

Addressing Cross-Lingual Word Sense Disambiguation on Low-Density Languages: Application to Persian

Navid Rekabsaz¹, Mihai Lupu¹, Allan Hanbury¹, and Andres Duque²

¹ Institute of Software Technology and Interactive Systems
TU WIEN

A-1040 Vienna, Austria

[last_name]@ifs.tuwien.ac.at

² Universidad Nacional de Educacin a Distancia (UNED)
Madrid 28040, Spain
aduque@lsi.uned.es

Abstract. We explore the use of unsupervised methods in Cross-Lingual Word Sense Disambiguation (CL-WSD) with the application of English to Persian. Our proposed approach targets the languages with scarce resources (low-density) by exploiting word embedding and semantic similarity of the words in context. We evaluate the approach on a recent evaluation benchmark and compare it with the state-of-the-art unsupervised system (CO-Graph). The results show that our approach outperforms both the standard baseline and the CO-Graph system in both of the task evaluation metrics (*Out-Of-Five* and *Best result*).

Keywords: Word Sense Disambiguation, cross-lingual, semantics, Word2Vec

1 Introduction

Word Sense Disambiguation (WSD) is the task of automatically selecting the most related sense for a word occurring in a context. It is considered as a main step in the course of approaching language understanding beyond the surface of the words.

Typically, WSD methods are classified into knowledge-based, supervised, and unsupervised. Knowledge-based approaches use available structured knowledge. Supervised approaches learn a computational model based on large amounts of annotated data. While these two approaches show competitive results in practice, they both have to face the knowledge acquisition bottleneck. This is a particular problem in specific domains or low-density languages. As an alternative, unsupervised approaches address WSD using only information extracted from existing corpora, such as various word co-occurrence indicators.

As two well-known benchmarks for CL-WSD, SemEval-2010 [9] and SemEval-2013 [10] provide an evaluation platform for word disambiguation from English to Dutch, German, Italian, Spanish, and French. Recently, Rekabsaz et al. [13] added the Persian (Farsi) language to this set by following the CL-WSD SemEval format to create the test collection.

Many participating systems in the SemEval tasks exploit parallel corpora, mainly Europarl [8], to overcome the knowledge acquisition bottleneck [11,14]. However, the

approaches used in the tasks are not applicable for many languages and domains due to the scarcity of bilingual corpora. Persian, for instance, suffers from the lack of reliable and comprehensive knowledge resources as well as parallel corpora. In such cases, unsupervised methods based on monolingual corpora (together with bilingual lexicon) are preferable, if not the only available option [18]. For example, Bungum et al. [2] find the probable translations of a context in the source language and identify the best translation using a language model of the target language. Duque et al. [5] build a co-occurrence graph in the target language, and test a variety of graph-based algorithms for identifying the best translation match.

In terms of combining Word Sense Disambiguation (WSD) and word embedding, Chen et al. [3] use knowledge-based WSD to identify distinct representations for different senses of the same term. Our approach for CL-WSD is the opposite of this: starting from word embedding representations, it identifies the similarity of the potential translations to the terms in their contexts and choose the translation with the highest semantic similarity to its context.

In order to evaluate our approach, we use the new benchmark of English to Persian CL-WSD [13], and compare our approach and the CO-Graph system [5], observing the advantages of using word embedding in CL-WSD.

In terms of related work addressing the CL-WSD problem in Persian, Sarrafzadeh et al. [15] follows a knowledge-based approach by exploiting FarsNet [17]. However, since their evaluation collection is not available, the results are impossible to compare with other possible approaches.

The remainder of this paper is organized as follows: Section 2 explains our unsupervised approach to English to Persian CL-WSD. We explain the outline of the experiments in Section 3, followed by discussing the result in Section 4. Finally, the study is concluded in Section 5.

2 Unsupervised CL-WSD Method

Our approach follows the main idea of the Lesk algorithm [12], namely that terms in a given context tend to share a common topic. We use word embedding to compute the semantic similarity between terms. We measure the relatedness of each possible translation of an ambiguous term to all possible translations of the terms in its context (the paragraph given by the task) and select the most similar translation to the context.

To formulate our CL-WSD approach, let us define T as the list of translation sets for the terms in a context: $T = \{T_1, T_2, \dots, T_n\}$ where n is the number of terms in the context, and T_i is the set of possible translations for the i^{th} term in the context. For each translation $t \in T_i$, we also have $P(t)$ —an indicator of how frequent this particular translation is.

