

Improving Supervised Bilingual Mapping of Word Embeddings

Armand Joulin and Piotr Bojanowski and Tomas Mikolov and Edouard Grave

Facebook AI Research

{ajoulin,bojanowski,tmikolov,egrave}@fb.com

Abstract

Continuous word representations, learned on different languages, can be aligned with remarkable precision. Using a small bilingual lexicon as training data, learning the linear transformation is often formulated as a regression problem using the square loss. The obtained mapping is known to suffer from the hubness problem, when used for retrieval tasks (e.g. for word translation). To address this issue, we propose to use a retrieval criterion instead of the square loss for learning the mapping. We evaluate our method on word translation, showing that our loss function leads to state-of-the-art results, with the biggest improvements observed for distant language pairs such as English-Chinese.

1 Introduction

Given continuous word representations for two languages, previous work (Mikolov et al., 2013) has shown that it is possible to learn word translations using a linear mapping between word vectors. Using a small bilingual lexicon as supervision, learning this mapping can be formulated as a regression problem. The learned transformation generalizes well to words that were not observed during training, allowing to extend the lexicon. Learning mappings between word vectors in two languages has many other applications. For example, it enables transferring predictive models from one language to others (Klementiev et al., 2012): one can train a sentiment analysis or spam detection system on English data and transfer it to low-resource languages effortlessly.

The simple method proposed by Mikolov et al. (2013) has been further improved by changing the problem parametrization. One interesting sugges-

tion consists in using ℓ_2 normalized word vectors, and constraining the linear mapping to be orthogonal (Xing et al., 2015). Using this parametrization, the regression problem can be efficiently solved using orthogonal Procrustes analysis (Artetxe et al., 2016; Smith et al., 2017). Despite discarding critical information in the form of vector norms, this has led to improved accuracies on standard benchmarks.

Recent results have suggested that using the square loss to learn the transformation is probably not optimal, as the resulting models suffer from the “hubness problem.” This limitation was addressed in the literature by using other criteria for the inference step, such as the inverted softmax (ISF, Smith et al., 2017) or the cross-domain similarity local scaling (CSLS, Conneau et al., 2017). These observations suggest that the square loss could favorably be replaced by a loss that is more adapted to retrieval or classification tasks. Moreover, this would allow to train the mapping using the same loss as the one used for inference.

Moreover, Artetxe et al. (2016) have shown that supervised mappings can be improved by using a “refinement procedure,” similar to the one presented by Conneau et al. (2017) for unsupervised models. This technique is a form of semi-supervised learning, and our loss function will also leverage the representations of words which do not appear in the training lexicon. Other forms of weak supervision can be considered, such as using exact string matches between the two vocabularies as additional examples.

Our main contribution is to propose a new objective for learning the mapping between word vectors. We show that this formulation is convex, and can be efficiently minimized using the projected subgradient descent algorithm. Using our technique, we can easily incorporate unsupervised information, leading to improvements comparable

to the ones obtained by refinement. We evaluate our approach on standard benchmarks for word translation, showing the effects of using an alternative loss function and unsupervised data.

2 Method

In this section, we describe our approach to learn a bilingual lexicon from a small set of pairs of words, called seeds. We assume that we are given two dictionaries, as well as continuous representations of each word, learned on monolingual data. We frame the problem of learning a bilingual lexicon as estimating a mapping between the word representations in the two languages. This mapping can then be used to infer translation of words which do not appear in the seed lexicon.

2.1 Learning a bilingual mapping

Let us start by introducing some notations. Each word $i \in \{1, \dots, N\}$ in the source language (respectively target language) is associated with a vector $\mathbf{x}_i \in \mathbb{R}^d$ (respectively $\mathbf{y}_i \in \mathbb{R}^d$). For simplicity, we assume that our initial lexicon, or seeds, corresponds to the first n pairs $(\mathbf{x}_i, \mathbf{y}_i)_{i \in \{1, \dots, n\}}$. The goal is to extend the lexicon to all source words $i \in \{n+1, \dots, N\}$, which are not part of the seeds. Following Mikolov et al. (2013), our goal is to learn a linear mapping $\mathbf{W} \in \mathbb{R}^{d \times d}$ between the word vectors of the seed lexicon, by minimizing a measure of discrepancy between mapped word vectors of the source language and word vectors of the target domain:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times d}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{W}\mathbf{x}_i, \mathbf{y}_i), \quad (1)$$

where ℓ is a loss, like the squared Euclidean norm: $\ell_2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2$ (Mikolov et al., 2013). Using the square loss leads to a linear least squares problem, which can be solved in closed form.

