

# Finding Word Sense Embeddings Of Known Meaning

Lyndon White, Roberto Togneri, Wei Liu, and Mohammed Bennisamoun

The University of Western Australia,  
35 Stirling Highway, Crawley, Western Australia  
lyndon.white@research.uwa.edu.au, roberto.togneri@uwa.edu.au,  
wei.liu@uwa.edu.au, mohammed.bennisamoun@uwa.edu.au

**Abstract.** Word sense embeddings are vector representations of polysemous words – words with multiple meanings. These induced sense embeddings, however, do not necessarily correspond to any dictionary senses of the word. To overcome this, we propose a method to find new sense embeddings with known meaning. We term this method refitting, as the new embedding is fitted to model the meaning of a target word in an example sentence. The new lexically refitted embeddings are learnt using the probabilities of the existing induced sense embeddings, as well as their vector values. Our contributions are threefold: (1) The refitting method to find the new sense embeddings; (2) a novel smoothing technique, for use with the refitting method; and (3) a new similarity measure for words in context, defined by using the refitted sense embeddings. We show how our techniques improve the performance of the Adaptive Skip-Gram sense embeddings for word similarity evaluation; and how they allow the embeddings to be used for lexical word sense disambiguation.

## 1 Introduction

Popular word embedding vectors, such as Word2Vec, represent a word’s semantic meaning and its syntactic role as a point in a vector space [1, 2]. As each word is only given one embedding, such methods are restricted to the representation of only a single combined sense, or meaning, of the word. *Word sense embeddings* generalise word embeddings to handle polysemous and homonymous words. Often these sense embeddings are learnt through unsupervised Word Sense Induction (WSI) [3–6]. The induced sense embeddings are unlikely to directly coincide with any set of human defined meaning at all, i.e. they will not match lexical senses such as those defined in a lexical dictionary, e.g. WordNet [7]. These induced senses may be more specific, more broad, or include the meanings of jargon not in common use.

One may argue that WSI systems can capture better word senses than human lexicographers do manually. However, this does not mean that induced senses can replace standard lexical senses. It is important to appreciate the vast wealth of existing knowledge defined around lexical senses. Methods to link induced senses to lexical senses allow us to take advantage of both worlds.

We propose a *refitting method* to generate a sense embedding vector that matches with a labelled lexical sense. Given an example sentence with the labelled lexical sense of a particular word, the refitting method algorithmically combines the induced sense embeddings of the target word such that the likelihood of the example sentence is maximised. We find that in doing so, the sense of the word in that sentence is captured. With the refitting, the induced sense embeddings are now able to be used in more general situations where standard senses, or user defined senses are desired.

Refitting word sense vectors to match a lexicographical sense inventory, such as WordNet or a translator’s dictionary, is possible if the sense inventory features at least one example of the target sense’s use. Our method allows this to be done very rapidly, and from only the single example of use this has with possible applications in low-resource languages.

Refitting can also be used to fit to a user provided example, giving a specific sense vector for that use. This has strong applications in information retrieval. The user can provide an example of a use of the word they are interested in. For example, searching for documents about “**banks**” as in “the river banks were very muddy”. By generating an embedding for that specific sense, and by comparing with the generated embeddings in the indexed documents, we can not only pick up on suitable uses of other-words for example “**beach**” and “**shore**”, but also exclude different usages, for example of a financial bank. The method we propose, using our refitted embeddings, has lower time complexity than AvgSimC [3], the current standard method for evaluating the similarity of words in context. This is detailed in Section 5.1.

We noted during refitting, that a single induced sense would often dominate the refitted representation. It is rare in natural language for the meaning to be so unequivocal. Generally, a significant overlap exists between the meaning of different lexical senses, and there is often a high level of disagreement when humans are asked to annotate a corpus [8]. We would expect that during refitting there would likewise be contention over the most likely induced sense. Towards this end, we develop a smoothing method, which we call *geometric smoothing* that de-emphasises the sharp decisions made by the (unsmoothed) refitting method. We found that this significantly improves the results. This suggests that the sharpness of sense decisions is an issue with the language model, which smoothing can correct. The geometric smoothing method is presented in Section 3.2.

We demonstrate the refitting method on sense embedding vectors induced using Adaptive Skip-Grams (AdaGram) [6], as well as our own simple greedy word sense embeddings. The method is applicable to any skip-gram-like language model that can take a sense vector as its input, and can output the probability of a word appearing in that sense’s context.

