A SIMPLE AND EFFICIENT METHOD TO GENERATE WORD SENSE REPRESENTATIONS

Luis Nieto Piña

Språkbanken, Department of Swedish University of Gothenburg luis.nieto.pina@svenska.gu.se

Richard Johansson

Språkbanken, Department of Swedish University of Gothenburg richard.johansson@svenska.gu.se

ABSTRACT

Distributed representations of words have boosted the performance of many Natural Language Processing tasks. However, usually only one representation per word is obtained, not acknowledging the fact that some words have multiple meanings. This has a negative effect on the individual word representations and the language model as a whole. In this paper we present a simple model that enables recent techniques for building word vectors to represent distinct senses of polysemic words. In our assessment of this model we show that it is able to effectively discriminate between words' senses and to do so in a computationally efficient manner.

1 Introduction

Distributed representations of words have helped obtain better language models (Bengio et al., 2003) and improve the performance of many natural language processing applications such as named entity recognition, chunking, paraphrasing, or sentiment classification (Turian et al., 2010; Socher et al., 2011; Glorot et al., 2011). Recently, Mikolov et al. (2013a;b) proposed the Skip-gram model, which is able to produce high-quality representations from large collections of text in an efficient manner.

Despite the achievements of distributed representations, polysemy or homonymy are usually disregarded even when word semantics may have a large influence on the models. This results in several distinct senses of one same word sharing a representation, and possibly influencing the representations of words related to those distinct senses under the premise that similar words should have similar representations.

There have been some recent attempts to address this issue. Reisinger & Mooney (2010) propose clustering feature vectors corresponding to occurrences of words into a predetermined number of clusters, whose centers provide multiple representations for each word, in an attempt to capture distinct usages of that word. Huang et al. (2012) cluster context vectors instead, built as weighted average of the words' vectors that surround a certain target word. Neelakantan et al. (2014) imports the idea of clustering context vectors into the Skip-gram model.

We present a simple method for obtaining sense representations directly during the Skip-gram training phase. It differs from Neelakantan et al. (2014)'s approach in that it does not need to create or maintain clusters to discriminate between senses, leading to a significant reduction in the model's complexity. In the following sections we describe the model and describe an initial assessment of the results it can produce.

2 Model Description

2.1 From Word Forms to Senses

The distributed representations for word forms that stem from a Skip-gram (Mikolov et al., 2013a;b) model are built on the premise that, given a certain target word, they should serve to predict its surrounding words in a text. I.e., the training of a Skip-gram model, given a target word w, is based on maximizing the log-probability of the context words of w, c_1, \ldots, c_n :

$$\sum_{i=1}^{n} \log p(c_i|w) \tag{1}$$

The training data usually consists of a large collection of sentences or documents, so that the role of target word w can be iterated over these sequences of words, while the context words c considered in each case are those that surround w within a window of a certain length, . The objective becomes then maximizing the average sum of the log-probabilities from Eq. 1.

We propose modify this model to include a sense s of the word w. Note that Eq. 1 equals

$$\log p(c_1, \dots, c_n | w) \tag{2}$$

if we assume the context words c_i to be independent of each other given a target word w. The notation in Eq. 2 allows us to consider the Skip-gram as a Naïve Bayes model parameterized by word embeddings (Mnih & Kavukcuoglu, 2013). In this scenario, including a sense would amount then to adding a latent variable s, and our model's behaviour given a target word w is to select a sense s, which is in its turn used to predict n context words c_1, \ldots, c_n . Formally:

$$p(s, c_1, \dots, c_n | w) = p(s | w) \cdot p(c_1, \dots, c_n | s) = p(s | w) \cdot p(c_1 | s) \dots p(c_n | s),$$
(3)

Thus, our training objective is to maximize the sum of the log-probabilities of context words c given a sense s of the target word w plus the log-probability of the sense s given the target word:

$$\log p(s|w) + \sum_{i=1}^{n} \log p(c_i|s) \tag{4}$$

We must now consider two distinct vocabularies: V containing all possible word forms (context and target words), and S containing all possible senses for the words in V, with sizes |V| and |S|, resp. Given a pre-set $D \in \mathbb{N}$, our ultimate goal is to obtain |S| dense, real-valued vectors of dimension D that represent the senses in our vocabulary S according to the objective function defined in Eq. 4.

The neural architecture of the Skip-gram model works with two separate representations for the same vocabulary of words. This double representation is not motivated in the original papers, but it stems from word2vec's code 1 that the model builds separate representations for context and target words, of which the former constitute the actual output of the system. (A note by Goldberg & Levy (2014) offers some insight into this subject.) We take advantage of this architecture and use one of these two representations to contain senses, rather than word forms: as our model only uses target words w as an intermediate step to select a sense s, we only do not need to keep a representation for them. In this way, our model builds a representation of the vocabulary V, for the context words, and another for the vocabulary S of senses, which contains the actual output. Note that the representation of context words is only used internally for the purposes of this work, and that context words are word forms; i.e., we only consider senses for the target words.

