

# Modal sense classification with task-specific context embeddings

Bo Li<sup>1,2</sup>, Mathieu Dehouck<sup>2</sup> and Pascal Denis<sup>2</sup> \*

1- Univ. Lille, CNRS, UMR 8163 - STL - Savoirs Textes Langage  
59000 Lille, France

2- Magnet, INRIA Lille - Nord Europe, 59650 Villeneuve d'Ascq, France  
{bo.li, pascal.denis, mathieu.dehouck}@inria.fr

**Abstract.** Sense disambiguation of modal constructions is a crucial part of natural language understanding. Framed as a supervised learning task, this problem heavily depends on an adequate feature representation of the modal verb context. Inspired by recent work on general word sense disambiguation, we propose a simple approach of modal sense classification in which standard shallow features are enhanced with task-specific context embedding features. Comprehensive experiments show that these enriched contextual representations fed into a simple SVM model lead to significant classification gains over shallow feature sets.

## 1 Introduction

Modal sense analysis plays an important role in NLP tasks such as sentiment analysis [1], hedge detection [2] and factuality recognition [3]. As an important part of modality analysis, sense classification of modal verbs has attracted notable interest of researchers in recent years. Modal sense classification (MSC) can be treated as a special type of Word Sense Disambiguation (WSD) [4], but differs from general WSD tasks in that MSC deals with a small set of modal verbs with a limited set of senses. For instance, the English modal verb *can* has a limited sense set containing *dynamic*, *deontic*, and *epistemic* senses.

MSC is typically modeled in a supervised learning framework, which in turn raises the question of how to build the most effective feature representation of the modal context. Existing studies have proposed various context representations. A first set of approaches rely on manually engineered features derived from sophisticated linguistic pre-processing. For instance, [5] uses syntax-based features in addition to shallow features such as POS and lemma. Pushing this line of research even further, [6, 7] propose to use semantic-rich features that are related to lexical, proposition-level and discourse-level semantic factors. Taking a drastically different approach, [8] employs Convolution Neural Networks (CNNs) to automatically extract modal sense features, and show that their approach is competitive with hand-crafted feature-based classifiers.

In this paper, we explore an intermediate line of research, attempting to leverage both shallow feature sets *and* pre-trained embeddings for the context words in a very simple learning architecture. That is, we propose to enrich a very

---

\*This work was supported by the ANR-16-CE93-0009 REM project.

simple feature set (basically, POS tags, uni-, bi- and tri-grams), known to be very effective for WSD, with embeddings for the context words. Inspired by the recent work [9] for general WSD, we propose several weighting schemes, including a new task-specific one, for constructing context embeddings from pre-trained word embeddings. Experimental results prove that this approach outperforms state-of-the-art methods over two standard MSC datasets.

## 2 Related work

**Modal sense classification.** Similar to WSD, MSC has been framed as a supervised learning problem. [5] presents an annotation schema for modal verbs in the MPQA corpus and represent modal senses with word-based and syntactic features. On top of feature sets, they perform machine learning with a maximum entropy classifier. Since the corpus used in [5] is relatively small and unbalanced, [6, 7] develop a paraphrase-driven, cross-lingual projection approach to create large and balanced corpora automatically. In addition, they develop semantically grounded features such as lexical features of the modified verb and subject-related features for sense representation. [8] casts modal sense classification as a novel semantic sentence classification task using CNNs to model propositions.

**Word sense disambiguation.** Word sense disambiguation is a long-standing challenge in natural language processing [4]. We focus here on most recent advances in WSD which come from the use of neural networks. One category of neural WSD works use neural networks to learn word embeddings [10, 9] or context representation [11, 12] that can be further used separately or as additional features to shallow sense features. The other category of approaches [13] model WSD as an end-to-end sequence learning problem. These two categories of approaches show similar performance as reported in [13]. Our work is inspired by [9] which has shown state-of-the-art performance in WSD and which is convenient for adaptation to MSC.

## 3 Approach

Following [9], we will use word embeddings trained on unlabeled data to enrich shallow features used in general WSD tasks. As far as we can tell, such a combination has never been investigated for MSC.

### 3.1 Shallow features (SF)

We refer to standard features used in WSD tasks and make use of three shallow feature sets as general features for MSC. For all feature sets, we compute with the same context window surrounding a target modal verb with size  $2n$  ( $n$  left,  $n$  right). The first feature set is POS tags of context words. The second feature set is context words excluding stop words. The third feature set is the local collocations of a target word which are ordered sequences of words appearing in the context window. We use several uni-gram, bi-gram and tri-gram patterns as

the collocation patterns in the context window. We use these simple yet efficient features to do away with heavy feature engineering.

