

Are Girls Neko or Shōjo? Cross-Lingual Alignment of Non-Isomorphic Embeddings with Iterative Normalization

Mozhi Zhang¹ Keyulu Xu² Ken-ichi Kawarabayashi³
Stefanie Jegelka² Jordan Boyd-Graber¹

¹University of Maryland, College Park, Maryland, USA

²Massachusetts Institute of Technology, Cambridge, Massachusetts, USA

³National Institute of Informatics, Tokyo, Japan

{mozhi, jbg}@umiacs.umd.edu {keyulu, stefje}@mit.edu k_keniti@nii.ac.jp

Abstract

Cross-lingual word embeddings (CLWE) underlie many multilingual natural language processing systems, often through orthogonal transformations of pre-trained monolingual embeddings. However, orthogonal mapping only works on language pairs whose embeddings are naturally isomorphic. For non-isomorphic pairs, our method (Iterative Normalization) transforms monolingual embeddings to make orthogonal alignment easier by simultaneously enforcing that (1) individual word vectors are unit length, and (2) each language’s average vector is zero. Iterative Normalization consistently improves word translation accuracy of three CLWE methods, with the largest improvement observed on English–Japanese (from 2% to 44% test accuracy).

1 Orthogonal Cross-Lingual Mappings

Cross-lingual word embedding (CLWE) models map words from multiple languages to a shared vector space, where words with similar meanings are close, regardless of language. CLWE is widely used in multilingual natural language processing (Klementiev et al., 2012; Guo et al., 2015; Zhang et al., 2016). Recent CLWE methods (Ruder et al., 2017; Glavas et al., 2019) independently train two monolingual embeddings on large monolingual corpora and then align them with a linear transformation. Previous work argues that these transformations should be *orthogonal* (Xing et al., 2015; Smith et al., 2017; Artetxe et al., 2016): for any two words, the dot product of their representations is the same as the dot product with the transformation. This preserves similarities and substructure of the original monolingual word embedding but enriches the embeddings with multilingual connections between languages.

Thus, many state-of-the-art mapping-based CLWE methods impose an orthogonal con-

straint (Artetxe et al., 2017; Conneau et al., 2018; Alvarez-Melis and Jaakkola, 2018; Artetxe et al., 2018; Ruder et al., 2018; Alvarez-Melis et al., 2019). The success of orthogonal methods relies on the assumption that embedding spaces are isomorphic; i.e., they have the same inner-product structures across languages, but this does not hold for all languages (Søgaard et al., 2018; Fujinuma et al., 2019). For example, English and Japanese fastText vectors (Bojanowski et al., 2017) have different substructures around “girl” (Figure 1 left). As a result, orthogonal mapping fails on some languages—when Hoshen and Wolf (2018) align fastText embeddings with orthogonal mappings, they report 81% English–Spanish word translation accuracy but only 2% for the more distant English–Japanese.

While recent work challenges the orthogonal assumption (Doval et al., 2018; Joulin et al., 2018; Jawanpuria et al., 2019), we focus on whether simple preprocessing techniques can *improve the suitability of orthogonal models*. Our iterative method normalizes monolingual embeddings to make their structures more similar (Figure 1), which improves subsequent alignment.

Our method is motivated by two desired properties of monolingual embeddings that support orthogonal alignment:

1. Every word vector has the same length.
2. Each language’s mean has the same length.

Standard preprocessing such as dimension-wise mean centering and length normalization (Artetxe et al., 2016) do not meet the two requirements at the same time. Our analysis leads to *Iterative Normalization*, an alternating projection algorithm that normalizes any word embedding to provably satisfy both conditions. After normalizing the monolingual embeddings, we then apply mapping-based CLWE algorithms on the transformed embeddings.