Oftentimes, the linear mapping \mathbf{W} is constrained to be orthogonal, *i.e.* such that $\mathbf{W}^\top \mathbf{W} = \mathbf{I}_d$, where \mathbf{I}_d is the identity matrix in dimension d . An argument in favor of using orthogonal matrices is that such mappings preserve distances between word vectors, and as a consequence also preserve word similarities. Moreover, previous work (Xing et al., 2015; Artetxe et al., 2016; Smith et al., 2017) experimentally observed that constraining the mapping in such a way improves the quality of the inferred lexicon. Finally, when using the square loss and

an orthogonal mapping \mathbf{W} , the optimization problem in Eq. (1) also admits a closed form solution (Gower and Dijkstra, 2004).

2.2 Inference

Once a mapping \mathbf{W} is learned by solving the problem from Eq. (1), one can infer word correspondences for words that are not in the initial lexicon. To this end, for each source word i in $\{n+1, \dots, N\}$, one can obtain its translation $t(i)$ in the target language by solving:

$$t(i) \in \arg \min_{j \in \{1, \dots, N\}} \ell(\mathbf{W}\mathbf{x}_i, \mathbf{y}_j). \quad (2)$$

When the squared loss is used, this amounts to computing $\mathbf{W}\mathbf{x}_i$ and performing a nearest neighbor search according to the Euclidean distance:

$$t(i) \in \arg \min_{j \in \{1, \dots, N\}} \|\mathbf{W}\mathbf{x}_i - \mathbf{y}_j\|_2^2. \quad (3)$$

A common observation in the literature is that using a nearest neighbor search to infer the bilingual lexicon suffers from the “hubness problem” (Dinu et al., 2014). Hubs are words that appear too frequently in the neighborhoods of other words, and their existence is an inherent property of data distributions in high dimensional space (Radovanović et al., 2010). This issue was first identified and well studied in the retrieval community (Doddington et al., 1998; Aucouturier and Pachet, 2008; Jegou et al., 2010). To mitigate this effect, a simple solution consists in replacing the square ℓ_2 -norm in the inference problem in Eq. (3) by another criterion. Several such criteria were proposed in the word translation literature, such as the inverted softmax (ISF, Smith et al., 2017) or the cross-domain similarity local scaling criterion (CSLS, Conneau et al., 2017).

This solution, both with ISF and CSLS criteria, is always applied with a transformation W learned using the square loss. Replacing the loss in Eq. (3) creates a discrepancy between the learning of the translation model and the inference. In this work, we propose to directly optimize the CSLS loss in Eq. (1) and therefore have coherent learning and inference procedures. Moreover, the translation model \mathbf{W} will be directly learned to avoid the hubness problem. We discuss the CSLS criterion in more details in the following section, and in Sec. 2.4 we will show that a simple convex relaxation leads to a tractable algorithm.

2.3 Word translation as a retrieval task.

The CSLS loss, as described by [Conneau et al. \(2017\)](#), can be written as:

$$\text{CSLS}(\mathbf{x}, \mathbf{y}) = -2 \cos(\mathbf{x}, \mathbf{y}) + \frac{1}{k} \sum_{\mathbf{y}' \in \mathcal{N}_Y(\mathbf{x})} \cos(\mathbf{x}, \mathbf{y}') + \frac{1}{k} \sum_{\mathbf{x}' \in \mathcal{N}_X(\mathbf{y})} \cos(\mathbf{x}', \mathbf{y}),$$

where $\mathcal{N}_Y(\mathbf{x})$ is the set of k nearest neighbors of point \mathbf{x} in the set of target word vectors $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$, and \cos is the cosine similarity. Note that, as far as the inference procedure is concerned, the second term in the expression of the CSLS loss is useless, as it is constant for all j . However, during training it allows to have a loss function which is symmetric with respect to its two arguments, which is a desirable property.