### 3.2 Word embedding features

Embeddings of context words can be combined in a way such that a comprehensive context embedding vector can be obtained to represent the modal context. The context vector can then be combined with those shallow features. Inspired by [9], we obtain context representation by taking a weighted sum of embedding vectors of all context words. Given a target modal verb  $w_0$  and a context window size  $2n$ , let us denote the context word at the relative position  $i$  ( $i \in [-n, 0) \cup (0, n]$ ) as  $w_i$ . The context embedding of  $w_0$  (denoted as  $\mathbf{w}_c \in R^d$ ) can be written as the weighted sum of embedding vector  $\mathbf{w}_i \in R^d$  of each context word  $w_i$ , which is:

$$\mathbf{w}_c = \sum_{i=-n \& i \neq 0}^n f(w_i) \mathbf{w}_i \quad (1)$$

where  $f(w_i)$  is a weighting function measuring the importance of each context word  $w_i$  w.r.t. the modal verb  $w_0$ . For simplicity, we assume that  $f(w_i)$  only depends on  $i$  (i.e., the relative position of  $w_i$  w.r.t.  $w_0$ ).

**Average weighting (WE<sub>av</sub>).** The most intuitive weighting schema simply treats all the context words equally, which corresponds to:

$$f(w_i) = \frac{1}{2n} \quad (2)$$

**Linear weighting (WE<sub>li</sub>).** In a linear manner, the importance of a word is assumed to be inversely proportional to its distance from the target word, which can be formulated as:

$$f(w_i) = 1 - \frac{|i - 1|}{n} \quad (3)$$

**Exponential weighting (WE<sub>ex</sub>).** Compared to linear manner, the exponential weighting schema put more emphasis on the nearest context words. The importance of a word decreases exponentially w.r.t. its distance from the target word, which is:

$$f(w_i) = \alpha^{\frac{i-1}{n-1}} \quad (4)$$

where  $\alpha$  is a hyper-parameter. The above weighting schema have been investigated in [9].

**Modal-specific weighting.** In order to adapt the weighting function  $f(w_i)$  to MSC, we want to assign more weights to context words which have closer connections to modal senses. Inspired by findings in [6], we note that the embedded verb in the scope of the modal and the subject noun phrase are two components acting as strong hints to determine the modal verb sense. Let us consider as an example the sentence *the springbok can jump very high*. When the embedded verb *jump* appears, we prefer a *dynamic* reading of *can*. Based on the intuition above, we assign more weights to context word which is the

embedded verb or which is the head in the subject noun phrase. Formally, we define a modal-specific weighting function  $g$  such that:

$$g(w_i) = \begin{cases} \beta_1 f(w_i), & w_i \text{ is embedded verb} \\ \beta_2 f(w_i), & w_i \text{ is the head in the subj noun phrase} \\ f(w_i), & \text{else} \end{cases}$$

where  $\beta_1$  and  $\beta_2$  are hyper-parameters indicating that words that are more relevant to modal verb senses take more responsibility, and  $f$  is one of the weighting schema defined above. We can obtain a context representation  $w_c$  that is optimized for MSC task by replacing  $f$  in equation 1 with  $g$  defined here.

## 4 Experiments

In this section, we perform comprehensive experiments in order to evaluate the usefulness of word embedding features in MSC.

**Datasets.** We make use here of two representative corpora for MSC. One is the small and unbalanced corpus annotated manually based on the MPQA corpus in [5]. The other one is the larger and balanced corpus EPOS which is constructed automatically via paraphrase-driven modal sense projection in [6]. We use two training/testing settings: (1) both training and testing sets picked from MPQA. (2) training set from EPOS and testing set from MPQA. Characteristics of training/testing sets are not given here due to space reasons.

**Experimental settings.** We have used Stanford parser to pre-process the original corpora. For classification, we employ SVM implemented in LibSVM [14]. The context window size  $2n$  is picked from  $\{2, 6, 10, 20\}$ . The parameter  $\alpha$  is set to  $\{0.1, 0.2, 0.5, 0.8\}$ , and the parameters  $\beta_1$  and  $\beta_2$  are selected from  $\{1, 2, 5, 10\}$ . The word embeddings are trained with word2vec[15] on Wikipedia. The hyper-parameters are optimized via 5-fold cross validation. We employ accuracy as the performance measure and McNemar’s test ( $p < 0.05$ ) to test significance.

**Competing Systems.** we use max frequent sense (MFS) baseline for unbalanced classifier and random sense (RS) baseline for balanced classifier. The other baseline systems are RR [5] and ZH [6]. In addition, we compare different combinations of features presented in section 3. SF is the shallow features, and  $+WE_{**}$  denotes the concatenation of word embedding features  $WE_{**}$  to SF.  $+WE_{**}'$  stands for the modal specific weighting version.