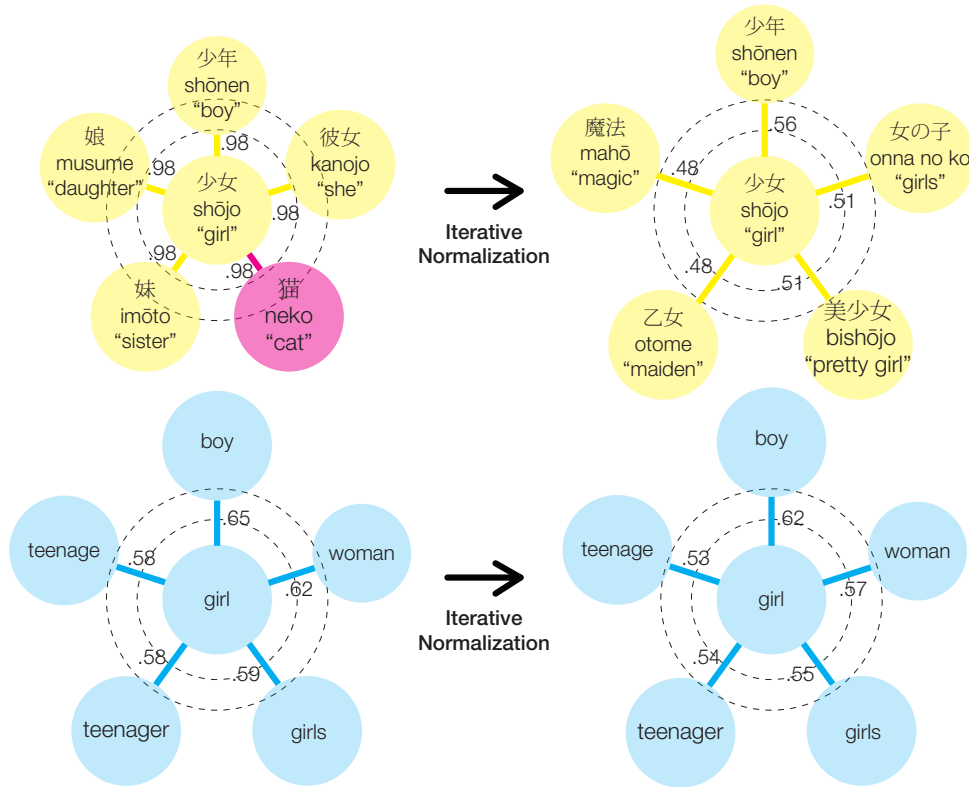


Figure 1: The most similar Japanese words for 少女 (shōjo “girl”) and English words for “girl”, measured by cosine similarity on Wikipedia fastText vectors, before (left) and after (right) Iterative Normalization. In the original embedding spaces, “boy” is the nearest neighbor for both languages but with a very different cosine similarity, and “cat” in English is not close to “girl”: both violate the isomorphism assumed by an orthogonal transformation for cross-lingual representations. Iterative Normalization replaces 猫 (neko “cat”) with the more relevant 美少女 (bishōjo “pretty girl”) and brings cosine similarities closer.

We empirically validate our theory by combining Iterative Normalization with three mapping-based CLWE methods. Iterative Normalization improves word translation accuracy on a dictionary induction benchmark across thirty-nine language pairs.

2 Learning Orthogonal Mappings

This section reviews learning orthogonal cross-lingual mapping between word embeddings and, along the way, introduces our notation.

We start with pre-trained word embeddings in a source language and a target language. We assume¹ all embeddings are d -dimensional, and the two languages have the same vocabulary size n . Let $\mathbf{X} \in \mathbb{R}^{d \times n}$ be the word embedding matrix for the source language, where each column $\mathbf{x}_i \in \mathbb{R}^d$ is the representation of the i -th word from the source language, and let $\mathbf{Z} \in \mathbb{R}^{d \times n}$ be the word embedding matrix for the target language. Our goal is to learn a transformation matrix $\mathbf{W} \in \mathbb{R}^{d \times d}$ that maps the source language vectors to the target lan-

¹Word translation benchmarks use the same assumptions.

guage space. While our experiments focus on the supervised case with a seed dictionary \mathcal{D} with translation pairs (i, j) , the analysis also applies to unsupervised projection.

One straightforward way to learn \mathbf{W} is by minimizing Euclidean distances between translation pairs (Mikolov et al., 2013a). Formally, we solve:

$$\min_{\mathbf{W}} \sum_{(i,j) \in \mathcal{D}} \|\mathbf{W}\mathbf{x}_i - \mathbf{z}_j\|_2^2. \quad (1)$$

Xing et al. (2015) further restrict \mathbf{W} to orthogonal transformations; i.e., $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$. The orthogonal constraint significantly improves word translation accuracy (Artetxe et al., 2016). However, this method still fails for some language pairs because word embeddings are not isomorphic across languages. To improve orthogonal alignment between non-isomorphic embedding spaces, we aim to transform monolingual embeddings in a way that helps orthogonal transformation.

3 When Orthogonal Mappings Work

When are two embedding spaces easily aligned? A good orthogonal mapping is more likely if word vectors have two properties: *length-invariance* and *center-invariance*.

Length-Invariance. First, all word vectors should have the same, constant length. Length-invariance resolves inconsistencies between monolingual word embedding and cross-lingual mapping objectives (Xing et al., 2015). During training, popular word embedding algorithms (Mikolov et al., 2013b; Pennington et al., 2014; Bojanowski et al., 2017) maximize *dot products* between similar words, but evaluate on *cosine similarity*. To make things worse, the transformation matrix minimizes a third metric, *Euclidean distance* (Equation 1). This inconsistency is naturally resolved when the lengths of word vectors are fixed. Suppose $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{v} \in \mathbb{R}^d$ have the same length, then

$$\mathbf{u}^\top \mathbf{v} \propto \cos(\mathbf{u}, \mathbf{v}) = 1 - \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_2^2.$$

Minimizing Euclidean distance is equivalent to maximizing both dot product and cosine similarity with constant word vector lengths, thus making objectives consistent.

