# Denoising Word Embeddings by Averaging in a Shared Space

**Avi Caciularu**[1]    **Ido Dagan**[1]    **Jacob Goldberger**[2]

[1]Computer Science Department, Bar-Ilan University
[2]Faculty of Engineering, Bar-Ilan University

avi.c33@gmail.com, dagan@cs.biu.ac.il, jacob.goldberger@biu.ac.il

## Abstract

We introduce a new approach for smoothing and improving the quality of word embeddings. We consider a method of fusing word embeddings that were trained on the same corpus but with different initializations. We project all the models to a shared vector space using an efficient implementation of the Generalized Procrustes Analysis (GPA) procedure, previously used in multilingual word translation. Our word representation demonstrates consistent improvements over the raw models as well as their simplistic average, on a range of tasks. As the new representations are more stable and reliable, there is a noticeable improvement in rare word evaluations.

## 1 Introduction

Continuous (non-contextualized) word embeddings have been introduced several years ago as a standard building block for NLP tasks. These models provide efficient ways to learn word representations in a fully self-supervised manner from text corpora, solely based on word co-occurrence statistics. A wide variety of methods now exist for generating word embeddings, with prominent methods including word2vec (Mikolov et al., 2013a), GloVe (Pennington et al., 2014), and FastText (Bojanowski et al., 2017). Recently, contextualized embeddings (Peters et al., 2018; Devlin et al., 2019), replaced the use of non-contextualized embeddings in many settings. Yet, the latter remain the standard choice for typical lexical-semantic tasks, e.g., semantic similarity (Hill et al., 2015), word analogy (Jurgens et al., 2012), relation classification (Barkan et al., 2020a), and paraphrase identification (Meged et al., 2020). These tasks consider the generic meanings of lexical items, given out of context, hence the use of non-contextualized embeddings is appropriate. Notably, FastText was shown to yield state-of-the-art results in most of these tasks (Bojanowski et al., 2017).

While word embedding methods proved to be powerful, they suffer from a certain level of noise, introduced by quite a few randomized steps in the embedding generation process, including embedding initialization, negative sampling, subsampling and mini-batch ordering. Consequently, different runs would yield different embedding geometries, of varying quality. This random noise might harm most severely the representation of rare words, for which the actual data signal is rather weak (Barkan et al., 2020b).

In this paper, we propose denoising word embedding models through generating multiple model versions, each created with different random seeds. Then, the resulting representations for each word should be fused effectively, in order to obtain a model with a reduced level of noise. Note, however, that simple averaging of the original word vectors is problematic, since each training session of the algorithm produces embeddings in a different space. In fact, the objective scores of both word2vec, Glove and FastText are invariant to multiplying all the word embeddings by an orthogonal matrix, hence, the algorithm output involves an arbitrary rotation of the embedding space.

For addressing this issue, we were inspired by recent approaches originally proposed for aligning multi-lingual embeddings (Chen and Cardie, 2018; Kementchedjhieva et al., 2018; Alaux et al., 2019; Jawanpuria et al., 2019; Taitelbaum et al., 2019). To obtain such alignments, these methods simultaneously project the original language-specific embeddings into a shared space, while enforcing (or at least encouraging) transitive orthogonal transformations. In our (monolingual) setting, we propose a related technique to project the different embedding versions into a shared space, while optimizing the projection towards obtaining an improved fused representation. We show that this results in im-

proved performance on a range of lexical-semantic tasks, with notable improvements for rare words, as well as on several sentence-level downstream tasks.

## 2 Word Averaging in a Shared Space

Assume we are given an ensemble of $k$ pre-trained word embedding sets, of the same word vocabulary of size $n$ and the same dimensionality $d$. In our setting, these sets are obtained by training the same embedding model using different random parameter initializations. Our goal is to fuse the $k$ embedding sets into a single "average" embedding that is hopefully more robust and would yield better performance on various tasks. Since each embedding set has its own space, we project the $k$ embedding spaces into a shared space, in which we induce averaged embeddings based on a mean squared error minimization objective.

