first co-references as $core.f_1$ and the second as $core.f_2$. We have included all the named entities in $core.f_1$.

The terminology we built and the references and co-references to individuals are used to compute a value of each sentence relevance. The computation is done as follows: we gather

- the individuals present in the topic and/or the texts;
- the terms and their various forms;
- the locations names present in the subject;
- the numbers (including the dates) present in the topic;
- the verbs and nouns present in the topic.

Each sentence is replaced by a summary containing only this information.

If a sentence contains none of the above information, its reduced form is empty - the sentence obtains a score of zero relevance. If not, the sum of the weights associated to each of these information provides a relevant score for the whole sentence. The last step consists in selecting among the whole set of sentences the relevant ones. The rule used considers sentences relevant if their scores are strictly higher than the average scores.

Table 5 shows the weights used for Run 1 and 2.

3.1.2 Semi-automatic approach

In this approach we used the concepts built by the expert with ACT. These concepts are combined with the informations gather by the automatic step. We did three runs (Run 3, 4 and 5) that used this approach.

We named P_a the sentences determined as relevant using the automatic approach, and P_c the sentences satisfying at least one rule in the set of concept’s based rules that have been determined by the expert. These rules allow us to determine if a sentence if relevant or not.

Table 6 shows the different weights used to combine P_a and P_c sentences.

	Relative Weights of	
	P_a	P_c
Run 3	0.1	0.9
Run 4	0.5	0.5
Run 5	0.9	0.1

Table 6: Weights used to combine automatic and concept's based approach.

topic	Ranks				
	Run 1	Run 2	Run 3	Run 4	Run 5
N51 (Event)	12	3	1	2	7
N54 (Event)	7	4	60	58	57
N55 (Event)	5	4	1	2	14
N61 (Opinion)	7	6	1	3	16
N70 (Opinion)	3	1	60	11	4
N76 (Opinion)	26	1	40	41	37
N78 (Opinion)	34	3	60	20	28
N82 (Event)	16	1	59	60	9
N85 (Event)	4	3	10	2	1
N94 (Opinion)	55	55	1	2	57
N95 (Event)	26	7	1	2	10
N96 (Opinion)	37	4	58	17	29

Table 7: Best results achieve for the relevant sentences retrieval task.

3.1.3 Results

We used the results given by TREC to evaluate our different approaches. Table 7 shows our best results (at least one run in the 5 first runs) for the relevant sentence retrieval task.

Given our results, it seems quite difficult to conclude. For this task, our better results are not link to a particular domain (Event/Opinion). As a global conclusion, we can see that Run 2 gives better results than Run 1. We can then conclude that it seems useful to take into account the persons identified in the sentences of the texts even if their don't appear in the definition of what is relevant for the topic. Unfortunately, we have not used this run as basic run for the semi-automatic approach. Results obtained for Runs 3, 4 and 5 have to be compare with those obtain with Run 1.

We can see that the use of concepts gives better results when the concepts were well defined (see topics N51, N55, N61, N94 and N95). Inversely, when the concepts were not useful, the semi-automatic approach gives worse results than the automatic one (see topics N54, N76, N78 and N96).

As we have already mention it, we don't have an expert for the topics under study and it seems interesting to notify that using our tool (ACT) to quickly analyse the texts, we found some useful concepts for the relevant task.

Table 8 shows the average Fscore obtain for each run on this task.

Run	average Fscore
Run 1	0.306
Run 2	0.356
Run 3	0.255
Run 4	0.299
Run 5	0.302

Table 8: Average Fscore associated to each run for the relevant task.

topic	Ranks				
	Run 1	Run 2	Run 3	Run 4	Run 5
N51 (Event)	4	12	1	2	7
N55 (Event)	6	2	1	3	53
N57 (Event)	2	3	1	3	3
N60 (Opinion)	59	50	1	26	60
N61 (Opinion)	9	5	1	2	2
N62 (Opinion)	17	4	54	45	1
N63 (Opinion)	54	4	10	12	60
N70 (Opinion)	1	2	29	6	3
N76 (Opinion)	59	20	2	1	59
N77 (Opinion)	1	3	57	29	2
N84 (Opinion)	1	6	58	7	3
N91 (Opinion)	54	2	46	33	8

Table 9: Best results for the Novelty detection (task 1).

