

This is the peer reviewed version of the following article:

Multilingual Word Sense Induction to Improve Web Search Result Clustering / Albano, Lorenzo; Beneventano, Domenico; Bergamaschi, Sonia. - STAMPA. - (2015), pp. 272-279. (Intervento presentato al convegno 23rd Italian Symposium on Advanced Database Systems, SEBD 2015 tenutosi a Gaeta, Italy nel 14-17 June 2015).

Curran Associates

Terms of use:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

14/05/2024 14:26

(Article begins on next page)

Multilingual Word Sense Induction to Improve Web Search Result Clustering (Extended abstract) *

Lorenzo Albano, Domenico Beneventano, and Sonia Bergamaschi

University of Modena and Reggio Emilia, DIF
Via P. Vivarelli 10, 41126 Modena, Italy
{lorenzo.albano,
domenico.beneventano,
sonia.bergamaschi}@unimore.it
<http://www.dbgroup.unimo.it>

Abstract. In [12] a novel approach to Web search result clustering based on Word Sense Induction, i.e. the automatic discovery of word senses from raw text was presented; key to the proposed approach is the idea of, first, automatically inducing senses for the target query and, second, clustering the search results based on their semantic similarity to the word senses induced. In [1] we proposed an innovative Word Sense Induction method based on multilingual data; key to our approach was the idea that a multilingual context representation, where the context of the words is expanded by considering its translations in different languages, may improve the WSI results; the experiments showed a clear performance gain. In this paper we give some preliminary ideas to exploit our multilingual Word Sense Induction method for Web search result clustering.

Keywords: Multilingual, Word Sense Induction, Web Search Result Clustering

1 Introduction

The use of word senses in place of surface word forms has been shown to improve performance on many computational tasks such as information extraction [5], information retrieval [20], data integration [3, 16], machine translation [21] and intelligent web search [14]. Two main techniques have been proposed to solve the word ambiguity problem from different perspectives: Word Sense Disambiguation (WSD) and Word Sense Induction (WSI). Word Sense Disambiguation is aimed at assigning word senses from a predefined *sense inventory* - such as

* Extended abstract of the paper L. Albano, D. Beneventano, S. Bergamaschi: Multilingual Word Sense Induction to Improve Web Search Result Clustering, to appear to the Multilingual Web Access workshop, at the 24th International World Wide Web Conference, Florence, Italy - May 18-22, 2015.

the WordNet [13] database - to words in context while Word Sense Induction is based on clustering words according to their meanings without the use of a Sense Inventory. In other words, given a target word used in a number of different contexts, WSI is the process of clustering these instances of the target word together by determining which contexts are the most similar to each other.

In a previous paper [1] we experimentally evaluated whether the use of more than one language at a time for representing the context of a word may have positive effect on the performance of a Word Sense Induction task or, conversely, the noise increase invalidates the benefits given by a multilingual context representation. The experiments showed a clear overall improvement of the performance: the single-language setting is outperformed by the multi-language settings on almost all the considered words. The registered performance gain reaches peaks of about 40% with certain words and even if we consider the average F-Measure, the difference between the single-language case and the five-language case is about 5%.

In this paper we give some preliminary ideas to exploit our multilingual Word Sense Induction method to Web search results clustering. Web search result clustering aims to facilitate information search on the Web; rather than the results of a query being presented as a list of text snippets, they are grouped on the basis of their similarity: each cluster is intended to represent a different meaning of the input query, thus taking into account the lexical ambiguity (i.e., polysemy) issue [12]. As observed in [12], many Web clustering engines group search results on the basis of their lexical similarity: as a result, text snippets with no word in common tend to be clustered separately even if they share the same meaning, whereas snippets with words in common may be grouped together even if they refer to different meanings of the input query.

To give an intuition of the approach, let us consider the following keyword query: *protect snow leopard*; Google search returns, among others, the following snippets (subscript *EN* denotes that we are using the English language as target language):

- $S1_{EN}$: World's leading authority on the study and protection of the endangered snow leopard.
- $S2_{EN}$: One of the hidden features of Snow Leopard is a built-in system to protect Mac users from malware
- $S3_{EN}$: The law especially protects Snow Leopards and hunting one is culpable of punishment by imprisonment and fines.

As observed in [12], although snippets 1 and 3 refer to the same meaning, they have no content word in common apart from our query words. As a result, a traditional Web clustering engine would most likely assign these snippets to different clusters. It has been shown in [12] that this problem can be addressed, thanks to a novel approach to Web search result clustering based on Word Sense Induction; the key of this approach is to first acquire the various senses (i.e., meanings) of an ambiguous query and then cluster the search results based on their semantic similarity to the word senses induced. The experiments, conducted

on data sets of ambiguous queries, have shown that this approach outperforms both Web clustering and search engines.