Let $x_{i,t} \in \mathbb{R}^d$ be the dense representation of the $t$-th word in the $i$-th embedding set. We model the mapping from the $i$-th set to the shared space by an orthogonal matrix denoted by $T_i$. Denote the sought shared space representation of the $t$-th word by $y_t \in \mathbb{R}^d$. Our goal is to find a set of transformations $T = \{T_1, ..., T_k\}$ and target word embeddings $y = \{y_1, ..., y_n\}$ in the shared space that minimize the following mean-squared error:

$$S(T, y) = \sum_{i=1}^{k} \sum_{t=1}^{n} \|T_i x_{i,t} - y_t\|^2. \qquad (1)$$

For this objective, it is easy to show that for a set of transformations $T_1, ..., T_k$, the optimal shared space representation is:

$$y_t = \frac{1}{k} \sum_{i=1}^{k} T_i x_{i,t}.$$

Hence, solving the optimization problem pertains to finding the $k$ optimal transformations.

In the case where $k = 2$, the optimal $T$ can be obtained in a closed form using the Procrustes Analysis (PA) procedure (Schönemann, 1966), which has been employed in recent bilingual word translation methods (Xing et al., 2015; Artetxe et al., 2016; Hamilton et al., 2016; Artetxe et al., 2017a,b; Conneau et al., 2017; Artetxe et al., 2018a,b; Ruder et al., 2018). In our setting, to obtain an improved embedding, we wish to average more than two embedding sets.

However, if $k > 2$ there is no closed form solution to (1) and thus, we need to find a solution using an iterative optimization process. To that end, we follow several works that suggested employing the General Procrustes Analysis (GPA) procedure, which is an extension of PA to multi-set alignment (Gower, 1975; Kementchedjhieva et al., 2018). Generally, the GPA consists of an alternate minimization procedure where we iterate between finding the orthogonal transformations and computing the shared space. The optimal transformation from each embedding space to the shared space is found by minimizing the following score,

$$S(T_i) = \sum_{t=1}^{n} \|T_i x_{i,t} - y_t\|^2, \qquad i = 1, ..., k.$$

The minimum of $S(T_i)$ can then be found by the closed form PA procedure. The updated transformation is $T_i = U_i V_i^\top$, where $U_i \Sigma_i V_i^\top$ is the singular value decomposition (SVD) of the $d \times d$ matrix $\sum_{t=1}^{n} y_t x_{i,t}^\top$. At each step in the iterative GPA algorithm, the score (1) is monotonically decreased until it converges to a local minimum point.

---

**Algorithm 1** Shared Space Embedding Averaging

1: **Input:** Ensemble of $k$ word embedding sets.
2: **Task:** Find the optimal average embedding.
3: **Preprocessing:**
4: Compute the cross-correlation matrices:
5: $C_{ij} = C_{ji}^\top = \sum_{t=1}^{n} x_{j,t} x_{i,t}^\top \quad 1 \le i < j \le k$
6: **Initialization:** $T_1 = \cdots = T_{k-1} = 0, T_k = I$
7: **while** not converged **do**
8:     **for** $i = 1, ..., k$ **do**
9:         $U\Sigma V^\top = \text{SVD} \left( \sum_{j \ne i} T_j C_{ij} \right)$
10:         $T_i \leftarrow UV^\top$
11:     **end for**
12: **end while**
13: **Compute the average embedding:**
14: $y_t \leftarrow \frac{1}{k} \sum_{i=1}^{k} T_i x_{i,t} \qquad t = 1, ..., n$

---

For large vocabularies, GPA is not efficient, because, in each iteration, when computing the SVD we need to sum over all the vocabulary words. To circumvent this computational cost, we adopt the optimization procedure from Taitelbaum et al. (2019), which we apply within each iteration. Instead of summing over the whole vocabulary, the following extension is proposed. Let $C_{ij} = \sum_t x_{j,t} x_{i,t}^\top$ be the cross-correlation matrix

|          | original          | denoised           |
|----------|-------------------|--------------------|
| word2vec | $0.40 \pm 0.005$  | $0.059 \pm 0.003$  |
| GloVe    | $0.38 \pm 0.006$  | $0.058 \pm 0.003$  |
| FastText | $0.35 \pm 0.003$  | $0.054 \pm 0.001$  |

Table 1: Average MSE scores of the embedding models with and without applying the SSEA algorithm.

for a pair $(i, j)$ of two original embedding spaces, which can be computed once, for all pairs of spaces, in a pre-processing step. Given the matrices $C_{ij}$ the computational complexity of the iterative averaging algorithm is independent of the vocabulary size, allowing us to compute efficiently the SVD. The resulting algorithm termed Shared Space Embedding Averaging (SSEA) is presented in Algorithm 1.[1]

## 3 Experimental Setup and Results

This section presents our evaluation protocol, datasets, data preparation, hyperparameter configuration and results.