The rest of the paper is organised as follows: Section 2 briefly discusses two areas of related works. Section 3 presents our refitting method, as well as our proposed geometric smoothing method. Section 4 describes the WSI embedding models used in the evaluations. Section 5 defines the RefittedSim measure for word similarity in context, and presents its results. Section 6 shows how the refitted sense vectors can be used for lexical WSD. Finally, the paper concludes in Section 7.

## 2 Related Works

### 2.1 Directly Learning Lexical Sense Embeddings

In this area of research, the induction of word sense embeddings is treated as a supervised, or semi-supervised task, that requires sense labelled corpora for training.

Iacobacci et al. [9] use a Continuous Bag of Word language model [1], using word senses as the labels rather than words. This is a direct application of word embedding techniques. To overcome the lack of a large sense labelled corpus, Iacobacci et al. use a 3rd party WSD tool, BabelFly [10], to add sense annotations to a previously unlabelled corpus.

Chen et al. [11] use a supervised approach to train sense vectors, with an unsupervised WSD labelling step. They partially disambiguate their training corpus, using word sense vectors based on WordNet; and use these labels to train their embeddings. This relabelled data is then used as training data, for finding sense embeddings using skip-grams.

Our refitting method learns a new sense embedding as a weighted sum of existing induced sense embeddings of the target word. Refitting is a one-shot learning solution, as compared to the approaches used in the works discussed above. A notable advantage is the time taken to add a new sense. Adding a new sense is practically instantaneous, and replacing the entire sense inventory, of several hundred thousand senses, is only a matter of a few hours. Whereas for the existing approaches this would require repeating the training process, which will often take several days. Refitting is a process done to word sense embeddings, rather than a method for finding sense embeddings from a large corpus.

### 2.2 Mapping induced senses to lexical senses

By defining a stochastic map between the induced and lexical senses, Agirre et al. [12], propose a general method for allowing WSI systems to be used for WSD. Their work was used in SemEval-2007 Task 02 [13] to evaluate all entries. Agirre et al. use a mapping corpus to find the probability of a lexical sense, given the induced sense according to the WSI system. This is more general than the approach we propose here, which only works for sense embedding based WSI. By exploiting the particular properties of sense embedding based WSI systems we propose a system that can better facilitate the use of this subset of WSI systems for WSD.

## 3 Proposed Refitting Framework

The key contribution of this work is to provide a way to synthesise a word sense embedding given only a single example sentence and a set of pretrained sense embedding vectors. We termed this *refitting* the sense vectors. By refitting the unsupervised vectors we define a new vector, that lines up with the specific meaning of the word from the example sentence.

This can be looked at as a one-shot learning problem, analogous to regression. The training of the induced sense, and of the language model, can be considered

an unsupervised pre-training step. The new word sense embedding should give a high value for the likelihood of the example sentence, according to the language model. It should also generalise to give a high likelihood of other contexts where this word sense occurs.

We initially attempted to directly optimise the sense vector to predict the example. We applied the L-BFGS [14] optimisation algorithm with the sense vector being the parameter being optimised over, and the objective being to maximise the probability of the example sentence according to the language model. This was found to generalise poorly, due to over-fitting, and to be very slow. Rather than a direct approach, we instead take inspiration from the locally linear relationship between meaning and vector position that has been demonstrated for word embeddings [1, 15, 16].

To refit the induced sense embeddings to a particular meaning of a word, we express that a new embedding as a weighted combination of the induced sense vectors. The weight is determined by the probability of each induced sense given the context.

Given a collection of induced (unlabelled) embeddings  $\mathbf{u} = u_1, \dots, u_{n_u}$ , and an example sentence  $\mathbf{c} = w_1, \dots, w_{n_c}$  we define a function  $l(\mathbf{u} | \mathbf{c})$  which determines the refitted sense vector, from the unsupervised vectors and the context as:

$$l(\mathbf{u} | \mathbf{c}) = \sum_{\forall u_i \in \mathbf{u}} u_i P(u_i | \mathbf{c}) \quad (1)$$

Bayes' Theorem can be used to estimate the posterior predictive distribution  $P(u_i | \mathbf{c})$ .

Bengio et al. [17] describe a similar method to Equation (1) for finding (single sense) word embeddings for words not found in their vocabulary. The formula they give is as per Equation (1), but summing over the entire vocabulary of words (rather than just  $\mathbf{u}$ ).