Length-invariance also satisfies a prerequisite for bilingual orthogonal alignment: the embeddings of translation pairs should have the same length. If a source word vector \mathbf{x}_i can be aligned to its target language translation $\mathbf{z}_j = \mathbf{W}\mathbf{x}_i$ with an orthogonal matrix \mathbf{W} , then

$$\|\mathbf{z}_j\|_2 = \|\mathbf{W}\mathbf{x}_i\|_2 = \|\mathbf{x}_i\|_2, \quad (2)$$

where the second equality follows from the orthogonality of \mathbf{W} . Equation (2) is trivially satisfied if all vectors have the same length. In summary, length-invariance not only promotes consistency between monolingual word embedding and cross-lingual mapping objective but also simplifies translation pair alignment.

Center-Invariance. Our second condition is that the mean vector of different languages should have the same length, which we prove is a pre-requisite for orthogonal alignment. Suppose two embedding matrices \mathbf{X} and \mathbf{Z} can be aligned with an orthogonal matrix \mathbf{W} such that $\mathbf{Z} = \mathbf{W}\mathbf{X}$. Let $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ and $\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$ be the mean vectors. Then $\bar{\mathbf{z}} = \mathbf{W}\bar{\mathbf{x}}$. Since \mathbf{W} is orthogonal,

$$\|\bar{\mathbf{z}}\|_2 = \|\mathbf{W}\bar{\mathbf{x}}\|_2 = \|\bar{\mathbf{x}}\|_2.$$

In other words, orthogonal mappings can *only* align embedding spaces with equal-magnitude centers.

A stronger version of center-invariance is zero-mean, where the mean vector of each language is zero. Artetxe et al. (2016) find that centering improves dictionary induction; our analysis provides an explanation.

4 Iterative Normalization

We now develop Iterative Normalization, which transforms monolingual word embeddings to satisfy both length-invariance and center-invariance. Specifically, we normalize word embeddings to simultaneously have unit-length and zero-mean. Formally, we produce embedding matrix \mathbf{X} such that

$$\|\mathbf{x}_i\|_2 = 1 \quad \text{for all } i, \quad (3)$$

and

$$\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}. \quad (4)$$

Iterative Normalization transforms the embeddings to make them satisfy both constraints *at the same time*. Let $\mathbf{x}_i^{(0)}$ be the initial embedding for word i . We assume that all word embeddings are non-zero.² For every word i , we iteratively transform each word vector \mathbf{x}_i by first making the vectors unit length,

$$\mathbf{y}_i^{(k)} = \mathbf{x}_i^{(k-1)} / \|\mathbf{x}_i^{(k-1)}\|_2, \quad (5)$$

and then making them mean zero,

$$\mathbf{x}_i^{(k)} = \mathbf{y}_i^{(k)} - \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^{(k)}. \quad (6)$$

Equation (5) and (6) project the embedding matrix \mathbf{X} to the set of embeddings that satisfy Equation (3) and (4). Therefore, our method is a form of alternating projection (Bauschke and Borwein, 1996), an algorithm to find a point in the intersection of two closed sets by alternatively projecting onto one of the two sets. Alternating projection guarantees convergence in the intersection of two convex sets at a linear rate (Gubin et al., 1967; Bauschke and Borwein, 1993). Unfortunately, the unit-length constraint is *non-convex*, ruling out the classic convergence proof. Nonetheless, we use recent results on alternating non-convex projections (Zhu and Li, 2018) to prove Iterative Normalization’s convergence (details in Appendix A).

²For such vectors, a small perturbation is an easy fix.








Method	Normalization	 JA	 ZH	 HI	 TR	 DA	 DE	 ES
Procrustes	None	1.7	32.5	33.3	44.9	54.0	73.5	81.4
	C+L	12.3	41.1	34.0	46.5	54.9	74.6	81.3
	IN	44.3	44.2	36.7	48.7	58.4	75.5	81.5
Procrustes + refine	None	1.7	32.5	33.6	46.3	56.8	74.3	81.9
	C+L	13.1	42.3	34.9	48.7	59.3	75.2	82.4
	IN	44.3	44.2	37.7	51.7	60.9	76.0	82.5
RCSLS	None	14.6	17.1	5.0	18.3	19.2	43.6	50.5
	C+L	16.1	45.1	36.2	50.7	58.3	77.5	83.6
	IN	56.3	48.6	38.0	52.4	60.5	78.1	83.9

Table 1: Word translation accuracy aligning English embeddings to seven languages. We combine three normalizations—no normalization (None), mean centering and length normalization (C+L), and Iterative Normalization (IN) for five rounds—with three CLWEs: Procrustes, Procrustes with refinement (Conneau et al., 2018), and RCSLS (Joulin et al., 2018). Procrustes with C+L is equivalent to Artetxe et al. (2016). The best result for each CLWE in each column in bold. Iterative Normalization has the best accuracy of the three normalization techniques.