### 3.1 Implementation Details and Data

We trained word2vec (Mikolov et al., 2013a), Fast-Text (Bojanowski et al., 2017) and GloVe (Pennington et al., 2014) embeddings. For word2vec we used the skip-gram model with negative sampling, which was shown advantageous on the evaluated tasks (Levy et al., 2015). We trained each of the models on the November 2019 dump of Wikipedia articles[2] for $k = 30$ times, with different random seeds, and used the default reported hyperparameters; we set the embedding dimension to $d = 200$, and considered each word within the maximal window $c_{max} = 5$, subsampling[3] threshold of $\rho = 10^{-5}$ and used 5 negative examples for every positive example. In order to keep a large amount of rare words in the corpus, no preprocessing was applied on the data, yielding a vocabulary size of $1.5 \cdot 10^6$. We then applied the SSEA algorithm to the embedding sets to obtain the average embedding. The original embedding sets and averaged embeddings were centered around the 0 vector and normalized to unit vectors.

### 3.2 Improved Embedding Stability

We next analyze how our method improves embedding quality and consistency, notably for rare
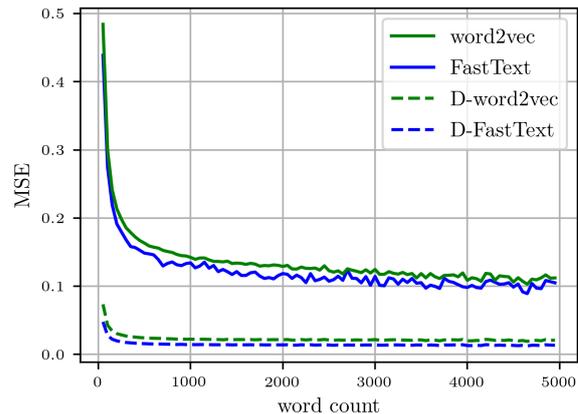


Figure 1: Average MSE for word embeddings vs their corpus occurrence count (binned with resolution of 50).

words. To that end, for any two embedding sets, $u$ and $v$, we can find the optimal mapping $Q$ between them using the PA algorithm and compute its mean square error (MSE), $\frac{1}{n} \sum_{t=1} \|Qu_t - v_t\|^2$. We define the stability of an embedding algorithm by the average MSE (over 10 random pairs of samples) between two instances of it. This score measures the similarity between the geometries of random instances generated by a particular embedding method, and thus reflects the consistency and stability of that method. The scores of the different models are depicted in Table 1. As observed, after applying SSEA the Average MSE drops by an order of magnitude, indicating much better stability of the obtained embeddings.

We can perform a similar analysis for each word separately. A consistent embedding of the $t$-th word in both sets $u$ and $v$ should result in a small mapping discrepancy $\|Qu_t - v_t\|^2$. Figure 1 depicts MSE for the models and their computed SSEA, as a function of the word's frequency in the corpus. The denoised version of the models is marked with a 'D-' prefix. For clarity of presentation, we did not include the results for GloVe (which are similar to word2vec). As expected, embedding stability always increases (MSE decreases) with word frequency. SSEA is notably more stable across the frequency range, with the error minimized early on and reduced most drastically for low frequencies.

### 3.3 Comparison of methods

We next compare our denoised model, denoted with a 'D-' prefix, with the original embedding models. As an additional baseline, we considered also the naïve averaged embedding model, denoted with a 'A-' prefix, where for every word we computed the simplistic mean embedding across all origi-

---

[1] The algorithm demonstration code is available at github.com/aviclu/SSEA. In practice, we utilized an efficient PyTorch implementation based on Taitelbaum et al. (2019).

[2] dumps.wikimedia.org/enwiki/latest/

[3] To speed up the training.

| Method | SimLex999 | MEN | WS353 | AP | Google | MSR | SemEval2012(2) | BLESS | RW |
|---|---|---|---|---|---|---|---|---|---|
| word2vec | 33.7 | 72.4 | 60.7 | 62.2 | 69.5 | 51.3 | 19.2 | 79 | 42 |
| A-word2vec | 33.1 | 72.3 | **60.8** | 61.9 | 69.2 | 51.2 | 19 | 78.1 | 41.7 |
| D-word2vec | **33.9** | **73.2** | **60.8** | **63.1** | **70.3** | **51.9** | **20** | **79.6** | **43.4** |
| GloVe | 34.4 | 73.4 | 62.3 | 63.3 | 75.1 | 54.5 | 19.7 | 79.2 | 47.2 |
| A-GloVe | 34.2 | 73.1 | 61.9 | 62.8 | 74.7 | 54.2 | 19.6 | 79 | 47.1 |
| D-GloVe | **34.8** | **75.1** | **62.7** | **64.3** | **75.9** | **55.2** | **20.1** | **79.9** | **48.5** |
| FastText | 41.2 | 78.6 | 70.7 | 72.2 | 75.7 | 63.4 | 19.8 | 81.5 | 47.1 |
| A-FastText | 41 | 78.1 | 69.7 | 72.1 | 74.1 | 62.8 | 19.4 | 80.8 | 46.6 |
| D-FastText | **42.2** | **79.3** | **71.8** | **72.9** | **77.4** | **63.8** | **20.2** | **82.7** | **50.3** |