The paper is organized as follows. In Section 2 the Multilingual Word Sense Induction (ML-WSI) method is outlined. Section 3 describes the proposed approach to exploit the ML-WSI method for Web search result clustering. Finally, in section 4, conclusions are drawn.

2 The Multilingual Word Sense Induction Method

In [1] we proposed the Multilingual Word Sense Induction (ML-WSI) method which performs Word Sense Induction on words of a *target language* by using *other languages* as support. The ML-WSI is based on the so-called *context clustering* approach [17] (the idea is that a given word, used in a specific sense, tends to co-occur with the same neighboring words [9]) composed by the following two steps: *A) Context Representation*, where each occurrence of a target word in a corpus is represented as a *vector of features*, and *B) Context clustering*, where context representation are clustered and, then, word's are grouped by meaning.

To give an intuition, we use the simplest features, which are the *unigrams* (individual words) composing the context, then each snippet of the before example is represented as a vector of words:

ctx_{1en} : { World, leading, authority, study, protection, endangered, snow leopard }
 ctx_{2en} : { hidden, features, Snow Leopard, built-in, system, protect, Mac, users, malware }
 ctx_{3en} : { law, especially, protects, Snow Leopards, hunting, culpable, punishment, imprisonment, fines }

The performance of the clustering step highly depends on the quality of the features used for the context representation [8]; in [1], we demonstrated that a *multilingual* context representation, where the context of the words is expanded by considering its translations in different languages, can improve the performance of the WSI process. In other words, the key of our ML-WSI method is to perform Word Sense Induction on a *target language* by using *other languages* as support.

To give an idea of our approach, by considering as target the English language and by using the Italian language as support language, for the above three snippets we consider their Italian translation:

$S1_{IT}$: Delle principali autorità del mondo sullo studio e la tutela del **leopardo delle nevi** in via di estinzione.
 $S2_{IT}$: Una delle caratteristiche nascoste di **Snow Leopard** e' un sistema integrato per proteggere gli utenti Mac da malware.
 $S3_{IT}$: La legge in particolare tutela **leopardi delle nevi** e caccia uno e' colpevole della pena con la reclusione e multe.

and their respective context vectors:

$ctx_{1it} : \{ \text{principali, autorità, mondo, studio, tutela, leopardo, nevi, via, estinzione} \}$
 $ctx_{2it} : \{ \text{caratteristiche, nascoste, Snow Leopard, sistema, integrato, proteggere, utenti, Mac, malware} \}$
 $ctx_{3it} : \{ \text{legge, particolare, tutela, leopardi, nevi, caccia, colpevole, pena, reclusione, multe} \}$

Intuitively, it can be seen that *snow leopard* has been translated in different ways ('Snow Leopard', 'leopardo delle nevi') depending on the meaning that it assumes in the sentence; thus a *multilingual context vector*, i.e., the union of the English and Italian context vectors, may improve the clustering process.

A WSI algorithm based on the *context clustering* approach uses a corpus, i. e. a large and structured set of texts. For our ML-WSI algorithm we need a multi-lingual parallel corpus, where each sentence is translated in different languages; the ML-WSI method relies on such translated sentences in order to obtain an increase in performance in the WSI task; a low quality translation may not only void the benefits given by the translation information but it may in addition lead to a loss of performance because of the noise introduced by wrongly translated words. In [1], we considered as target the English language and as support four languages, Italian, French, Spanish and Portuguese; we used a multi-lingual parallel corpus JRC-Acquis [18]; this domain-specific corpus (it is a collection of European Union laws) is available in 22+ languages and the translation of the original text in all this languages has been conducted by expert translators and, thus, the results the translations are very accurate.

3 Web Search Results Clustering with ML-WSI

In this section we outline the main steps of the approach to Web search result clustering based on WSI proposed in [12]; then, we describe how such steps are modified to adopt our ML-WSI method. Our purpose is to extend the approach presented in [12] with our ML-WSI approach in order verify if the Web Search Clustering scenario can get benefits from the additional information introduced by a multilingual representation of the text.

The aim of our approach is to perform Web Search Results Clustering on queries of a target languages by using other languages as support. In order to simplify the presentation, and without loss of generality, we consider the English language as target and the Italian language as support. Let Q_{EN} be a query in the target language.