### 3.1 A General WSD method

Using the language model and application of Bayes' theorem, we define a general word sense disambiguation method that can be used for refitting (Equation (1)), and for lexical word sense disambiguation (see Section 6). This is a standard approach of using Bayes' theorem [5, 6]. We present it here for completeness.

The context is used to determine which sense is the most suitable for this use of the *target word* (the word being disambiguated). Let  $\mathbf{s} = (s_1, \dots, s_n)$ , be the collection of senses for the target word<sup>1</sup>.

Let  $\mathbf{c} = (w_1, \dots, w_{n_c})$  be a sequence of words making up the context of the target word. For example for the target word *kid*, the context could be  $\mathbf{c} = (\text{wow the wool from the, is, so, soft, and, fluffy})$ , where *kid* is the central word taken from between *the* and *fluffy*.

<sup>1</sup> As this part of our method is used with both the unsupervised senses and the lexical senses, referred to as  $\mathbf{u}$  and  $\mathbf{l}$  respectively in other parts of the paper, here we use a general sense  $\mathbf{s}$  to avoid confusion.

For any particular sense,  $s_i$ , the multiple sense skip-gram language model can be used to find the probability of a word  $w_j$  occurring in the context:  $P(w_j | s_i)$ . By assuming the conditional independence of each word  $w_j$  in the context, given the sense embedding  $s_i$ , the probability of the context can be calculated:

$$P(\mathbf{c} | s_i) = \prod_{\forall w_j \in \mathbf{c}} P(w_j | s_i) \quad (2)$$

The correctness of the conditional independence assumption depends on the quality of the representation – the ideal sense representation would fully capture all information about the contexts it can appear in – thus the other contexts elements would not present any additional information, and so  $P(w_a | w_b, s_i) = P(w_a | s_i)$ . Given this, we have an estimate of  $P(\mathbf{c} | s_i)$  which can be used to find  $P(s_i | \mathbf{c})$ . However, a false assumption of independence contributes towards overly sharp estimates of the posterior distribution [18], which we seek to address in Section 3.2 with geometric smoothing.

Bayes’ Theorem is applied to this context likelihood function  $P(\mathbf{c} | s_i)$  and a prior for the sense  $P(s_i)$  to allow the posterior probability to be found:

$$P(s_i | \mathbf{c}) = \frac{P(\mathbf{c} | s_i)P(s_i)}{\sum_{s_j \in \mathbf{s}} P(\mathbf{c} | s_j)P(s_j)} \quad (3)$$

This is the probability of the sense given the context.

### 3.2 Geometric Smoothing for General WSD

During refitting, we note that often one induced sense would be calculated as having much higher probability of occurring than the others (according to Equation (3)). This level of certainty is not expected to occur in natural languages, ambiguity is almost always possible. To resolve such dominance problems, we propose a new *geometric smoothing* method. This is suitable for smoothing posterior probability estimates derived from products of conditionally independent likelihoods. It smooths the resulting distribution, by shifting all probabilities to be closer to the uniform distribution.

We hypothesize that the sharpness of probability estimates from Equation (3) is a result of data sparsity, and of a false independence assumption in Equation (2). This is well known to occur for n-gram language models [18]. Word-embeddings language models largely overcome the data sparsity problem due to weight sharing effects [17]. We suggest that the problem remains for word sense embeddings, where there are many more classes. Thus the training data must be split further between each sense than it was when split for each word. The power law distribution of word use [19] is compounded by word senses within those used also following the a power law distribution [20]. Rare senses are liable to over-fit to the few contexts they do occur in, and so give disproportionately high likelihoods to contexts that those are similar to. We propose to handle these issues through additional smoothing.

We consider replacing the unnormalised posterior with its  $n_c$ -th root, where  $n_c$  is the length of the context. We replace the likelihood of Equation (2) with

$P_S(\mathbf{c} | s_i) = \prod_{w_j \in \mathbf{c}} \sqrt[n_c]{P(w_j | s_i)}$ . Similarly, we replace the prior with:  $P_S(s_i) = \sqrt[n_c]{P(w_j | s_i)}$ . When this is substituted into Equation (3), it becomes a smoothed version of  $P(s_i | \mathbf{c})$ .

$$P_S(s_i | \mathbf{c}) = \frac{\sqrt[n_c]{P(\mathbf{c} | s_i)P(s_i)}}{\sum_{s_j \in \mathbf{s}} \sqrt[n_c]{P(\mathbf{c} | s_j)P(s_j)}} \quad (4)$$