Given an embedding model in the target language, we compute the similarity of two translation terms t and \bar{t} using their embedding vectors. However, in some cases the translation t of one term in English may be two or more words in Persian (multi-word term), and since our word embedding model is generally created on words level, we will have more than one vector. Therefore, assuming every term t as a set words w ,

we define a general similarity function between two translation terms as follows:

$$\text{sim}(t, \bar{t}) = \max_{w \in t, \bar{w} \in \bar{t}} (\cos(V_w, V_{\bar{w}})) \quad (1)$$

where V_w is the vector representation of the word, and \cos is the cosine function.

Having a definition of similarity between two translation terms, we now move to defining the similarity between a translation candidate of the ambiguous term and the list of translation sets T the terms in the context. We consider two ways to approach it:

The first, denoted as *RelAgg*, uses the *ContextVec* function to create a vector, representing the translated context terms in the target language. The *ContextVec* function is defined in Algorithm 1.

Algorithm 1: ContextVec

Input: translation candidate t , and the list of translation sets T
Output: vector representation of the context
 $sumVec \leftarrow []$;
for $T_i \in T$ **do**
 $t^* \leftarrow \arg \max_{\bar{t} \in T_i} (\text{sim}(t, \bar{t}))$;
 $maxVec \leftarrow V_{t^*}$;
 $sumVec \leftarrow sumVec + maxVec$;
return $\text{norm}(sumVec)$;

The *norm* function in Algorithm 1 applies the Euclidean norm.

Given the vector representation of the context, *RelAgg* calculates the cosine between the vector of each translation candidate t to the context vector, multiplied by the probability of the translation candidate $P(t)$, shown as follows:

$$\text{RelAgg}(t, T) = \cos(V_t, \text{ContextVec}(t, T))P(t) \quad (2)$$

The second approach, denoted as *RelGreedy*, searches among all the translation terms in all the sets T_i , and returns the value of the most similar translation term to the translation candidate. Similar to *RelAgg*, the final score is multiplied by the probability of the translation candidate. The *RelGreedy* approach is defined as follows:

$$\text{RelGreedy}(t, T) = \max_{T_i \in T} \left(\max_{\bar{t} \in T_i} (\text{sim}(t, \bar{t})) \right) P(t) \quad (3)$$

Finally, given the score of the similarity of each translation candidate t_i to its context using either *RelAgg* or *RelGreedy*, we can select the best translation among the candidates, as follows:

$$\text{Result} = \arg \max_{t_i} (\text{Rel}^*(t_i, T)) \quad (4)$$

where t_i is a translation candidate for the term with ambiguity, and Rel^* is either *RelAgg* or *RelGreedy*.

3 Experiments Setup

Resources Similar to Jadidinejad et al. [7], we use the PerStem tool [4] for stemming and TagPer [16] for POS tagging of Persian language. We create a word2vec SkipGram model on a stemmed corpus of Hamshahri collection [1] with sub-sampling at $t = 10^{-4}$, the context windows of 5 terms, epochs of 25, terms count threshold of 5, and dimension of 200.

Beside the monolingual word embedding, a bilingual lexicon is required for our unsupervised CL-WSD approach. While using parallel corpora is considered as a more effective method for creating lexica [6], due to the lack of reliable parallel corpora, we have to use a simple English to Persian dictionary. To have it in digital form, we use the online API of one of the existing translation services³.

Benchmark As mentioned before, we use the novel English to Persian CL-WSD collection [13], which follows the format of SemEval-2013 test collection. The collection consists of 20 nouns, each with 50 cases (paragraphs) in English where the sense of each noun in its corresponding paragraphs is ambiguous. The aim of the benchmark is to find the correct Persian translations of the ambiguous terms.

Evaluation As the official evaluation measure of the SemEval 2013 CL-WSD task [10], we use the F score (harmonic mean of precision and recall), applied in two settings:

- *Best Result* (`Best`), in which a system suggests any number of translations for each target term, and the final score is divided by the number of these translations.
- *Out-Of-Five* (`OO5`) as a more relaxed evaluation setting, in which the system provides up to five different translations, and the best one is selected as the final score.

Baselines The first—STD—is introduced in the SemEval 2013 CL-WSD task as a basic baseline. Similar to the original collection paper, to create the baseline we select the most common and the five most common translations for the `Best` and `OO5` settings respectively.