Theorem 1. *If the embeddings are non-zero after each iteration; i.e., $\mathbf{x}_i^{(k)} \neq \mathbf{0}$ for all i and k , then the sequence $\{\mathbf{X}^{(k)}\}$ produced by Iterative Normalization is convergent.*

All embeddings in our experiments satisfy the non-zero assumption; it is violated only when all words have the same embedding. In degenerate cases, the algorithm might converge to a solution that does not meet the two requirements. Empirically, our method always satisfy both constraints.

Previous approach and differences. Artetxe et al. (2016) also study the unit-length and zero-mean constraints, but our work differs in two aspects. First, they motivate the zero-mean condition based on the heuristic argument that two randomly selected word types should not be semantically similar (or dissimilar) in expectation. While this statement is attractive at first blush, some word types have more synonyms than others, so we argue that word types might not be evenly distributed in the semantic space. We instead show that zero-mean is helpful because it satisfies center-invariance, a *necessary condition* for orthogonal mappings. Second, Artetxe et al. (2016) attempt to enforce the two constraints by a single round of dimension-wise mean centering and length normalization. Unfortunately, this often fails to meet the constraints *at the same time*—length normalization can change the mean, and mean centering can change vector length. In contrast, Iterative Normalization simultaneously meets both constraints and is empirically better (Table 1) on dictionary induction.

5 Dictionary Induction Experiments

On a dictionary induction benchmark, we combine Iterative Normalization with three CLWE methods and show improvement in word translation accuracy across languages.

5.1 Dataset and Methods

We train and evaluate CLWE on MUSE dictionaries (Conneau et al., 2018) with default split. We align English embeddings to thirty-nine target language embeddings, pre-trained on Wikipedia with fastText (Bojanowski et al., 2017). The alignment matrices are trained from dictionaries of 5,000 source words. We report top-1 word translation accuracy for 1,500 source words, using cross-domain similarity local scaling (Conneau et al., 2018, CSLS). We experiment with the following CLWE methods.³

Procrustes Analysis. Our first algorithm uses Procrustes analysis (Schönemann, 1966) to find the orthogonal transformation that minimizes Equation 1, the total distance between translation pairs.

Post-hoc Refinement. Orthogonal mappings can be improved with refinement steps (Artetxe et al., 2017; Conneau et al., 2018). After learning an initial mapping \mathbf{W}_0 from the seed dictionary \mathcal{D} , we build a synthetic dictionary \mathcal{D}_1 by translating each word with \mathbf{W}_0 . We then use the new dictionary \mathcal{D}_1 to learn a new mapping \mathbf{W}_1 and repeat the process.

³We only report accuracy for one run, because these CLWE methods are deterministic.

Relaxed CSLS Loss (RCSLS). Joulin et al. (2018) optimize CSLS scores between translation pairs instead of Equation (1). RCSLS has state-of-the-art supervised word translation accuracies on MUSE (Glavas et al., 2019). For the ease of optimization, RCSLS does not enforce the orthogonal constraint. Nevertheless, Iterative Normalization also improves its accuracy (Table 1), showing it can help linear non-orthogonal mappings too.

5.2 Training Details

We use the implementation from MUSE for Procrustes analysis and refinement (Conneau et al., 2018). We use five refinement steps. For RCSLS, we use the same hyperparameter selection strategy as Joulin et al. (2018)—we choose learning rate from $\{1, 10, 25, 50\}$ and number of epochs from $\{10, 20\}$ by validation. As recommended by Joulin et al. (2018), we turn off the spectral constraint. We use ten nearest neighbors when computing CSLS.

5.3 Translation Accuracy

For each method, we compare three normalization strategies: (1) no normalization, (2) dimension-wise mean centering followed by length normalization (Artetxe et al., 2016), and (3) five rounds of Iterative Normalization. Table 1 shows word translation accuracies on seven selected target languages. Results on other languages are in Appendix B.

As our theory predicts, Iterative Normalization increases translation accuracy for Procrustes analysis (with and without refinement) across languages. While centering and length-normalization also helps, the improvement is smaller, confirming that one round of normalization is insufficient. The largest margin is on English-Japanese, where Iterative Normalization increases test accuracy by more than 40%. Figure 1 shows an example of how Iterative Normalization makes the substructure of an English-Japanese translation pair more similar.