Table 2: Results for lexical-semantic benchmarks. Best performance is bolded.

nal spaces. Note that we did not compare other proposed embeddings or meta-embedding learning methods, but rather restricted our analysis to empirically verifying our embedding aggregation method and validating the assumptions behind the empirical analysis we performed.

### 3.4 Evaluations on Lexical Semantic Tasks

We evaluated the performance of our method over lexical-semantic tasks, including word similarity, analogy solving, and concept categorization: **SimLex999** (Hill et al., 2015), **MEN** (Bruni et al., 2014), **WS353** (Finkelstein et al., 2002), **AP** (Almuhareb and Poesio, 2004), **Google** (Mikolov et al., 2013b), **MSR** (Mikolov et al., 2013c), **SemEval-2012** (Jurgens et al., 2012), **BLESS** (Baroni and Lenci, 2011) and **RW** (Luong et al., 2013), (focusing on rare words). For the analogy task, we reported the accuracy. For the remaining tasks, we computed Spearman's correlation between the cosine similarity of the embeddings and the human judgments.

**Results** The results of the lexical-semantic tasks are depicted in Table 2, averaged over 30 runs for each method. Our method obtained better performance than the other methods, substantially for FastText embeddings. As shown, the naïve averaging performed poorly, which highlights the fact that simply averaging different embedding spaces does not improve word representation quality. The most notable performance gain was in the rare-words task, in line with the analysis in Fig. 1, suggesting that on rare words the raw embedding vectors fit the data less accurately.

### 3.5 Evaluations On Downstream Tasks

For completeness, we next show the relative advantage of our denoising method also when applied to several sentence-level downstream benchmarks. While contextualized embeddings domi-

nate a wide range of sentence- and document- level NLP tasks (Peters et al., 2018; Devlin et al., 2019; Caciularu et al., 2021), we assessed the relative advantage of our denoising method when utilizing (non-contextualized) word embeddings in sentence- an document- level settings. We applied the exact procedure proposed in Li et al. (2017) and Rogers et al. (2018), as an effective benchmark for the quality of static embedding models. We first used sequence labeling tasks. The morphological and syntactic performance was evaluated using part-of-speech tagging, **POS**, and chunking, **CHK**. Both named entity recognition, **NER**, and multi-way classification of semantic relation classes, **RE**, tasks were used for evaluating semantic information at the word level. For the above POS, NER and CHK sequence labeling tasks, we used the CoNLL 2003 dataset (Sang and Meulder, 2003) and for the RE task, we used the SemEval 2010 task 8 dataset (Hendrickx et al., 2010). The neural network models employed for these downstream tasks are fully described in (Rogers et al., 2018). Next, we evaluated the following semantic level tasks: document-level polarity classification, **PC**, using the Stanford IMDB movie review dataset (Maas et al., 2011), sentence level sentiment polarity classification, **SEN**, using the MR dataset of short movie reviews (Pang and Lee, 2005), and classification of subjectivity and objectivity task, **SUB**, that uses the Rotten Tomatoes user review snippets against official movie plot summaries (Pang and Lee, 2004). Similarly to the performance results in Table 2, the current results show that the suggested denoised embeddings obtained better overall performance than the other methods, substantially for FastText embeddings.