For the second baseline, we evaluate the Persian benchmark on the state-of-the-art unsupervised CL-WSD system, called CO-Graph [5]. The CO-Graph system offers competitive results in the SemEval 2010 and SemEval 2013 CL-WSD tasks, for all the proposed languages. It outperforms all of the unsupervised participating systems using only monolingual corpora, and even most of the ones which use parallel corpora or knowledge resources. To evaluate the CO-Graph system on the Persian benchmark, we first create the graph using the articles of the Hamshahri collection, each as a document. In the construction of the graph, we only take into account the nouns by POS tagging. After evaluating various algorithms, we find the Dijkstra algorithm together with $p\text{-value}=10^{-6}$ as the best performing approach.

Table 1: Results of F-measure on `OOF` and `Best` evaluation settings.

Setting	Method	F-measure
<code>OOF</code>	RelAgg	0.502
	RelGreedy	0.493
	CO-Graph Dijkstra [5]	0.441
	STD	0.418
<code>Best</code>	RelAgg	0.188
	RelGreedy	0.183
	CO-Graph Dijkstra [5]	0.174
	STD	0.158

4 Results and Discussion

To find the translation for each ambiguous noun, we first apply POS tagging on the English sentences of the SemEval 2013 CL-WSD task and only select the verbs and nouns as the context of the ambiguous terms. We then lemmatize the context terms using WordNetLemmatizer of the NLTK toolkit and find their translations in the bilingual lexicon. Using the word embedding of the translated terms, we finally calculate the relatedness score of each translation candidate to its context using RelAgg and RelGreedy. The translation probability rate in our lexica is used as the $P(t)$ value in Eq. 2 and Eq. 3.

Table 1 shows the F-measure results of RelAgg and RelGreedy as well as the baselines on the `OOF` and `Best` evaluation settings. The results for both evaluation settings show that our approach outperforms the standard and the CO-Graph baselines. Comparing the approaches, we observe similar results for the RelAgg and RelGreedy methods, while RelAgg has slightly better performance, specially in the `OOF` evaluation setting.

In Figure 1, we compare the effectiveness of our methods per each term with the baselines in the `Best` setting as the more challenging one. The results show that while for most terms, our approach outperforms the standard baseline as well as the CO-Graph system, none of the systems can outperform the standard baseline for the terms ‘mood’ and ‘side’. Analyzing the results of these terms, we observe that in some sentences, none of the nouns and verbs in the context share any common topic with senses of the ambiguous terms. For example, using only the semantics of the nouns and verbs in the context, the correct sense of ‘mood’ cannot be distinguished in either of the sentences: ‘it reflected the *mood* of the moment’ (state of the feeling) and ‘a general *mood* in Whitehall’ (inclination, tendency). Similar cases are observed for the term ‘side’: e.g., ‘both *sides* reaffirmed their commitment’ (groups opposing each other) in comparison to ‘at the *side* of the cottage’ (a position to the left or right of a place). While these examples show the limitations of the context-based methods, the overall results show the ability of word embedding and statistical-based approaches for the CL-WSD tasks, specially in the absence of reliable resources.

³ Available in https://github.com/navid-rekabsaz/wsd_persian/tree/master/resources/dictionary

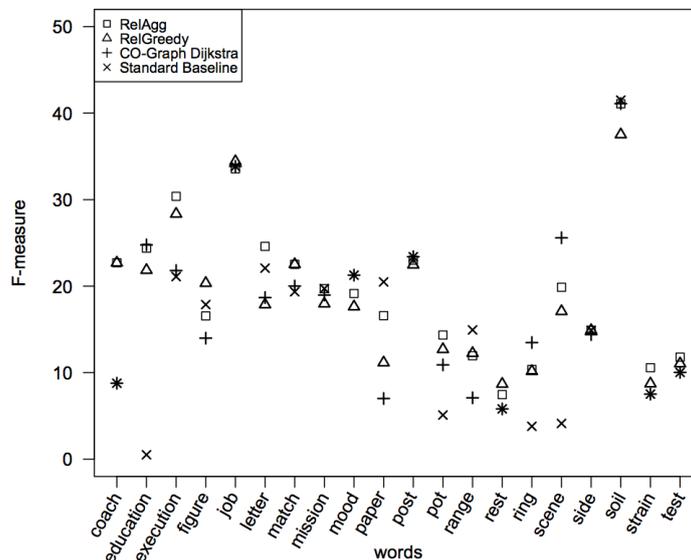


Fig. 1: F-measure results (multiplied by 100) of the Best evaluation setting, per each term