Surprisingly, normalization is even more important for RCSLS, a CLWE method without orthogonal constraint. RCSLS combined with Iterative Normalization has state-of-the-art accuracy, but RCSLS is much worse than Procrustes analysis on unnormalized embeddings, suggesting that length-invariance and center-invariance are also helpful for learning linear non-orthogonal mappings.

Dataset	Before	After
WS-353	73.9	73.7
MC	81.2	83.9
RG	79.7	80.0
YP-130	53.3	57.6

Table 2: Correlations before and after applying Iterative Normalization on four English word similarity benchmarks: WS-353 (Finkelstein et al., 2002), MC (Miller and Charles, 1991), RG (Rubenstein and Goodenough, 1965), and YP-130 (Yang and Powers, 2006). The scores are similar, which shows that Iterative Normalization retains useful structures from the original embeddings.

5.4 Monolingual Word Similarity

Many trivial solutions satisfy both length-invariance and center-invariance; e.g., we can map half of words to e and the rest to $-e$, where e is any unit-length vector. A meaningful transformation should also preserve useful structure in the original embeddings. We confirm Iterative Normalization does not hurt scores on English word similarity benchmarks (Table 2), showing that Iterative Normalization produces meaningful representations.

6 Conclusion

We identify two conditions that make cross-lingual orthogonal mapping easier: length-invariance and center-invariance, and provide a simple algorithm that transforms monolingual embeddings to satisfy both conditions. Our method improves word translation accuracy of different mapping-based CLWE algorithms across languages. In the future, we will investigate whether our method helps other downstream tasks.

Acknowledgments

We thank the anonymous reviewers for comments. Boyd-Graber and Zhang are supported by DARPA award HR0011-15-C-0113 under subcontract to Raytheon BBN Technologies. Jegelka and Xu are supported by NSF CAREER award 1553284. Xu is also supported by a Chevron-MIT Energy Fellowship. Kawarabayashi is supported by JST ERATO JPMJER1201 and JSPS Kakenhi JP18H05291. Any opinions, findings, conclusions, or recommendations expressed here are those of the authors and do not necessarily reflect the view of the sponsors.

References

- David Alvarez-Melis and Tommi S. Jaakkola. 2018. [Gromov-wasserstein alignment of word embedding spaces](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- David Alvarez-Melis, Stefanie Jegelka, and Tommi S Jaakkola. 2019. [Towards optimal transport with global invariances](#). In *Proceedings of Artificial Intelligence and Statistics*.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2016. [Learning principled bilingual mappings of word embeddings while preserving monolingual invariance](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2017. [Learning bilingual word embeddings with \(almost\) no bilingual data](#). In *Proceedings of the Association for Computational Linguistics*.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. [A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings](#). In *Proceedings of the Association for Computational Linguistics*.
- Heinz H. Bauschke and Jonathan M. Borwein. 1993. On the convergence of von Neumann’s alternating projection algorithm for two sets. *Set-Valued Analysis*, 1(2):185–212.
- Heinz H. Bauschke and Jonathan M. Borwein. 1996. [On projection algorithms for solving convex feasibility problems](#). *SIAM review*, 38(3):367–426.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Felix E. Browder. 1967. Convergence of approximants to fixed points of nonexpansive nonlinear mappings in Banach spaces. *Archive for Rational Mechanics and Analysis*, 24(1):82–90.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. [Word translation without parallel data](#). In *Proceedings of the International Conference on Learning Representations*.
- Yerai Doval, Jose Camacho-Collados, Luis Espinosa-Anke, and Steven Schockaert. 2018. [Improving cross-lingual word embeddings by meeting in the middle](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppin. 2002. [Placing search in context: The concept revisited](#). *ACM Transactions on information systems*, 20(1):116–131.
- Yoshinari Fujinuma, Jordan Boyd-Graber, and Michael J. Paul. 2019. A resource-free evaluation metric for cross-lingual word embeddings based on graph modularity. In *Proceedings of the Association for Computational Linguistics*.
- Goran Glavas, Robert Litschko, Sebastian Ruder, and Ivan Vulic. 2019. How to (properly) evaluate cross-lingual word embeddings: On strong baselines, comparative analyses, and some misconceptions. In *Proceedings of the Association for Computational Linguistics*.
- L.G. Gubin, B.T. Polyak, and E.V. Raik. 1967. The method of projections for finding the common point of convex sets. *USSR Computational Mathematics and Mathematical Physics*, 7(6):1–24.
- Jiang Guo, Wanxiang Che, David Yarowsky, Haifeng Wang, and Ting Liu. 2015. [Cross-lingual dependency parsing based on distributed representations](#). In *Proceedings of the Association for Computational Linguistics*.
- Yedid Hoshen and Lior Wolf. 2018. [Non-adversarial unsupervised word translation](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- Pratik Jawanpuria, Arjun Balgovind, Anoop Kunchukuttan, and Bamdev Mishra. 2019. [Learning multilingual word embeddings in latent metric space: a geometric approach](#). *Transactions of the Association for Computational Linguistics*, 7:107–120.
- Armand Joulin, Piotr Bojanowski, Tomas Mikolov, Hervé Jégou, and Edouard Grave. 2018. [Loss in translation: Learning bilingual word mapping with a retrieval criterion](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- Alexandre Klementiev, Ivan Titov, and Binod Bhattarai. 2012. [Inducing crosslingual distributed representations of words](#). *Proceedings of International Conference on Computational Linguistics*.
- Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. 2013a. [Exploiting similarities among languages for machine translation](#). *arXiv preprint arXiv:1309.4168*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013b. [Distributed representations of words and phrases and their compositionality](#). In *Proceedings of Advances in Neural Information Processing Systems*.
- George A. Miller and Walter G. Charles. 1991. [Contextual correlates of semantic similarity](#). *Language and Cognitive Processes*, 6(1):1–28.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of Empirical Methods in Natural Language Processing*.