2.2 Selecting a Sense

In the description of our model above we have considered that for each target word w we are able to select a sense s. We now explain the mechanism used for this purpose. The probability of a context

¹http://code.google.com/p/word2vec/

word c_i given a sense s, as they appear in the model's objective function defined in Eq. 4, $p(c_i|s)$, $\forall i \in [1, n]$, can be calculated using the *softmax* function:

$$p(c_i|s) = \frac{e^{v_{c_i}^{\mathsf{T}} \cdot v_s}}{\sum_{i=1}^{|V|} e^{v_{c_j}^{\mathsf{T}} \cdot v_s}} = \frac{e^{v_{c_i}^{\mathsf{T}} \cdot v_s}}{Z(s)},\tag{5}$$

where v_{c_i} (resp. v_s) denotes the vector representing context word c_i (resp. sense s), v^\intercal denotes the transposed vector v, and in the last equality we have used Z(s) to identify the normalizer over all context words. With respect to the probability of a sense s given a target word w, for simplicity we assume that all senses are equally probable; i.e., $p(s|w) = \frac{1}{K}$ for any of the K senses s of word w, senses s0.

Using Bayes formula on Eq. 3, we can now obtain the posterior probability of a sense s given the target word w and the context words c_1, \ldots, c_n :

$$p(s|c_1, ..., c_n, w) = \frac{p(s|w) \cdot p(c_1, ..., c_n|s)}{\sum_{s_k \in \text{senses}(w)} p(s_k|w) \cdot p(c_1, ..., c_n|s_k)}$$

$$= \frac{e^{(c_1 + \dots + c_n) \cdot s} \cdot Z(s)^{-n}}{\sum_{s_k \in \text{senses}(w)} e^{(c_1 + \dots + c_n) \cdot s_k} \cdot Z(s_k)^{-n}}$$
(6)

During training, thus, given a target word w and context words $c_1, \ldots c_n$, the most probable sense $s \in \text{senses}(w)$ is the one that maximizes Eq. 6. Unfortunately, in most cases it is computationally impractical to explicitly calculate Z(s). From a number of possible approximations, we have empirically found that considering Z(s) to be constant yields the best results; this is not an unreasonable approximation if we expect the context word vectors to be densely and evenly spread out in the vector space. Under this assumption, the most probable sense s of w is the one that maximizes

$$\frac{e^{(c_1+\dots+c_n)\cdot s}}{\sum_{s_k \in \text{senses}(w)} e^{(c_1+\dots+c_n)\cdot s_k}}$$
(7)

For each word occurrence, we propose to select and train only its most probable sense. This approach of *hard sense assignments* is also taken in Neelakantan et al. (2014)'s work and we follow it here, although it would be interesting to compare it with a *soft* updates of all senses of a given word weighted by the probabilities obtained with Eq. 6.

The training algorithm, thus, iterates over a sequence of words, selecting each one in turn as a target word w and its context words as those in a window of a maximum pre-set size. For each target word, a number K of senses s is considered, and the most probable one selected according to Eq. 7. (Note that, as the number of senses needs to be informed (using, for example, a lexicon), monosemic words need only have one representation.) The selected sense s substitutes the target word w in the original Skip-gram model, and any of the known techniques used to train it can be subsequently applied to obtain sense representations. The training process is drafted in Algorithm 1 using Skip-gram with Negative Sampling.

Negative Sampling (Mikolov et al., 2013b), based on Noise Contrastive Estimation (Mnih & Teh, 2012), is a computationally efficient approximation for the original Skip-gram objective function (Eq. 1). In our implementation it learns the sense representations by sampling N_{neg} words from a

noise distribution and using logistic regression to distinguish them from a certain context word c of a target word w. This process is also illustrated in Algorithm 1.