Let us now write the optimization problem corresponding to learning the bilingual mapping with CSLS. We follow previous work and constrain the linear mapping \mathbf{W} to belong to the set of orthogonal matrices \mathcal{O}_d . Without loss of generality, we assume that word vectors are normalized, *i.e.* $\|\mathbf{x}_i\|_2 = 1$ and $\|\mathbf{y}_i\|_2 = 1$. With these assumptions, we have $\cos(\mathbf{W}\mathbf{x}_i, \mathbf{y}_i) = \mathbf{x}_i^\top \mathbf{W}^\top \mathbf{y}_i$, and our optimization problem can be written as:

$$\begin{aligned} \min_{\mathbf{W} \in \mathcal{O}_d} \frac{1}{n} \sum_{i=1}^n -2 \mathbf{x}_i^\top \mathbf{W}^\top \mathbf{y}_i \\ + \frac{1}{k} \sum_{\mathbf{y}_j \in \mathcal{N}_Y(\mathbf{W}\mathbf{x}_i)} \mathbf{x}_i^\top \mathbf{W}^\top \mathbf{y}_j \\ + \frac{1}{k} \sum_{\mathbf{W}\mathbf{x}_j \in \mathcal{N}_X(\mathbf{y}_i)} \mathbf{x}_j^\top \mathbf{W}^\top \mathbf{y}_i. \end{aligned} \quad (4)$$

2.4 Optimization

Our problem, as formulated in Eq. (4) corresponds to the minimization of a non-smooth cost function over the manifold of orthogonal matrices \mathcal{O}_d . As such, it can be solved using manifold optimization tools ([Boumal et al., 2014](#)). However, as we will show it here, we can instead derive convex relaxations that lead to a simple and tractable minimization algorithm.

In this work, we consider two relaxations of the set \mathcal{O}_d , that we compare empirically. The first one consists in replacing the set of orthogonal matrices \mathcal{O}_d by its convex hull \mathcal{C}_d , that is the set of matrices with singular values smaller than 1 (*i.e.* the unit ball of the spectral norm). The second one, which

is looser, consists in considering the ball of radius \sqrt{d} in Frobenius norm that we denote by \mathcal{B}_d .

Having a convex domain allows us to reason about the convexity of the cost function. We observe that the second and third terms in the CSLS loss can be rewritten as follows:

$$\sum_{\mathbf{y}_j \in \mathcal{N}_k(\mathbf{W}\mathbf{x}_i)} \mathbf{x}_i^\top \mathbf{W}^\top \mathbf{y}_j = \max_{S \in \mathcal{S}_k(n)} \sum_{j \in S} \mathbf{x}_i^\top \mathbf{W}^\top \mathbf{y}_j,$$

where $\mathcal{S}_k(n)$ denotes the set of all subsets of $\{1, \dots, n\}$ of size k . This term, seen as a function of \mathbf{W} , is a maximum of linear functions of \mathbf{W} , which is convex. This observation shows that the CSLS loss is a convex function with respect to the mapping \mathbf{W} , and that it is piecewise linear (hence non-smooth). We minimize this objective function over the two convex sets \mathcal{C}_d and \mathcal{B}_d by using projected subgradient descent.

The projection on the set \mathcal{C}_d consists in taking the singular value decomposition (SVD) of the matrix, and thresholding the singular values to one. Projecting on the set \mathcal{B}_d corresponds to dividing the matrix by its Frobenius norm. For both sets, an orthogonal mapping can be obtained by rounding the solution, *i.e.* computing the SVD and setting the singular values to one.

2.5 Extended Normalization

Usually, the number of word pairs in the seed lexicon n is small with respect to the size of the dictionaries N . To benefit from unlabeled data, it is common to add an iterative “refinement procedure” ([Artetxe et al., 2016](#)) when learning the translation model W . Given a model \mathbf{W}_t , this procedure consists in augmenting the training lexicon, by keeping the best-inferred translation in Eq. (3), and learning a new mapping \mathbf{W}_{t+1} by solving the problem in Eq. (1). This strategy is similar to standard semi-supervised approaches where the training set is augmented over time.

In practice, state-of-the-art methods use different losses in the two steps of the refinement procedure, learning \mathbf{W} with the square loss and inferring the lexicon with the CSLS or ISF loss. This strategy has a risk of diverging as more unlabeled pairs are added, and in general does not have any guarantee of convergence. In this work, we propose to use the unpaired words in the dictionaries as “negatives” in the CSLS loss: instead of computing the k -nearest neighbors $\mathcal{N}_Y(\mathbf{W}\mathbf{x}_i)$ amongst the annotated words $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, we do it over the whole dictionary $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$.