The motivation for taking the  $n_c$ -th root comes from considering the case of the uniform prior. In this case  $P_S(\mathbf{c} | s_i)$  is the geometric mean of the individual word probabilities  $P_S(w_j | s_i)$ . Consider, if one has two context sentences,  $\mathbf{c} = \{w_1, \dots, w_{n_c}\}$  and  $\mathbf{c}' = \{w'_1, \dots, w'_{n_{c'}}\}$ , such that  $n_{c'} > n_c$  then using Equation (2) to calculate  $P(\mathbf{c} | s_i)$  and  $P(\mathbf{c}' | s_i)$  will result in incomparable results as additional number of probability terms will dominate – often significantly more than the relative values of the probabilities themselves. The number of words that can occur in the context of any given sense is very large – a large portion of the vocabulary. We would expect, averaging across all words, that each addition word in the context would decrease the probability by a factor of  $\frac{1}{V}$ , where  $V$  is the vocabulary size. The expected probabilities for  $P(\mathbf{c} | s_i)$  is  $\frac{1}{V^{n_c}}$  and for  $P(\mathbf{c}' | s_i)$  is  $\frac{1}{V^{n_{c'}}}$ . As  $n_{c'} > n_c$ , thus we expect  $P(\mathbf{c}' | s_i) \ll P(\mathbf{c} | s_i)$ . Taking the  $n_c$ -th and  $n_{c'}$ -th roots of  $P(\mathbf{c} | s_i)$  and  $P(\mathbf{c}' | s_i)$  normalises these probabilities so that they have the same expected value; thus making a context-length independent comparison possible. When this normalisation is applied to Equation (3), we get the smoothing effect.

## 4 Experimental Sense Embedding Models

We trained two sense embedding models, AdaGram [6] and our own Greedy Sense Embedding method. During training we use the Wikipedia dataset as used by Huang et al. [4]. However, we do not perform the extensive preprocessing used in that work.

Most of our evaluations are carried out on Adaptive SkipGrams (AdaGram) [6]. AdaGram is a non-parametric Bayesian extension of Skip-gram. It learns a number of different word senses, as are required to properly model the language.

We use the implementation<sup>2</sup> provided by the authors with minor adjustments for Julia [21] v0.5 compatibility.

The AdaGram model was configured to have up to 30 senses per word, where each sense is represented by a 100 dimension vector. The sense threshold was set to  $10^{-10}$  to encourage many senses. Only words with at least 20 occurrences are kept, this gives a total vocabulary size of 497,537 words.

To confirm that our techniques are not merely a quirk of the AdaGram method or its implementation, we implemented a new simple baseline word sense embedding method. This method starts with a fixed number of randomly initialised embeddings, then greedily assigns each training case to the sense which predicts it with the highest probability (using Equation (3)). The task remains the same: using skip-grams with hierarchical softmax to predict the context words for the input

<sup>2</sup> <https://github.com/sbos/AdaGram.jl>

word sense. This is similar to [22], however it is using collocation probability, rather than distance in vector-space as the sense assignment measure. Our implementation is based on a heavily modified version of Word2Vec.jl<sup>3</sup>.

This method is intrinsically worse than AdaGram. Nothing in the model encourages diversification and specialisation of the embeddings. Manual inspection reveals that a variety of senses are captured, though with significant repetition of common senses, and with rare senses being missed. Regardless of its low quality, it is a fully independent method from AdaGram, and so is suitable for our use in checking the generalisation of the refitting techniques.

The vocabulary used is smaller than for the AdaGram model. Words with at least 20,000 occurrences are allocated 20 senses. Words with at least 250 occurrences are restricted to a single sense. The remaining rare words are discarded. This results in a vocabulary size of 88,262, with 2,796 words having multiple senses. We always use a uniform prior, as the model does not facilitate easy calculation of the prior.

## 5 Similarity of Words in Context

Estimating word similarity with context is the task of determining how similar words are, when presented with the context they occur in. The goal of this task is to match human judgements of word similarity. For each of the target words and contexts; we use refitting on the target word to create a word sense embedding specialised for the meaning in the context provided. Then the similarity of the refitted vectors can be measured using cosine distance (or similar). By measuring similarity this way, we are defining a new similarity measure.

Reisinger and Mooney [3] define a number of measures for word similarity suitable for use with sense embeddings. The most successful was AvgSimC, which has become the gold standard method for use on similarity tasks. It has been used with great success in many works [4, 11, 5].