5 Conclusion

We study the application of word embedding-based methods in unsupervised Cross Language Word Sense Disambiguation (CL-WSD) when translating an English noun, appeared in a paragraph, to Persian. Our semantic approach uses embedding of the candidate translations as well as translated context terms to calculate the semantic similarity of each translation to its context. The proposed approach outperforms both the CO-Graph system—a state-of-the-art system in unsupervised CL-WSD—as well as the standard baseline.

We however observe fundamental limitations of the methods based exclusively on context as bag of words. Despite this fact, the current work offers a possible solution for all languages/domains with scarce knowledge-based or parallel corpora resources, by exploiting the use of a monolingual corpus together with a simple bilingual lexicon.

References

1. A. AleAhmad, H. Amiri, E. Darrudi, M. Rahgozar, and F. Oroumchian. Hamshahri: A standard persian text collection. *Knowledge-Based Systems Journal*, 2009.
2. L. Bungum, B. Gambäck, A. Lynam, and E. Marsi. Improving word translation disambiguation by capturing multiword expressions with dictionaries. In *Proc. of NAACL HLT*, 2013.

3. X. Chen, Z. Liu, and M. Sun. A unified model for word sense representation and disambiguation. In *Proc. of EMNLP*, 2014.
4. J. Dehdari and D. Lonsdale. A link grammar parser for persian. *aspects of iranian linguistics*, 2008.
5. A. Duque, L. Araujo, and J. Martinez-Romo. Co-graph: A new graph-based technique for cross-lingual word sense disambiguation. *Natural Language Engineering Journal*, 2015.
6. A. Duque, J. Martinez-Romo, and L. Araujo. Choosing the best dictionary for cross-lingual word sense disambiguation. *Knowledge-Based Systems Journal*, 2015.
7. A. H. Jadidinejad, F. Mahmoudi, and J. Dehdari. Evaluation of perstem: a simple and efficient stemming algorithm for persian. In *Multilingual Information Access Evaluation I. Text Retrieval Experiments*. 2010.
8. P. Koehn. Europarl: a parallel corpus for statistical machine translation. In *MT Summit*, 2005.
9. E. Lefever and V. Hoste. Semeval-2010 task 3: Cross-lingual word sense disambiguation. In *Proc. of the International Workshop on Semantic Evaluation*, 2010.
10. E. Lefever and V. Hoste. Semeval-2013 task 10: Cross-lingual word sense disambiguation. *Proc. of SemEval*, pages 158–166, 2013.
11. E. Lefever, V. Hoste, and M. De Cock. Parasense or how to use parallel corpora for word sense disambiguation. In *Proc. of ACL-HLT*, 2011.
12. M. Lesk. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proc. of Conference on Systems documentation*, 1986.
13. N. Rekabsaz, S. Sabetghadam, M. Lupu, L. Andersson, and A. Hanbury. Standard test collection for english-persian cross-lingual word sense disambiguation. In *Proc. of LREC*, 2016.
14. A. Rudnick, C. Liu, and M. Gasser. Hltidi: cl-wsd using markov random fields for semeval-2013 task 10. In *Joint Conference on Lexical and Computational Semantics (*SEM)*, 2013.
15. B. Sarrafzadeh, N. Yakovets, N. Cercone, and A. An. Cross-lingual word sense disambiguation for languages with scarce resources. In *Advances in Artificial Intelligence: Proc. of 24th Canadian Conference on Artificial Intelligence*. Springer Berlin Heidelberg, 2011.
16. M. Seraji, B. Megyesi, and J. Nivre. A basic language resource kit for persian. In *LREC*, pages 2245–2252. Citeseer, 2012.
17. M. Shamsfard, A. Hesabi, H. Fadaei, N. Mansoory, A. Famian, S. Bagherbeigi, E. Fekri, M. Monshizadeh, and S. M. Assi. Semi automatic development of farsnet; the persian wordnet. In *Proc. of Global WordNet Conference*, 2010.
18. S. Sofianopoulos, M. Vassiliou, and G. Tambouratzis. Implementing a language-independent mt methodology. In *Proc. of the First Workshop on Multilingual Modeling*. ACL, 2012.