Method	en-es	es-en	en-fr	fr-en	en-de	de-en	en-ru	ru-en	en-zh	zh-en	avg.
Adversarial + refine	81.7	83.3	82.3	82.1	74.0	72.2	44.0	59.1	32.5	31.4	64.3
ICP + refine	82.2	83.8	82.5	82.5	74.8	73.1	46.3	61.6	-	-	-
Procrustes	81.4	82.9	81.1	82.4	73.5	72.4	51.7	63.7	42.7	36.7	66.8
Procrustes + refine	82.4	83.9	82.3	83.2	75.3	73.2	50.1	63.5	40.3	35.5	66.9
CSLS (spectral)	83.0	84.9	82.7	84.1	78.2	75.8	56.4	66.3	44.4	45.6	70.1
CSLS (Frobenius)	84.5	86.4	83.1	84.1	79.1	75.9	57.0	67.1	44.6	41.9	70.4

Table 1: Comparison between our approach, Procrustes and unsupervised approaches. All the methods used the CSLS criterion for the retrieval step. “Refine” refers to the Refinement method of [Conneau et al. \(2017\)](#). Note that “Adversarial” and “ICP” are unsupervised while Procrustes and our approach use a bilingual lexicon ([Conneau et al., 2017](#); [Hoshen and Wolf, 2018](#)).

3 Experiments

Implementation details. For all pairs of languages, we fix the number of epochs to 10 and choose a learning rate in $\{1, 10, 25, 50\}$. We divide the learning by 2 when the loss does not decrease. For the English-Chinese pair (en-zh), we also run our algorithm with and without centering the word vectors. Both parameters are selected on a validation set. All the word vectors are ℓ_2 unit normalized. We set the number of nearest neighbors k in the CSLS loss to 10. When using exact string matches, we exclude pairs that are found in the test set. Unless otherwise specified, we use the fastText word vectors trained on Wikipedia ([Bojanowski et al., 2016](#)).¹

Supervised learning. First, we carry out a comparison of our approach to Orthogonal Procrustes, with and without refinement and two unsupervised approaches. As mentioned in Sec. 2.5, the refinement procedure iteratively trains a mapping and adds pairs of points to the training set, based on their score. Our method only uses unlabeled data in the “negative terms” of the CSLS loss, avoiding the risks of divergence. We compare methods on a set of 5 language pairs (in both directions): en-es, en-fr, en-de, en-ru and en-zh. For our method, we report the performance obtained with both relaxations, namely the convex hull \mathcal{C}_d (CSLS spectral) and the unit ball in Frobenius norm \mathcal{B}_d (CSLS Frobenius). We report results for this comparison in Table 1.

We see that our approach outperforms Orthogonal Procrustes with and without refinement by several percents (+3.5% on average). This shows that

using the same criterion for learning the mapping \mathbf{W} and doing inference is important. Moreover, the refinement procedure only slightly helps Orthogonal Procrustes (+0.1% on average). For languages which are considered distant, it even worsens the accuracy (−1.8% on average for English and Chinese).

Table 1 also compares the performance of solutions obtained with different relaxations to the set of orthogonal matrices. We observe that, on average, using Frobenius CSLS leads to a better performance (+0.3%), with an improved accuracy on all pairs except zh-en. Contrary to observations made in previous work, this experiment suggests that preserving the distance between word vectors is not essential to obtain good word translations. Moreover, we have observed that doing a rounding by setting the singular values of \mathbf{W} to one was further decreasing performance. A regular linear mapping \mathbf{W} with bounded Frobenius norm works well when learned with a retrieval criterion. In the rest of the experiments we report results obtained with the \mathcal{C}_d relaxation (CSLS spectral).

Impact of extended normalization. As mentioned in Sec. 2.5, instead of carrying out a refinement procedure, we instead extend the nearest neighbor search in the CSLS loss to the whole dictionary, including words outside of the lexicon. In order to evaluate the impact of this extension, we train a mapping while searching within n or N candidates. We report results of this comparison in Table 2. We observe that using extended nearest neighbors improves the performance of our method by a few percents (by +2.4% on average). For comparison, when excluding Chinese, the refinement procedure only helps Orthogonal

¹<https://fasttext.cc/docs/en/pretrained-vectors.html>

	Full	Seeds
en-es	83.0	80.7
es-en	84.9	83.9
en-fr	82.7	81.7
fr-en	84.1	83.2
en-de	78.2	75.1
de-en	75.8	72.1
en-ru	56.4	51.1
ru-en	66.3	63.8
avg.	76.4	74.0

Table 2: Comparison between our approach with and without the use of all the words for the nearest neighbor search in the CSLS criterion. Full uses all the words in the dictionary, while Seeds only uses annotated pairs from the seed lexicon for nearest neighbors.