AvgSimC is defined using distance metric  $d$  (normally cosine distance) as:

$$\text{AvgSimC}((\mathbf{u}, \mathbf{c}), (\mathbf{u}', \mathbf{c}')) = \frac{1}{n \times n'} \sum_{u_i \in \mathbf{u}} \sum_{u'_j \in \mathbf{u}'} P(u_i | \mathbf{c}) P(u'_j | \mathbf{c}') d(u_i, u'_j) \quad (5)$$

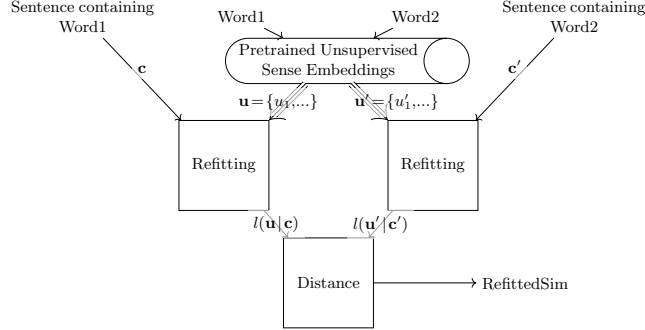
for contexts  $\mathbf{c}$  and  $\mathbf{c}'$ , the contexts of the two words to be compared, and for  $\mathbf{u} = \{u_1, \dots, u_n\}$  and  $\mathbf{u}' = \{u'_1, \dots, u'_{n'}\}$  the respective sets of induced senses of the two words.

### 5.1 A New Similarity Measure: RefittedSim

We define a new similarity measure, RefittedSim, as the distance between the refitted sense embeddings. As shown in Figure 1 the example contexts are used to refit the induced sense embeddings of each word. This is a direct application of Equation (1).

<sup>3</sup> <https://github.com/tanmaykm/Word2Vec.jl/>

Fig. 1: Block diagram for RefittedSim similarity measure



Using the same definitions as in Equation (5), RefittedSim is defined as:

$$\text{RefittedSim}((\mathbf{u}, \mathbf{c}), (\mathbf{u}', \mathbf{c}')) = d(l(\mathbf{u}|\mathbf{c}), l(\mathbf{u}'|\mathbf{c}')) = d\left(\sum_{u_i \in \mathbf{u}} u_i P(u_i|\mathbf{c}), \sum_{u'_j \in \mathbf{u}'} u'_j P(u'_j|\mathbf{c}')\right) \quad (6)$$

AvgSimC is a probability weighted average of pairwise computed distances for each sense vector. Whereas RefittedSim is a single distance measured between the two refitted vectors – which are the probability weighted averages of the original unsupervised sense vectors.

There is a notable difference in time complexity between AvgSimC and RefittedSim. AvgSimC has time complexity  $O(n\|\mathbf{c}\| + n'\|\mathbf{c}'\| + n \times n')$ , while RefittedSim has  $O(n\|\mathbf{c}\| + n'\|\mathbf{c}'\|)$ . The product of the number of senses of each word  $n \times n'$ , may be small for dictionary senses, but it is often large for induced senses. Dictionaries tend to define only a few senses per word – the average<sup>4</sup> number of senses per word in WordNet is less than three [7]. For induced senses, however, it is often desirable to train many more senses, to get better results using the more fine-grained information. Reisinger and Mooney [3] found optimal results in several evaluations near 50 senses. In such cases the  $O(n \times n')$  is significant, avoiding it with RefittedSim makes the similarity measure more useful for information retrieval.

## 5.2 Experimental Setup

We evaluate our refitting method using Stanford’s Contextual Word Similarities (SCWS) dataset [4]. During evaluation, each context paragraph is limited to 5 words to either side of the target word, as in the training.

## 5.3 Results

Table 1a shows the results of our evaluations on the SCWS similarity task. A significant improvement can be seen by applying our techniques.

<sup>4</sup> It should be noted, though, that the number of meanings is not normally distributed [23].



Table 1: Spearman rank correlation  $\rho \times 100$  when evaluated on the SCWS task.