1. all the possible word senses of Q_{EN} are induced by a WSI algorithm from a text corpus;

With the ML-WSI algorithm, this step is performed by the following two steps:

- 1a) Q_{EN} is translated in the support language;
- 1b) all the possible word senses of $Q_{EN} \cup Q_{IT}$ are induced by the ML-WSI algorithm from a bilingual text corpus;

2. Q_{EN} is executed by the web search engine in order to get the search result snippets: $R = (S1_{EN}, ..., Sn_{EN})$;
It should be noted that the query is only executed in the target language (English).
3. Each search snippet S_{EN} is processed and mapped to the most appropriate meaning by the WSI algorithm.
In order to use the ML-WSI algorithm, this step requires that each snippet S_{EN} is translated in the support language
4. The resulting clustering of snippets in $R = (S1_{EN}, ..., Sn_{EN})$ is returned.

To exploit the ML-WSI method for Web search result clustering two fundamental ingredients are required:

1. Query and Snippet Translation;
2. Multilingual Text Corpus.

3.1 Query and Snippet Translation

To translate in different languages both the query (step 1a) and the search snippets (step 3) a naïve approach is to use a Machine Translation (MT) system. MT systems have worse performance compared to human translators; however, the performance of a state-of-the-art MT system should be enough for our purpose since we do not need a faultless translation. Moreover the fact that we translate the snippets and the query with the same MT system could be an advantage because the translations will probably be consistent. More specifically, translation is applied to snippets and not to the single words of a context features as phrase-based translation models outperform word-based translation models for almost all language pairs [10]. In our intuitive examples, we used one of the most common Machine Translation systems, the Google Translator (GT), which adopt a phrase-based translation model.

The approach based on Machine Translation systems is only a naïve approach, since such MT systems themselves have a word sense disambiguation component built in a component which is integral to the MT processes itself. Future works will be focused on the use of the so-called Cross Lingual Snippet Generation systems [11], which generate snippets in multiple languages starting from documents available only in one language. Another interesting scenario to consider is when the same web page is available in more languages; in this case we will evaluate whether the Snippet Translation process may be performed by the so-called Multi Lingual Snippet Generation systems which are able to generate snippets in multiple languages [15].

3.2 Multilingual Text Corpus

As said in section 2 the ML-WSI method relies on a multi-lingual parallel corpus; the web search clustering is a general-domain application and the use of a domain-specific corpus (such as the multi-lingual parallel corpus JRC-Acquis we

used in [1]) can adversely affect the clustering performance. As a consequence, a multilingual and domain-independent corpus is needed to perform step 1b of our method. While there are many single language corpora (especially in english language), the multilingual corpora are few and their size is often not comparable with the typical size of a single language corpus.

Future works will follow two directions. First, we will consider general multilingual corpus, like DBPedia [2], and other multilingual corpus, like [19, 6]. Second, we will consider the automatic translation in different languages of the existent single language corpora, i.e., we will experimentally verify if the use of the multilingual corpus could be effectively replaced by the use of a MT system without a significative loss in performance for the specific application. In particular, to evaluate the impact of the ML-WSI method in the Web search result clustering, we will also consider the same two corpora used in [12]:

- **Google Web1T** [4]: This corpus is a large collection of n -grams ($n = 1, \dots, 5$)-namely, windows of n consecutive tokens-occurring in one terabyte of Web documents as collected by Google.
- **ukWaC** [7]: This corpus was constructed by crawling the .uk domain and obtaining a large sample of Web pages that were automatically part-of-speech tagged using the TreeTagger tool. For this corpus we considered all the co-occurrences of WordNet lemmas that appear in the same sentence.

4 Conclusions

Word Sense Induction has been shown useful in many scenario. In [12], Word Sense Induction was proposed as a novel approach to Web search result clustering, in the context of an Intelligent Web Search scenario. In our previous paper [1] we proposed the ML-WSI method, an innovative Word Sense Induction method based on multilingual data. In the present paper, the adoption of the ML-WSI method to improve the performance of Web search result clustering is proposed and some issues for the future work are individuated.

We conclude with some considerations concerning the use of the ML-WSI method for Multilingual Web Access concerned with retrieval from the Web, where documents in multiple languages co-exist and need to be retrieved to a query in any language. It should be noted that, even if we talked about target and support language and our clustering considered only web search results in English, the method can be used also to retrieve documents in any language, because we translate both the query and the resulting snippets. This mean that we can execute one query for each considered language, and then we can translate the obtained snippets for that language, in all the other languages. In this way we obtain snippets written in the same set of languages, independently from the language used for the query and from the language of the retrieved document, and we can cluster by topic or meaning all the documents written in different languages.