3.2 Task 1: Novelty Detection

For this task, we simply used the summary of each sentence to determine if a sentence was novel or not. We consider that a sentence was novel if it contains at least one information that have not been previously seen. As for the previous task, we have two automatic runs: Runs 1 and 2, and three semi-automatic runs: Runs 3, 4 and 5. In the two first runs, sentences can not contain concepts and in the other ones, sentences can contain such informations.

3.3 Results

Table 9 shows our best results for this task.

Quite surprisingly, some of these results were obtain for topics in which we were not able to achieve good results for the relevant task. For instance, for topics N57, N62, N63, N77, N84 and N91, our automatic approaches didn't appear in the five first results for the relevant task but are in the five first for the novelty one.

More generally, we can see that the use of concepts improve the results for many runs (see topics N51, N55, N57, N60, N61, N63, N76 and N91).

Finally, we can see that our approaches perform better on Opinion topics than on Event ones.

Table 10 shows the average Fscore obtained for each run on this task.

Run	average Fscore
Run 1	0.066
Run 2	0.108
Run 3	0.098
Run 4	0.098
Run 5	0.072

Table 10: Average Fscore associated to each run for the novelty sub-task.

4 Task 2 of Novelty Track

For this second task, we used the given relevant sentences for each topic and we changed our novelty detection system.

First of all, we applied some filters on each sentence. The following informations were removed from the sentences.

	no concept	using concepts
using co-references resolution	Run 1	Run 3
without co-references resolution	Run 2	Run 4

Table 11: Combination used for the different runs.

topic	Ranks				
	Run 1	Run 2	Run 3	Run 4	Run 5
N56 (Event)	2	2	39	39	41
N65 (Opinion)	14	14	4	9	25
N70 (Opinion)	3	3	5	5	2
N73 (Event)	2	1	11	10	7
N75 (Opinion)	2	2	5	5	8
N86 (Opinion)	12	12	2	2	2
N88 (Event)	3	2	25	21	23
N89 (Opinion)	2	2	38	38	30
N91 (Opinion)	4	4	36	37	31
N93 (Opinion)	3	3	35	35	11
N96 (Opinion)	3	3	6	6	22
N99 (Opinion)	3	3	6	5	23
N100 (Opinion)	18	18	8	20	2

Table 12: Best results for the Novelty detection task.

- punctuation character
- all words except non modal verbs, nouns and adjectives
- and words: **be have get do can must should would could need shall will say says said relevant irrelevant**

For each filtered sentence, we have computed the $TF \times IDF$ score and we have normalized it by the number of words of the sentence. This information ($TF \times IDF$) was the basic information used to determine if a sentence was novel or not.

To determine the novelty, we combined two different strategies: the use of a static threshold and the use of a dynamic threshold. Thus, a sentence is novel if

- its score is higher than the static threshold.
- its score is higher than the one of the previous novel sentence: dynamic threshold

Two different approaches have been used to summarize the sentences.

- using or not the co-references resolution step
- using or not the concepts.

Table 11 shows the different combination we used. Run 5 is the same as Run 4, but we changed the static threshold in order to be more permissive.

4.1 Results

Table 12 shows our best results for this task.

In this task, our approach reach the first place only one time. But, surprisingly the automatic runs provide the best results. For this task, our results seem to be stable.

As for the first novelty detection task, the results are better for the texts relative to Opinion. Finally, it seems the resolution of the coreferences doesn't affect the global results.

Table 13 shows the average Fscore obtained for each run on this task.

Run	average Fscore
Run 1	0.614
Run 2	0.614
Run 3	0.598
Run 4	0.597
Run 5	0.602

Table 13: Average Fscore associated to each run for the novelty task.