- Herbert Rubenstein and John B Goodenough. 1965. [Contextual correlates of synonymy](#). *Communications of the ACM*, 8(10):627–633.
- Sebastian Ruder, Ryan Cotterell, Yova Kementchedjhiya, and Anders Søgaard. 2018. [A discriminative latent-variable model for bilingual lexicon induction](#). In *Proceedings of Empirical Methods in Natural Language Processing*.
- Sebastian Ruder, Ivan Vulić, and Anders Søgaard. 2017. [A survey of cross-lingual embedding models](#). *arXiv preprint arXiv:1706.04902*.
- Peter H. Schönemann. 1966. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10.
- Samuel L. Smith, David H. P. Turban, Steven Hamblin, and Nils Y. Hammerla. 2017. [Offline bilingual word vectors, orthogonal transformations and the inverted softmax](#). In *Proceedings of the International Conference on Learning Representations*.
- Anders Søgaard, Sebastian Ruder, and Ivan Vulić. 2018. [On the limitations of unsupervised bilingual dictionary induction](#). In *Proceedings of the Association for Computational Linguistics*.
- Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. 2015. [Normalized word embedding and orthogonal transform for bilingual word translation](#). In *Conference of the North American Chapter of the Association for Computational Linguistics*.
- Dongqiang Yang and David M. Powers. 2006. Verb similarity on the taxonomy of wordnet. In *International WordNet Conference*.
- Yuan Zhang, David Gaddy, Regina Barzilay, and Tommi Jaakkola. 2016. [Ten pairs to tag – multilingual POS tagging via coarse mapping between embeddings](#). In *Conference of the North American Chapter of the Association for Computational Linguistics*.
- Zhihui Zhu and Xiao Li. 2018. [Convergence analysis of alternating nonconvex projections](#). *arXiv preprint arXiv:1802.03889*.

A Proof for Theorem 1

Our convergence analysis is based on a recent result on alternating non-convex projections. Theorem 1 in the work of [Zhu and Li \(2018\)](#) states that the convergence of alternating projection holds even if the constraint sets are non-convex, as long as the two constraint sets satisfy the following assumption:

Assumption 1. *Let \mathbb{X} and \mathbb{Y} be any two closed semi-algebraic sets, and let $\{(\mathbf{x}_k, \mathbf{y}_k)\}$ be the sequence of iterates generated by the alternating projection method (e.g., Iterative Normalization). Assume the sequence $\{(\mathbf{x}_k, \mathbf{y}_k)\}$ is bounded and the sets \mathbb{X} and \mathbb{Y} obey the following properties:*

- (i) *three-point property of \mathbb{Y} : there exists a non-negative function $\delta_\alpha : \mathbb{Y} \times \mathbb{Y} \rightarrow \mathbb{R}$ with $\alpha > 0$ such that for any $k \geq 1$, we have*

$$\delta_\alpha(\mathbf{y}_k, \mathbf{y}_{k-1}) \geq \alpha \|\mathbf{y}_k - \mathbf{y}_{k-1}\|_2$$

and

$$\delta_\alpha(\mathbf{y}_{k-1}, \mathbf{y}_k) + \|\mathbf{x}_k - \mathbf{y}_k\|_2^2 \leq \|\mathbf{x}_k - \mathbf{y}_{k-1}\|_2^2,$$

- (ii) *local contraction property of \mathbb{X} : there exist $\epsilon > 0$ and $\beta > 0$ such that when $\|\mathbf{y}_k - \mathbf{y}_{k-1}\|_2 \leq \epsilon$, we have*

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 &= \|\mathcal{P}_{\mathbb{X}}(\mathbf{y}_k) - \mathcal{P}_{\mathbb{X}}(\mathbf{y}_{k-1})\|_2 \\ &\leq \beta \|\mathbf{y}_k - \mathbf{y}_{k-1}\|_2 \end{aligned}$$

where $\mathcal{P}_{\mathbb{X}}$ is the projection onto \mathbb{X} .