(a) For varying hyper-parameters.					(b) Compared to other methods .			
Method	Geometric Smoothing	Use Prior	AvgSimC	RefittedSim	RefittedSim-S is with smoothing, and RefittedSim-SU is with uniform prior			
AdaGram	T	T	<b>53.8</b>	64.8	Paper	Embedding	Similarity	$\rho \times 100$
AdaGram	T	F	36.1	<b>65.0</b>	This paper	AdaGram	AvgSimC	43.8
AdaGram	F	T	43.8	47.8	This paper	AdaGram	RefittedSim-S	64.8
AdaGram	F	F	20.7	24.1	This paper	AdaGram	RefittedSim-SU	65.0
Greedy	T	F	23.6	49.7	[4]	Huang et al.	AvgSimC	65.7
Greedy	F	F	22.2	40.7	[4]	Pruned tf-idf	AvgSimC	60.5
					[11]	Chen et al.	AvgSimC	68.9
					[5]	Tian et al.	AvgSimC	65.4
					[5]	Tian et al.	MaxSim	<b>65.6</b>
					[9]	SenseEmbed	Min Tanimoto	58.9
					[9]	SenseEmbed	Weighted Tanimoto	62.4

The RefittedSim method consistently outperforms AvgSimC across all configurations. Similarly geometric smoothing consistently improves performance both for AvgSimC and for RefittedSim. The improvement is significantly more for RefittedSim than for AvgSimC results. In general using the unsupervised sense prior estimate from the AdaGram model, improves performance – particularly for AvgSimC. The exception to this is with RefittedSim with smoothing, where it makes very little difference. Unsurprisingly, given its low quality, the Greedy embeddings are always outperformed by AdaGram. It is not clear if these improvements will transfer to clustering based methods due to the differences in how the sense probability is estimated, compared to the language model based method evaluated on in Table 1a.

Table 1b compares our results with those reported in the literature using other methods. These results are not directly comparable, as each method uses a different training corpus, with different preprocessing steps, which can have significant effects on performance. It can be seen that by applying our techniques we bring the results of our AdaGram model from very poor ( $\rho \times 100 = 43.8$ ) when using normal AvgSimC without smoothing, up to being competitive with other models, when using RefittedSim with smoothing. The method of Chen et al. [11], has a significant lead on the other results presented. This can be attributed to its very effective semi-supervised fine-tuning method. This suggests a possible avenue for future development in using refitted sense vectors to relabel a corpus, and then performing fine-tuning similar to that done by Chen et al.

## 6 Word Sense Disambiguation

### 6.1 Refitting for Word Sense Disambiguation

Once refitting has been used to create sense vectors for lexical word senses, an obvious use of them is to perform word sense disambiguation. In this section we refer to the lexical word sense disambiguation problem, i.e. to take a word and find its dictionary sense; whereas the methods discussed in Equations (3)

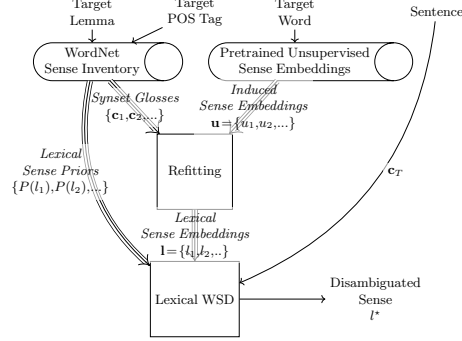


Fig. 2: Block diagram for performing WSD using refitting.

and (4) consider the more general problem, as applicable to disambiguating lexical or induced word senses depending on the inputs. Our overall process shown in Figure 2 uses both: first disambiguating the induced senses as part of refitting, then using the refitted sense vectors to find the most likely dictionary sense.

First, refitting is used to transform the induced sense vectors into lexical sense vectors. We use the targeted word’s lemma (i.e. base form), and part of speech (POS) tag to retrieve all possible definitions of the word (Glosses) from WordNet; there is one gloss per sense. These glosses are used as the example sentence to perform refitting (see Section 3). We find embeddings,  $\mathbf{l} = \{l_1, \dots, l_{n_l}\}$  for each of the lexical word senses using Equation (1). These lexical word senses are still supported by the language model, which means one can apply the general WSD method to determine the posterior probability of a word sense, given an observed context.

When given a sentence  $\mathbf{c}_T$ , containing a target word to be disambiguated, the probability of each lexical word sense  $P(l_i | \mathbf{c}_T)$ , can be found using Equation (3) (or the smoothed version Equation (4)), over the lexically refitted sense embeddings. Then, selecting the correct sense is simply selecting the most likely sense:

$$l^*(\mathbf{l}, \mathbf{c}_T) = \underset{\forall l_i \in \mathbf{l}}{\operatorname{argmax}}: P(l_i | \mathbf{c}_T) = \underset{\forall l_i \in \mathbf{l}}{\operatorname{argmax}}: \frac{P(\mathbf{c}_T | l_i) P(l_i)}{\sum_{\forall l_j \in \mathbf{l}} P(\mathbf{c}_T | l_j) P(l_j)} \quad (7)$$

## 6.2 Lexical Sense Prior

WordNet includes frequency counts for each word sense based on Semcor [24]. These form a prior for  $P(l_i)$ . The comparatively small size of Semcor means that many word senses do not occur at all. We apply add-one smoothing to remove any zero counts. This is in addition to using our proposed geometric smoothing as an optional part of the general WSD. Geometric smoothing serves a different (but related) purpose, of decreasing the sharpness of the likelihood function – not of removing impossibilities from the prior.