Algorithm 1: Selection of senses and training using Skip-gram with Negative Sampling. (Note that v_x denotes the vector representation of word x.)

```
Input: Sequence of words w_1, \ldots, w_N, window size n, learning rate \alpha, number of negative
              words N_{neq}
   Output: Updated vectors for each sense of words w_i, i = 1, ..., N
 1 for t = 1, ..., N do
 2
         w = w_i
         K \leftarrow number of senses of w
 3
         context(w) = \{c_1, \ldots, c_n \mid c_i = w_{t+i}, i = -n, \ldots, n, i \neq 0\}
 4
        for k = 1, ..., K do p_k = \frac{e^{(v_{c_1} + \dots + v_{c_n}) \cdot v_{s_k}}}{\sum_{j=1}^K e^{(v_{c_1} + \dots + v_{c_n}) \cdot v_{s_k}}}
 5
 6
         s = \arg \max_{k=1,\dots,K} p_k
7
        for i=1,\ldots,n do
 8
              f = \frac{1}{1 + e^{v_{c_i} \cdot v_s}}
 9
              g_c = \alpha(1 - f)
10
              Update v_{c_i} with g_d
11
12
              for d=1,\ldots,N_{neg} do
13
                   c \leftarrow \text{word sampled from noise distribution}, c \neq c_i
14
                   f = \frac{1}{1 + e^{v_c \cdot v_s}}
15
                   g_d = -\alpha \cdot f
16
                   Update v_c with g_d
17
                   g = g + g_d
18
              Update v_s with g
19
```

3 EXPERIMENTS

We trained the model described in Section 2 on Swedish text using a context window of 10 words and vectors of 200 dimensions. The model requires the number of senses to be specified for each word; as a heuristic, we used the number of senses listed in the SALDO lexicon (Borin et al., 2013).

As a training corpus, we created a corpus of 1 billion words downloaded from Språkbanken, the Swedish language bank.² The corpora are distributed in a format where the text has been tokenized, part-of-speech-tagged and lemmatized. Compounds have been segmented automatically and when a lemma was not listed in SALDO, we used the parts of the compounds instead. The input to the software computing the embeddings consisted of lemma forms with concatenated part-of-speech tags, e.g. *dricka*-verb for the verb 'to drink' and *dricka*-noun for the noun 'drink'.

The training time of our model on this corpus was 22 hours. For the sake of time performance comparison, we run an off-the-shelf word2vec execution on our corpus using the same parameterization described above; the training of word vectors took 20 hours, which illustrates the little complexity that our model adds to the original Skip-gram.

3.1 Inspection of nearest neighbors

We evaluate the output of the algorithm qualitatively by inspecting the nearest neighbors of the senses of a number of example words, and comparing them to the senses listed in SALDO. We leave a quantitative evaluation to future work.

Table 1 shows the nearest neighbor lists of the senses of two words where the algorithm has been able to learn the distinctions used in the lexicon. The verb flyga 'to fly' has two senses listed in SALDO: to travel by airplane and to move through the air. The adjective $\ddot{o}m$ 'tender' also has

²http://spraakbanken.gu.se

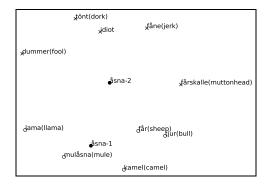


Figure 1: 2D projections of the two senses of *åsna* ('donkey' and 'slow-witted person') and their nearest neighbors.

two senses, similar to the corresponding English word: one emotional and one physical. The lists are semantically coherent, although we note that they are topical rather than substitutional; this is expected since the algorithm was applied to lemmatized and compound-segmented text and we use a fairly wide context window.

```
flyg 'flight'
                               flaxa 'to flap wings'
                                                          kärleksfull 'loving'
                                                                                        svullen 'swollen'
flygning 'flight'
                               studsa 'to bounce'
                                                          ömsint 'tender'
                                                                                        ömma 'to be sore'
flygplan 'airplane'
                               sväva 'to hover'
                                                          smek 'caress'
                                                                                        värka 'to ache'
charterplan 'charter plane'
                               skjuta 'to shoot'
                                                          kärleksord 'word of love'
                                                                                        mörbulta 'to bruise'
SAS-plan 'SAS plane'
                               susa 'to whiz'
                                                          ömtålig 'delicate'
                                                                                        ont 'pain'
                  (a) flyga 'to fly'
                                                                            (b) öm 'tender'
```

Table 1: Examples of nearest neighbors of the two senses of two example words.

In a related example, Figure 1 shows the projections onto a 2D space ³ of the representations for the two senses of *åsna*: 'donkey' or 'slow-witted person', and those of their corresponding nearest neighbors.

For some other words we have inspected, we fail to find one or more of the senses. This is typically when one sense is very dominant, drowning out the rare senses. For instance, the word *rock* has two senses, 'rock music' and 'coat', where the first one is much more frequent. While one of the induced senses is close to some pieces of clothing, most of its nearest neighbors are styles of music.