[Zhu and Li \(2018\)](#) only consider sets of vectors, but our constraint are sets of matrices. For ease of exposition, we treat every embedding matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$ as a vector by concatenating the column vectors: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$. The l^2 -norm of the concatenated vector $\|\mathbf{X}\|_2$ is equivalent to the Frobenius norm of the original matrix $\|\mathbf{X}\|_F$.

The two operations in Iterative Normalization, Equation (5) and (6), are projections onto two constraint sets, unit-length set $\mathbb{Y} = \{\mathbf{X} \in \mathbb{R}^{d \times n} : \forall i, \|\mathbf{x}_i\|_2 = 1\}$ and zero-mean set $\mathbb{X} = \{\mathbf{X} \in \mathbb{R}^{d \times n} : \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}\}$. To prove convergence of Iterative Normalization, we show that \mathbb{Y} satisfies the three-point property, and \mathbb{X} satisfies the local contraction property.

Three-point property of \mathbb{Y} . For any $\mathbf{Y}' \in \mathbb{Y}$ and $\mathbf{X} \in \mathbb{R}^{n \times d}$, let \mathbf{Y} be the projection of \mathbf{X} onto the constraint set \mathbb{Y} with Equation (5). The columns of \mathbf{Y} and \mathbf{Y}' have the same length, so we have

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}'\|_2^2 - \|\mathbf{X} - \mathbf{Y}\|_2^2 &= \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{y}'_i\|_2^2 - \|\mathbf{x}_i - \mathbf{y}_i\|_2^2 \\ &= \sum_{i=1}^n 2\mathbf{x}_i^\top \mathbf{y}_i - 2\mathbf{x}_i^\top \mathbf{y}'_i. \end{aligned} \quad (7)$$

Since \mathbf{Y} is the projection of \mathbf{X} onto the unit-length set with Equation (5); i.e., $\mathbf{y}_i = \mathbf{x}_i / \|\mathbf{x}_i\|_2$, we can rewrite Equation (7).

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}'\|_2^2 - \|\mathbf{X} - \mathbf{Y}\|_2^2 &= \sum_{i=1}^n \|\mathbf{x}_i\|_2 (2\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}'_i). \end{aligned} \quad (8)$$

All columns of \mathbf{Y} and \mathbf{Y}' are unit-length. Therefore, we can further rewrite Equation (8).

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}'\|_2^2 - \|\mathbf{X} - \mathbf{Y}\|_2^2 &= \sum_{i=1}^n \|\mathbf{x}_i\|_2 (2 - 2\mathbf{y}_i^\top \mathbf{y}'_i) \\ &= \sum_{i=1}^n \|\mathbf{x}_i\|_2 \|\mathbf{y}_i - \mathbf{y}'_i\|_2^2. \end{aligned}$$

Let $l = \min_i \{\|\mathbf{x}_i\|_2\}$ be the minimum length of the columns in \mathbf{X} . We have the following inequality:

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}'\|_2^2 - \|\mathbf{X} - \mathbf{Y}\|_2^2 &\geq \sum_{i=1}^n l \|\mathbf{y}_i - \mathbf{y}'_i\|_2^2 \\ &= l \|\mathbf{Y} - \mathbf{Y}'\|_2^2. \end{aligned}$$

From our non-zero assumption, the minimum column length l is always positive. Let l_k be the minimum column length of the embedding matrix $\mathbf{X}^{(k)}$ after the k -th iteration. It follows that \mathbb{Y} satisfies the three-point property with $\alpha = \min_k \{l_k\}$ and $\delta_\alpha(\mathbf{Y}, \mathbf{Y}') = \alpha \|\mathbf{Y} - \mathbf{Y}'\|_2^2$.

Local contraction property of \mathbb{X} . The zero-mean constraint set \mathbb{X} is convex and closed: if two matrices \mathbf{X} and \mathbf{Y} both have zero-mean, their linear interpolation $\lambda \mathbf{X} + (1 - \lambda) \mathbf{Y}$ must also have zero-mean for any $0 < \lambda < 1$. Projections onto convex

sets in a Hilbert space are contractive (Browder, 1967), and therefore \mathbb{X} satisfies the local contraction property with any positive ϵ and $\beta = 1$.

In summary, the two constraint sets that Iterative Normalization projects onto satisfy Assumption 1. Therefore, Iterative Normalization converges following the analysis of Zhu and Li (2018).