References

1. Albano, L., Beneventano, D., Bergamaschi, S.: Word sense induction with multilingual features representation. In: Web Intelligence (WI) and Intelligent Agent Technologies (IAT), 2014 IEEE/WIC/ACM International Joint Conferences on. vol. 2, pp. 343–349. IEEE (2014)
2. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: Dbpedia: A nucleus for a web of open data. Springer (2007)
3. Beneventano, D., Bergamaschi, S., Sorrentino, S.: Extending wordnet with compound nouns for semi-automatic annotation in data integration systems. In: Natural Language Processing and Knowledge Engineering, 2009. NLP-KE 2009. International Conference on. pp. 1–8 (Sept 2009)
4. Brants, T., Franz, A.: {Web 1T 5-gram Version 1}. Linguistic Data Consortium, Philadelphia (2006)
5. Chai, J., Biermann, A.: The use of word sense disambiguation in an information extraction system. Proceedings of Sixteenth National Conference in Artificial Intelligence and Eleventh Annual Conference on Innovative Applications of Artificial Intelligence (July 1999)
6. Eermak, F., Rosen, A.: The case of intercorp, a multilingual parallel corpus. International Journal of Corpus Linguistics 17(3), 411–427 (2012), <http://www.scopus.com/inward/record.url?eid=2-s2.0-84873188562&partnerID=40&md5=cb266efd4b16aa377595dc7bc6dac5c2>, cited By 1
7. Ferraresi, A., Zanchetta, E., Baroni, M., Bernardini, S.: Introducing and evaluating ukwac, a very large web-derived corpus of english. In: Proceedings of the 4th Web as Corpus Workshop (WAC-4) Can we beat Google. pp. 47–54 (2008)
8. H., S.: Automatic word sense discrimination. Computational Linguistics (1998)
9. Harris, Z.: Distributional structure (1954)
10. Koehn, P., Och, F.J., Marcu, D.: Statistical phrase-based translation. In: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1. pp. 48–54. NAACL '03, Association for Computational Linguistics, Stroudsburg, PA, USA (2003), <http://dx.doi.org/10.3115/1073445.1073462>
11. Lohar, P., Bhaskar, P., Pal, S., Bandyopadhyay, S.: Cross lingual snippet generation using snippet translation system. In: Computational Linguistics and Intelligent Text Processing, pp. 331–342. Springer (2014)
12. Marco, A.D., Navigli, R.: Clustering and diversifying web search results with graph-based word sense induction. Computational Linguistics (July 2012) (2013), http://www.mitpressjournals.org/doi/abs/10.1162/COLI_a_00148
13. Miller, G.A.: Wordnet: A lexical database for english. Commun. ACM 38(11), 39–41 (Nov 1995), <http://doi.acm.org/10.1145/219717.219748>
14. Navigli, R., Lapata, M.: An experimental study on graph connectivity for unsupervised word sense disambiguation. IEEE Transactions on Pattern Analysis and Machine Intelligence (2010)
15. Pinaki Bhaskar, S.B.: Cross lingual query dependent snippet generation. (IJCSIT) International Journal of Computer Science and Information Technologies 3 (4), 4603–4609 (2014)
16. Po, L., Sorrentino, S.: Automatic generation of probabilistic relationships for improving schema matching. Information Systems, Volume 36, Issue 2 (2011), pp. 192–208 (2011)

17. Purandare, A., Pedersen, T.: Senseclusters - finding clusters that represent word senses. Proceedings of Fifth Annual Meeting of the North American Chapter of the Association for Computational Linguistics (NAACL-04) (May 2004)
18. Steinberger, R., et al.: The jrc-acquis: A multilingual aligned parallel corpus with 20+ languages. Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC'2006). Genoa, Italy (May 2006)
19. Steinberger, R., Ebrahim, M., Poulis, A., Carrasco-Benitez, M., Schlter, P., Przybyszewski, M., Gilbro, S.: An overview of the european unions highly multilingual parallel corpora. *Language Resources and Evaluation* 48(4), 679–707 (2014), <http://www.scopus.com/inward/record.url?eid=2-s2.0-84913548891&partnerID=40&md5=fb5b59a0c25d89293225f5e0e74ae68b>, cited By 0
20. Uzuner, O., Katz, B.: Word sense disambiguation for information retrieval. Proceedings of AAAI/IAAI1999 (July 1999)
21. Vickrey, Biewald, Teyssier, Koller: Word-sense disambiguation for machine translation. HLT/EMNLP (2005)
22. Zhang, M.: Word sense induction for machine translation (invited talk). Proceedings of the 28th Pacific Asia Conference on Language, Information and Computing (December 2014)