In other cases, the algorithm has come up with meaningful sense distinctions, but not exactly as in the lexicon. For instance, the lexicon lists two senses for the noun *böna*: 'bean' and 'girl'; the algorithm has instead created two bean senses: bean as a plant part or bean as food. In some other cases, the algorithm finds genre-related distinctions instead of sense distinctions. For instance, for the verb *älska*, with two senses 'to love' or 'to make love', the algorithm has found two stylistically different uses of the first sense: one standard, and one related to informal words frequently used in social media. Similarly, for the noun *svamp* 'sponge' or 'mushroom'/fungus', the algorithm does not find the sponge sense but distinguishes taxonomic, cooking-related, and nature-related uses of the mushroom/fungus sense. It's also worth mentioning that when some frequent foreign word is homographic with a Swedish word, it tends to be assigned to a sense. For instance, for the adjective *sur* 'sour', the lexicon lists one taste and one chemical sense; the algorithm conflates those two senses but creates a sense for the French preposition.

4 CONCLUSIONS AND FUTURE WORK

In this paper, we present a model for automatically building sense vectors based on the Skip-gram method. In order to learn the sense vectors, we modify the Skip-gram model to take into account the

³The projection was computed using scikit-learn (Pedregosa et al., 2011) using multidimensional scaling of the distances in a 200-dimensional vector space.

number of senses of each target word. By including a mechanism to select the most probable sense given a target word and its context, only slight modifications to the original training algorithm are necessary for it to learn distinct representations of word senses from unstructured text.

To evaluate our model we train it on a 1-billion-word Swedish corpus and use the SALDO lexicon to inform the number of senses associated to each word. Over a series of examples in which we analyse the nearest neighbors of some of the represented senses, we show how the obtained sense representations are able to replicate the senses defined in SALDO, or to make novel sense distinctions in others. On instances in which a sense is dominant we observe that the obtained representations favour this sense in detriment of other, less common ones.

We have used a lexicon just for setting the number of senses of a given word, and showed that with that information we are able to obtain coherent sense representations. An interesting line of research lies in further exploiting existing knowledge resources for learning better sense vectors. E.g., integrating in this model the network topology of a lexicon such as SALDO, that links together senses of related words, could arguably help improve the representations for those rare senses with which our model currently struggles by learning their representations taking into account those of neighboring senses.

We also hope to provide a more systematic evaluation of our model so that a more accurate assessment of its qualities can be made, and its performance more easily compared against that of similar recent work.

REFERENCES

- Bengio, Yoshua, Ducharme, Réjean, Vincent, Pascal, and Jauvin, Christian. A neural probabilistic language model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- Borin, Lars, Forsberg, Markus, and Lönngren, Lennart. SALDO: a touch of yin to WordNet's yang. *Language Resources and Evaluation*, 47:1191–1211, 2013.
- Glorot, Xavier, Bordes, Antoine, and Bengio, Yoshua. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 513–520, 2011.
- Goldberg, Yoav and Levy, Omer. word2vec explained: deriving Mikolov et al.'s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*, 2014.
- Huang, Eric H, Socher, Richard, Manning, Christopher D, and Ng, Andrew 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*, pp. 873–882. Association for Computational Linguistics, 2012.
- Mikolov, Tomas, Chen, Kai, Corrado, Greg, and Dean, Jeffrey. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013a.
- Mikolov, Tomas, Sutskever, Ilya, Chen, Kai, Corrado, Greg S, and Dean, Jeff. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pp. 3111–3119, 2013b.
- Mnih, Andriy and Kavukcuoglu, Koray. Learning word embeddings efficiently with noise-contrastive estimation. In Advances in Neural Information Processing Systems, pp. 2265–2273, 2013.
- Mnih, Andriy and Teh, Yee Whye. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*, 2012.
- Neelakantan, Arvind, Shankar, Jeevan, Passos, Alexandre, and McCallum, Andrew. Efficient non-parametric estimation of multiple embeddings per word in vector space. In *Proceedings of EMNLP*, 2014.

- Pedregosa, Fabian, Varoquaux, Gaël, Gramfort, Alexandre, Michel, Vincent, Thirion, Bertrand, Grisel, Olivier, Blondel, Mathieu, Prettenhofer, Peter, Weiss, Ron, Dubourg, Vincent, VanderPlas, Jake, Passos, Alexandre, Cournapeau, David, Brucher, Matthieu, Perrot, Matthieu, and Duchesnay, Edouard. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Reisinger, Joseph and Mooney, Raymond 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*, pp. 109–117. Association for Computational Linguistics, 2010.
- Socher, Richard, Huang, Eric H, Pennin, Jeffrey, Manning, Christopher D, and Ng, Andrew Y. Dynamic pooling and unfolding recursive autoencoders for paraphrase detection. In *Advances in Neural Information Processing Systems*, pp. 801–809, 2011.
- Turian, Joseph, Ratinov, Lev, and Bengio, Yoshua. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pp. 384–394. Association for Computational Linguistics, 2010.