B Results on All Languages

Table 3 shows word translation accuracies on all target languages. Iterative Normalization improves accuracy on all languages.

Target	Procrustes			Procrustes + refine			RCSLS		
	None	C+L	IN	None	C+L	IN	None	C+L	IN
AF	26.3	28.3	29.7	27.7	28.7	30.4	9.3	28.6	29.3
AR	36.5	37.1	37.9	36.5	37.1	37.9	18.4	40.5	41.5
BS	22.3	23.5	24.4	23.3	23.9	26.6	5.4	25.5	26.6
CA	65.9	67.6	68.9	66.5	67.6	68.9	43.0	68.9	69.5
CS	54.0	54.7	55.3	54.0	54.7	55.7	29.9	57.8	58.2
DA	54.0	54.9	58.4	56.8	59.3	60.9	19.2	58.3	60.5
DE	73.5	74.6	75.5	74.3	75.2	76.0	43.6	77.5	78.1
EL	44.0	44.9	47.5	44.6	45.9	47.9	14.0	47.1	48.5
ES	81.4	81.3	81.5	81.9	82.1	82.5	50.5	83.6	83.9
ET	31.9	34.5	36.1	31.9	35.3	36.4	8.1	37.3	39.4
FA	33.1	33.7	37.3	33.1	34.1	37.3	5.9	37.5	38.3
FI	47.6	48.5	50.9	47.6	50.1	51.1	20.9	52.3	53.3
FR	81.1	81.3	81.7	82.1	82.7	82.4	53.1	83.9	83.9
HE	40.2	43.1	43.7	40.2	43.1	43.7	13.1	49.7	50.1
HI	33.3	34.0	36.7	33.6	34.9	37.7	5.0	36.2	38.0
HR	37.0	37.8	40.2	37.6	37.8	40.2	14.5	41.1	42.6
HU	51.8	54.1	55.5	53.3	54.1	56.1	11.7	57.3	58.2
ID	65.6	65.7	67.9	67.7	68.4	70.3	24.8	68.9	70.0
IT	76.2	76.6	76.6	77.5	78.1	78.1	48.4	78.8	79.1
JA	1.7	13.1	44.3	1.7	13.1	44.3	14.6	16.1	56.3
KO	31.5	32.1	33.9	31.5	32.1	33.9	6.4	37.5	37.5
LT	22.5	22.8	23.2	22.5	22.8	23.3	7.6	23.3	23.5
LV	23.6	24.9	26.1	23.6	24.9	26.1	10.1	28.3	28.7
MS	44.0	45.4	48.9	46.5	48.3	51.1	19.9	49.1	50.2
NL	72.8	73.7	74.1	73.8	75.1	75.8	46.7	75.6	75.8
PL	58.2	60.2	60.1	58.5	60.2	60.4	39.4	62.4	62.5
PT	79.5	79.7	79.9	79.9	81.0	81.2	63.1	81.1	81.7
RO	58.1	60.5	61.8	59.9	60.5	62.5	27.1	61.9	63.3
RU	51.7	52.1	52.1	51.7	52.1	52.1	26.6	57.1	57.9
SK	38.0	39.3	40.4	38.0	39.3	41.7	13.3	41.5	42.3
SL	32.5	34.3	36.7	32.5	34.4	36.7	12.3	36.0	37.9
SQ	23.5	25.1	27.3	23.5	25.1	27.3	4.4	26.5	27.3
SV	58.7	59.6	60.7	60.9	61.2	62.6	35.6	63.8	63.9
TA	15.1	15.5	16.8	15.1	15.5	17.7	6.7	16.3	17.1
TH	22.5	23.3	22.9	22.5	23.3	22.9	9.4	23.7	23.9
TR	44.9	46.5	48.7	46.3	48.7	51.7	18.3	50.7	52.4
UK	34.8	35.9	36.3	35.5	35.9	36.5	18.8	40.7	40.8
VI	41.3	42.1	43.7	42.1	42.7	44.2	14.2	43.3	43.9
ZH	32.5	42.3	44.2	32.5	42.3	44.2	17.1	45.1	48.6
Average	44.7	46.3	48.4	45.3	47.0	49.1	21.8	49.0	50.9

Table 3: Word translation accuracy aligning English embeddings to thirty-nine languages. We combine three normalizations—no normalization (None), mean centering and length normalization (C+L), and Iterative Normalization (IN) for five rounds—with three CLWEs: Procrustes, Procrustes with refinement (Conneau et al., 2018), and RCSLS (Joulin et al., 2018). Procrustes with C+L is equivalent to Artetxe et al. (2016). The best result for each CLWE in each column **in bold**. Iterative Normalization has the best accuracy of the three normalization techniques.