## 4 Related Work

A similar situation of aligning different word embeddings into a shared space occurs in multi-lingual

| Method | POS | CHK | NER | RE | PC | SEN | SUB |
|---|---|---|---|---|---|---|---|
| word2vec | 81.5 | 80.1 | 93.3 | 71.4 | 89.2 | 73.9 | 76.4 |
| A-word2vec | 78 | 77.5 | 90.9 | 67.4 | 86.4 | 64.3 | 75.6 |
| D-word2vec | **81.6** | **80.2** | **93.6** | **73.1** | **89.7** | **74** | **77.4** |
| GloVe | 77.5 | 70.4 | 85.2 | 66.7 | 80.2 | 70.2 | 72.7 |
| A-GloVe | 77.1 | 70.2 | 84.9 | 62.3 | 77.7 | 62.2 | 71.8 |
| D-GloVe | **77.8** | **71.1** | **86.6** | **68.2** | **80.8** | **71.3** | **73.9** |
| FastText | 80.6 | 79.1 | 92.2 | 74 | 88.9 | 74.9 | 73.9 |
| A-FastText | 78.4 | 78.8 | 90.2 | 73.6 | 89 | 74.1 | 73.3 |
| D-FastText | **82.4** | **81.2** | **94.9** | **75.2** | **90.5** | **77.3** | **76.7** |

Table 3: Results for downstream task. Best performance is bolded.

word translation tasks which are based on distinct monolingual word embeddings. Word translation is performed by transforming each language word embeddings into a shared space by an orthogonal matrix, for creating a "universal language", which is useful for the word translation process. Our setting may be considered by viewing each embedding set as a different language, where our goal is to find the shared space where embedding averaging is meaningful.

The main challenge in multilingual word translation is to obtain a reliable multi-way word correspondence in either a supervised or unsupervised manner. One problem is that standard dictionaries contain multiple senses for words, which is problematic for bilingual translation, and further amplified in a multilingual setting. In our case of embedding averaging, the mapping problem vanishes since we are addressing a single language and the word correspondences hold trivially among different embeddings of the same word. Thus, in our setting, there are no problems of wrong word correspondences, neither the issue of having different word translations due to multiple word senses. Studies have shown that for the multi-lingual translation problem, enforcing the transformation to be strictly orthogonal is too restrictive and performance can be improved by using the orthogonalization as a regularization (Chen and Cardie, 2018) that yields matrices that are close to be orthogonal. In our much simpler setting of a single language, with a trivial identity word correspondence, enforcing the orthogonalization constraint is reasonable.

Another related problem is *meta-embedding* (Yin and Schütze, 2016), which aims to fuse information from different embedding models. Various methods have been proposed for embedding fusion, such as concatenation, simple averaging, weighted averaging (Coates and Bollegala, 2018;

Kiela et al., 2018) and autoencoding (Bollegala and Bao, 2018). Some of these methods (concatenation and autoencoding) are not scalable when the goal is to fuse many sets, while others (simple averaging) yield inferior results, as described in the above works. Note that our method is not intended to be a competitor of meta-embedding, but rather a complementary method.

An additional related work is the recent method from (Muromägi et al., 2017). Similarly to our work, they proposed a method based on the Procrustes Analysis procedure for aligning and averaging sets of word embedding models. However, the mapping algorithm they used is much more computationally demanding, as it requires to go over all the dictionary words in every iteration. Instead, we propose an efficient optimization algorithm, which requires just one such computation during each iteration, and is theoretically guaranteed to converge to a local minimum point. While their work focuses on improving over the Estonian language, we suggest evaluating this approach on English data and on a range of different downstream tasks. We show that our method significantly improves upon rare words, which is beneficial for small sized / domain-specific corpora.

## 5 Conclusions

We presented a novel technique for creating better word representations by training an embedding model several times, from which we derive an averaged representation. The resulting word representations proved to be more stable and reliable than the raw embeddings. Our method exhibits performance gains in lexical-semantic tasks, notably over rare words, confirming our analytical assumptions. This suggests that our method may be particularly useful for training embedding models in low-resource settings. Appealing future research may extend our approach to improving sentence-level representations, by fusing several contextualized embedding models.

## Acknowledgments

# References

Jean Alaux, Edouard Grave, Marco Cuturi, and Armand Joulin. 2019. Unsupervised hyperalignment for multilingual word embeddings. In *International Conference on Learning Representations (ICLR)*.

Abdulrahman Almuhareb and Massimo Poesio. 2004. Attribute-based and value-based clustering: An evaluation. In *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2016. Learning principled bilingual mappings of word embeddings while preserving monolingual invariance. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2017a. Learning bilingual word embeddings with (almost) no bilingual data. In *Annual Meeting of the Association for Computational Linguistics (ACL)*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018a. Generalizing and improving bilingual word embedding mappings with a multi-step framework of linear transformations. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*.

Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018b. A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. *arXiv preprint arXiv:1805.06297*.

Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2017b. Unsupervised neural machine translation. *arXiv preprint arXiv:1710.11041*.

Oren Barkan, Avi Caciularu, and Ido Dagan. 2020a. Within-between lexical relation classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3521–3527, Online. Association for Computational Linguistics.

Oren Barkan, Idan Rejwan, Avi Caciularu, and Noam Koenigstein. 2020b. Bayesian hierarchical words representation learning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3871–3877, Online. Association for Computational Linguistics.

Marco Baroni and Alessandro Lenci. 2011. How we blessed distributional semantic evaluation. In *Proceedings of the GEMS 2011 Workshop on GEometrical Models of Natural Language Semantics*.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics (TACL)*.

Danushka Bollegala and Cong Bao. 2018. Learning word meta-embeddings by autoencoding. In *International Conference on Computational Linguistics, (COLING)*.

Elia Bruni, Nam-Khanh Tran, and Marco Baroni. 2014. Multimodal distributional semantics. *Journal of Artificial Intelligence Research*.

Avi Caciularu, Arman Cohan, Iz Beltagy, Matthew E Peters, Arie Cattan, and Ido Dagan. 2021. Cross-document language modeling. *arXiv e-prints*, pages arXiv–2101.

Xilun Chen and Claire Cardie. 2018. Unsupervised multilingual word embeddings. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Joshua Coates and Danushka Bollegala. 2018. Frustratingly easy meta-embedding – computing metaembeddings by averaging source word embeddings. In *North American Chapter of the Association for Computational Linguistics (NAACL)*.

Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

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*.

John C Gower. 1975. Generalized procrustes analysis. *Psychometrika*.

William L Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. *arXiv preprint arXiv:1605.09096*.

Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 33–38, Uppsala, Sweden. Association for Computational Linguistics.

Felix Hill, Roi Reichart, and Anna Korhonen. 2015. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*.

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 (TACL)*.

David A Jurgens, Peter D Turney, Saif M Mohammad, and Keith J Holyoak. 2012. Semeval-2012 task 2: Measuring degrees of relational similarity. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*.

Yova Kementchedjhieva, Sebastian Ruder, Ryan Cotterell, and Anders Søgaard. 2018. Generalizing procrustes analysis for better bilingual dictionary induction. In *The SIGNLL Conference on Computational Natural Language Learning (CoNLL)*.

Douwe Kiela, Changhan Wang, and Kyunghyun Cho. 2018. Dynamic meta-embeddings for improved sentence representations. In *Proceedings of thr Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics (TACL)*.

Bofang Li, Tao Liu, Zhe Zhao, Buzhou Tang, Aleksandr Drozd, Anna Rogers, and Xiaoyong Du. 2017. Investigating different syntactic context types and context representations for learning word embeddings. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2421–2431, Copenhagen, Denmark. Association for Computational Linguistics.

Thang Luong, Richard Socher, and Christopher Manning. 2013. Better word representations with recursive neural networks for morphology. In *Proceedings of Conference on Computational Natural Language Learning*.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

Yehudit Meged, Avi Caciularu, Vered Shwartz, and Ido Dagan. 2020. Paraphrasing vs coreferring: Two sides of the same coin. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4897–4907, Online. Association for Computational Linguistics.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013a. Distributed representa-

tions of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems (NIPS)*.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems (NIPS)*.

Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013c. Linguistic regularities in continuous space word representations. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*.

Avo Muromägi, Kairit Sirts, and Sven Laur. 2017. Linear ensembles of word embedding models. In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, pages 96–104, Gothenburg, Sweden. Association for Computational Linguistics.

Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 271–278, Barcelona, Spain.

Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.

Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.

Anna Rogers, Shashwath Hosur Ananthakrishna, and Anna Rumshisky. 2018. What's in your embedding, and how it predicts task performance. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2690–2703, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Sebastian Ruder, Ryan Cotterell, Yova Kementchedjhieva, and Anders Søgaard. 2018. A discriminative latent-variable model for bilingual lexicon induction. *arXiv preprint arXiv:1808.09334*.

Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*.

Peter Schönemann. 1966. A generalized solution of the orthogonal procrustes problem. *Psychometrika*.

Hagai Taitelbaum, Gal Chechik, and Jacob Goldberger. 2019. A multi-pairwise extension of procrustes analysis for multilingual word translation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

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: Human Language Technologies (NAACL)*.

Wenpeng Yin and Hinrich Schütze. 2016. Learning word meta-embeddings. In *Annual Meeting of the Association for Computational Linguistics (ACL)*.