Procrustes with CSLS by +0.6% (see Table 1). Using the full dictionary for the normalization benefits from seeing more data, making the method more robust.

Comparison to the state of the art. Finally, we also report the accuracy of our methods when trained and tested with the setting of Dinu et al. (2014). For fair comparison, we use word vectors learned on the WaCky datasets (Baroni et al., 2009). In this experiment, we select the number of epochs between $\{1, 2, 5, 10\}$ based on the performance on a validation set. We report state-of-the-art results and our performance in Table 3. We observe that we reach state-of-the-art performance on the en-it pair, and perform comparably on it-en. Also, we obtain such good results despite the fact that the seed lexicon in this setup is noisy and small.

4 Conclusion

In this paper, we show that using a retrieval criterion instead of the square loss improves the supervised learning of bilingual mappings. We prove that the CSLS criterion, previously used for retrieval, is convex in \mathbf{W} and thus can be used for learning. In that setup, the same criterion is used for learning the mapping and for inferring the lexicon. Using CSLS, we can also include words that do not appear in the supervised lexicon in the training objective, leading to further improvements. Finally, we experimentally show that when using

	en-it	it-en
Adversarial + refine + CSLS	45.1	38.3
Mikolov et al. (2013)	33.8	24.9
Dinu et al. (2014)	38.5	24.6
Artetxe et al. (2016)	39.7	33.8
Smith et al. (2017)	43.1	38.0
Procrustes + CSLS	44.9	38.5
CSLS (spectral)	45.3	37.9

Table 3: Accuracy on English and Italian with the setting of Dinu et al. (2014). “Adversarial” is an unsupervised technique. The adversarial and Procrustes results are from Conneau et al. (2017). Artetxe et al. (2016) is weakly supervised, as it only uses number tokens as a source of supervision.

this objective, constraining the mapping to be orthogonal does not improve the quality of translations.

References

- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2016. Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2289–2294.
- Jean-Julien Aucouturier and Francois Pachet. 2008. A scale-free distribution of false positives for a large class of audio similarity measures. *Pattern recognition*, 41(1):272–284.
- Marco Baroni, Silvia Bernardini, Adriano Ferraresi, and Eros Zanchetta. 2009. The wacky wide web: a collection of very large linguistically processed web-crawled corpora. *Language resources and evaluation*, 43(3):209–226.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. 2014. [Manopt, a Matlab toolbox for optimization on manifolds](#). *Journal of Machine Learning Research*, 15:1455–1459.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.

- Georgiana Dinu, Angeliki Lazaridou, and Marco Baroni. 2014. Improving zero-shot learning by mitigating the hubness problem. *arXiv preprint arXiv:1412.6568*.
- George Doddington, Walter Liggett, Alvin Martin, Mark Przybocki, and Douglas Reynolds. 1998. Sheep, goats, lambs and wolves: A statistical analysis of speaker performance in the nist 1998 speaker recognition evaluation. Technical report.
- John C Gower and Garmt B Dijksterhuis. 2004. *Procrustes problems*, volume 30. Oxford University Press on Demand.
- Yedid Hoshen and Lior Wolf. 2018. An iterative closest point method for unsupervised word translation. *arXiv preprint arXiv:1801.06126*.
- Herve Jegou, Cordelia Schmid, Hedi Harzallah, and Jakob Verbeek. 2010. Accurate image search using the contextual dissimilarity measure. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1):2–11.
- Alexandre Klementiev, Ivan Titov, and Binod Bhattarai. 2012. Inducing crosslingual distributed representations of words. *Proceedings of COLING 2012*, pages 1459–1474.
- Tomas Mikolov, Quoc V Le, and Ilya Sutskever. 2013. Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*.
- Miloš Radovanović, Alexandros Nanopoulos, and Mirjana Ivanović. 2010. Hubs in space: Popular nearest neighbors in high-dimensional data. *Journal of Machine Learning Research*, 11(Sep):2487–2531.
- Samuel L Smith, David HP Turban, Steven Hamblin, and Nils Y Hammerla. 2017. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. *arXiv preprint arXiv:1702.03859*.
- Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. 2015. Normalized word embedding and orthogonal transform for bilingual word translation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1006–1011.