5 Conclusions

The text-mining system we are building deals with the specific problem of identifying the instances of relevant concepts found in the texts. This has several consequences. We develop a chain of linguistic treatment such that the n -th module improves the semantic tagging of the $(n-1)$ -th. This chain has to be friendly towards at least two kinds of experts: a linguistic expert, especially for the modules dealing mostly with linguistic problems (such as correcting wrong grammatical tagging), and a field expert for the modules dealing mostly with the meaning of group of words. Our definition of friendliness includes also developing learning procedures adapted to various steps of the linguistic treatment, mainly for grammatical tagging, terminology, and concept learning. In our view, concept learning requires a special learning procedure that we called Extensional Induction. Our interaction with the expert differs from classical supervised learning, in that the expert is not simply a resource who is only able to provide examples, and unable to provide the formalized knowledge underlying these examples. This is why we are developing specific programming languages which enable the field expert to intervene directly in some of the linguistic tasks. Our approach is thus not particularly well adapted to the TREC competition, but our results show that the whole system is functional and that it provides usable information. In this TREC competition we worked at two levels of our complete chain. In one level, we stopped the linguistic treatment at the level of terminology (i.e., detecting the collocations relevant to the text). Relevance was then defined as the appearance of the same terms in the task definition (exactly as given by the TREC competition team) and in the texts. Our relatively poor results show that we should have been using relevance definitions extended by human-provided comments. Novelty was defined by a $TF \times IDF$ measurement which seems to work quite correctly, but that could be improved by using the expert-defined concepts as we shall now see. The second level stopped the linguistic treatment after the definition of the concepts. Relevance was then defined presence of a relevant concept and novelty as presence of a new concept. For each of the 5 runs, this approach proved to be less efficient than the simpler first one. We noticed however that the use of concepts enabled us to obtain excellent results on specific topics (and extremely bad ones as well) in different runs. We explain these very irregular results by our own lack of ability to define properly the relevant concepts for all the 50 topics since we got our best results on topics that either we understood well (e.g., Pinochet, topic N51) or that were found interesting (e.g., Lt-Col Collins, topic N85).

References

- [1] A. Amrani, Y. Kodratoff, and O. Matte-Tailliez. A semi-automatic system for tagging specialized corpora. *Proceedings of the Eighth Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD'04)*, 3056:670–681, 2004.
- [2] E. Brill. Some advances in transformation-based part of speech tagging. In *AAAI*, volume 1, pages 722–727, 1994.
- [3] B. Daille, E. Gaussier, and J.M. Lang. An evaluation of statistical scores for word association. In *J.Ginzburg, Z. Khasidashvili, C. Vogel, J.-J. Levy, and E. Vallduvi (eds) The Tbilisi Symposium on Logic, Language and Computation: Selected Papers, CSLI Publications*, pages 177–188, 1998.
- [4] E. Diday, J. Lemaire, J. Poujet, and F. Testu. *Eléments d'analyse des données*. Paris:Dunod, 1982.

- [5] S. Evert and H. Kermes. Experiments on Candidate Data for Collocation Extraction. In *Proceedings of the 10th Conference of The European Chapter of the Association for Computational Linguistics*, pages 83–86, 2003.
- [6] D. Faure and T. Poibeau. First experiments of using semantic knowledge learned by asium for information extraction task using intex. *Ontology Learning, ECAI-2000 Workshop*, pages 7–12, 2000.
- [7] C. Jacquemin. Syntagmatic and paradigmatic representations of term variation. In *Proceedings, 37th Annual Meeting of the Association for Computational Linguistics (ACL'99)*, p. 341–348., 1999.
- [8] Y. Kodratoff. Comparing machine learning and knowledge discovery in databases: An application to knowledge discovery in texts. *Machine Learning and its Applications*, pages 1–21, 2001.
- [9] R.S. Michalski and R.E. Stepp. Learning from observation: Conceptual clustering. *Machine Learning: An Artificial Intelligence Approach*, page 1983, 331-363.
- [10] George A. Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J. Miller. Introduction to wordnet: an on-line lexical database. *International Journal of Lexicography* 3 (4), pages 235–244, 1990. revised august 1993.
- [11] M. Roche, J. Azé, O. Matte-Tailliez, and Y. Kodratoff. Mining texts by association rules discovery in a technical corpus. In *Proceedings of IIPWM'04 (Intelligent Information Processing and Web Mining)*, Springer Verlag series "Advances in Soft Computing", pages 89–98, 2004.
- [12] M. Roche, T. Heitz, O. Matte-Tailliez, and Y. Kodratoff. EXIT: Un système itératif pour l'extraction de la terminologie du domaine à partir de corpus spécialisés. In *Proceedings of JADT'04 (International Conference on Statistical Analysis of Textual Data)*, volume 2, pages 946–956, 2004.
- [13] F. Xu, D. Kurz, J. Piskorski, and S. Schmeier. A Domain Adaptive Approach to Automatic Acquisition of Domain Relevant Terms and their Relations with Bootstrapping. In *LREC 2002, the third international conference on language resources and evaluation*, 2002.