### 6.3 Experimental Setup

The WSD performance is evaluated on the SemEval 2007 Task 7.

We use the weighted mapping method of Agirre et al. [12], (see Section 2.2) as a baseline alternative method for using WSI senses for WSD. We use Semcor as the mapping corpus, to derive the mapping weights.

The second baseline we use is the Most Frequent Sense (MFS). This method always disambiguates any word as having its most common meaning. Due to the power law distribution of word senses, this is a very effective heuristic [20]. We also consider the results when using a backoff to MSF when a method is unable to determine the word sense the method can report the MFS instead of returning no result (a non-attempt).

### 6.4 Word Sense Disambiguation Results

Method	Attempted	Precision	Recall	F1
Refitted-S AdaGram	99.91%	0.799	0.799	<b>0.799</b>
Refitted AdaGram	99.91%	0.774	0.773	0.774
Refitted-S Greedy	79.95%	0.797	0.637	0.708
Refitted-S Greedy *	100.00%	0.793	0.793	0.793
Refitted Greedy	79.95%	0.725	0.580	0.645
Refitted Greedy *	100.00%	0.793	0.793	0.793
Mapped AdaGram	84.31%	0.776	0.654	0.710
Mapped AdaGram *	100.00%	0.736	0.736	0.736
MFS baseline	100.00%	0.789	0.789	0.789

Table 2: Results on SemEval 2007 Task 7 – course-all-words disambiguation. The -S marks results using geometric smoothing. The \* marks results with MSF backoff.

The results of employing our method for WSD, are shown in Table 2. Our results using smoothed refitting, both with AdaGram and Greedy Embeddings with backoff, outperform the MSF baseline [25] – noted as a surprisingly hard baseline to beat [11].

The mapping method [12] was not up to the task of mapping unsupervised senses to supervised senses, on this large scale task. The Refitting method works better. Though refitting is only usable for language-model embedding WSI, the mapping method is suitable for all WSI systems.

While not directly comparable due to the difference in training data, we note that our Refitted results, are similar in performance, as measured by F1 score, to the results reported by Chen et al. [11]. AdaGram with smoothing, and Greedy embeddings with backoff have close to the same result as reported for L2R with backoff – with the AdaGram slightly better and the Greedy embeddings slightly worse. They are exceeded by the best method reported in that paper: S2C method with backoff. Comparison to non-embedding based methods is not discussed here for brevity. Historically state of the art systems have functioned very differently; normally by approaching the WSD task by more direct means.

Our results are not strong enough for Refitted AdaGram to be used as a WSD method on its own, but do demonstrate that the senses found by refitting are capturing the information from lexical senses. It is now evident that the refitted

sense embeddings are able to perform WSD, which was not possible with the unsupervised senses.

## 7 Conclusion

A new method is proposed for taking unsupervised word embeddings, and adapting them to align to particular given lexical senses, or user provided usage examples. This refitting method thus allows us to find word sense embeddings with known meaning. This method can be seen as a one-shot learning task, where only a single labelled example of each class is available for training. We show how our method can be used to create embeddings to evaluate the similarity of words, given their contexts.

This allows us to propose a new similarity measuring method, RefittedSim. The performance of RefittedSim on AdaGram is comparable to the results reported by the researchers of other sense embeddings techniques using AvgSimC, but its time complexity is significantly lower. We also demonstrate how similar refitting principles can be used to create a set of vectors that are aligned to the meanings in a sense inventory, such as WordNet.

We show how this can be used for word sense disambiguation. On this difficult task, it performs marginally better than the hard to beat MFS baseline, and significantly better than a general mapping method used for working with WSI senses on lexical WSD tasks. As part of our method for refitting, we present a geometric smoothing to overcome the issues of overly dominant senses probability estimates. We show that this significantly improves the performance. Our refitting method provides effective bridging between the vector space representation of meaning, and the traditional discrete lexical representation. More generally it allows a sense embedding to be created to model the meaning of a word in any given sentence. Significant applications of sense embeddings in tasks such as more accurate information retrieval thus become possible.

## References

1. Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv:1301.3781 (2013)
2. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP 2014). (2014) 1532–1543
3. Reisinger, J., Mooney, R.J.: Multi-prototype vector-space models of word meaning. In: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Association for Computational Linguistics (2010) 109–117
4. Huang, E.H., Socher, R., Manning, C.D., Ng, A.Y.: Improving word representations via global context and multiple word prototypes. In: Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1, Association for Computational Linguistics (2012) 873–882
5. Tian, F., Dai, H., Bian, J., Gao, B., Zhang, R., Chen, E., Liu, T.Y.: A probabilistic model for learning multi-prototype word embeddings. In: COLING. (2014) 151–160

6. Bartunov, S., Kondrashkin, D., Osokin, A., Vetrov, D.P.: Breaking sticks and ambiguities with adaptive skip-gram. *CoRR* **abs/1502.07257** (2015)
7. Miller, G.A.: Wordnet: a lexical database for english. *Communications of the ACM* **38** (1995) 39–41
8. Véronis, J.: A study of polysemy judgements and inter-annotator agreement. In: Programme and advanced papers of the Senseval workshop. (1998) 2–4
9. Iacobacci, I., Pilehvar, M.T., Navigli, R.: Senseembed: learning sense embeddings for word and relational similarity. In: *Proceedings of ACL*. (2015) 95–105
10. Moro, A., Raganato, A., Navigli, R.: Entity Linking meets Word Sense Disambiguation: a Unified Approach. *Transactions of the Association for Computational Linguistics (TACL)* **2** (2014) 231–244
11. Chen, X., Liu, Z., Sun, M.: A unified model for word sense representation and disambiguation. In: *EMNLP, Citeseer* (2014) 1025–1035
12. Agirre, E., Martínez, D., De Lacalle, O.L., Soroa, A.: Evaluating and optimizing the parameters of an unsupervised graph-based wsd algorithm. In: *Proceedings of the first workshop on graph based methods for natural language processing*, Association for Computational Linguistics (2006) 89–96
13. Agirre, E., Soroa, A.: Semeval-2007 task 02: Evaluating word sense induction and discrimination systems. In: *Proceedings of the 4th International Workshop on Semantic Evaluations. SemEval '07*, Stroudsburg, PA, USA, Association for Computational Linguistics (2007) 7–12
14. Nocedal, J.: Updating quasi-newton matrices with limited storage. *Mathematics of computation* **35** (1980) 773–782
15. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in Neural Information Processing Systems*. (2013) 3111–3119
16. Mikolov, T., Yih, W.t., Zweig, G.: Linguistic regularities in continuous space word representations. In: *HLT-NAACL*. (2013) 746–751
17. Bengio, Y., Ducharme, R., Vincent, P., Janvin, C.: A neural probabilistic language model. *The Journal of Machine Learning Research* (2003) 137–186
18. Rosenfeld, R.: Two decades of statistical language modeling: Where do we go from here? *Proceedings of the IEEE* **88** (2000) 1270–1278
19. Zipf, G.: Human behavior and the principle of least effort: an introduction to human ecology. Addison-Wesley Press (1949)
20. Kilgariff, A. In: *How Dominant Is the Commonest Sense of a Word?* Springer Berlin Heidelberg, Berlin, Heidelberg (2004) 103–111
21. Bezanson, J., Edelman, A., Karpinski, S., Shah, V.B.: Julia: A fresh approach to numerical computing. (2014)
22. Neelakantan, A., Shankar, J., Passos, A., McCallum, A.: Efficient non-parametric estimation of multiple embeddings per word in vector space. *arXiv preprint arXiv:1504.06654* (2015)
23. Zipf, G.K.: The meaning-frequency relationship of words. *The Journal of general psychology* **33** (1945) 251–256
24. Tengi, R.L.: Design and implementation of the WordNet lexical database and searching software. In: *WordNet: an electronic lexical database*, The MIT Press, Cambridge, Massachusetts. (1998) 105
25. Navigli, R., Litkowski, K.C., Hargraves, O.: Semeval-2007 task 07: Coarse-grained english all-words task. In: *Proceedings of the 4th International Workshop on Semantic Evaluations. SemEval '07*, Stroudsburg, PA, USA, Association for Computational Linguistics (2007) 30–35