**Results on unbalanced corpus.** We firstly show the results on the unbalanced MPQA corpus in table 1. The performance of SF is comparable to RR and slightly inferior to ZH, given that SF contains similar features as RR. Furthermore, when we consider word embeddings as extra features, we obtain results that are at least as good as SF. The best performance comes from  $+WE_{ex}'$  where, compared to SF, we have gotten significant improvement on *can* and insignificant improvement on *could*, and slightly decrease that is not significant on *should*. However, with any of  $+WE_{av}$ ,  $+WE_{li}$  and  $+WE_{ex}$ , we cannot achieve any significant improvement. We further note that MFS is a strong baseline on

the unbalanced corpus. Neither RR nor ZH is able to beat the MFS baseline with ZH being slightly better than RR, which is coincident with conclusions in previous work [6].

Table 1: Accuracy with various features on MPQA. + or - indicates that the improvements or degradations with respect to SF are statistically significant. The highest value in each row is marked in bold.

	MFS	RR	ZH	SF	+WE <sub>av</sub>	+WE <sub>li</sub>	+WE <sub>ex</sub>	+WE <sub>ex'</sub>
can	69.9	66.6	66.1	63.5	66.2	66.7	67.1	<b>70.1</b> <sup>+</sup>
could	65.0	62.5	67.9	67.5	67.8	68.9	70.4	<b>70.5</b>
may	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>	<b>93.6</b>
must	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>	<b>94.3</b>
shall	<b>84.6</b>	83.3	83.3	83.3	83.3	83.3	83.3	83.3
should	90.8	90.8	<b>92.9</b>	91.1	90.9	90.8	90.8	90.8

**Results on balanced corpus.** Since the MPQA training set is small in size and unbalanced in sense distribution, machine learning algorithms tend to over fit the training set. To circumvent this problem, we perform additional experiments and train on the balanced EPOS corpus. The results are reported in table 2. Not surprisingly, the RS baseline is inferior to any other system in the table. When word embedding features are combined to SF, we note improvement on most of the six modal verbs, and the best results are obtained with +WE<sub>ex'</sub>. The improvements of +WE<sub>ex'</sub> over SF are significant on *can*. Both results on balanced and unlabeled corpora reveal that task-specific embedding features on top of WE<sub>ex</sub> are superior to their counterparts without task-specific information, which shows the best performance overall.

Table 2: Accuracy with various features on EPOS+MPQA. + or - indicates that the improvements or degradations with respect to SF are statistically significant. The highest value in each row is marked in bold.

	RS	RR	ZH	SF	+WE <sub>av</sub>	+WE <sub>li</sub>	+WE <sub>ex</sub>	+WE <sub>ex'</sub>
can	33.3	57.8	60.4	65.6	64.4	64.7	64.7	<b>67.3</b> <sup>+</sup>
could	33.3	49.2	56.3	56.7	57.9	57.9	<b>60.0</b>	59.6
may	50.0	92.1	92.1	93.6	94.2	93.6	94.2	<b>95.0</b>
must	50.0	71.7	<b>85.6</b>	76.8	78.7	79.3	79.9	80.6
shall	50.0	53.9	53.9	<b>61.5</b>	46.2	46.2	46.2	53.8
should	50.0	76.3	88.3	<b>90.8</b>	<b>90.8</b>	<b>90.8</b>	<b>90.8</b>	<b>90.8</b>

## 5 Conclusion

In this paper, we make use of word embeddings learned from unlabeled data as additional features to more adequately represent the context of modal verbs. Our main experimental result is that the best weighting scheme for combining

pre-trained word embeddings into context embeddings is a corrected version of exponentially decaying weighting which attributes higher weights to the verb modified by the modal and its subject.

## References

- [1] Yang Liu, Xiaohui Yu, Zhongshuai Chen, and Bing Liu. Sentiment analysis of sentences with modalities. In *Proceedings of the 2013 International Workshop on Mining Unstructured Big Data Using Natural Language Processing*, UnstructureNLP '13, pages 39–44, 2013.
- [2] Roser Morante and Walter Daelemans. Learning the scope of hedge cues in biomedical texts. In *Proceedings of the Workshop on Current Trends in Biomedical Natural Language Processing*, BioNLP '09, pages 28–36, 2009.
- [3] Roser Sauri and James Pustejovsky. Are you sure that this happened? assessing the factuality degree of events in text. *Computational Linguistics*, 38(2):261–299, 2012.
- [4] Roberto Navigli. Word sense disambiguation: A survey. *ACM Computing Surveys*, 41(2):10:1–10:69, February 2009.
- [5] Josef Ruppenhofer and Ines Rehbein. Yes we can!? annotating english modal verbs. In *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC)*, pages 1538–1545, Istanbul, Turkey, 2012.
- [6] Mengfei Zhou, Anette Frank, Annemarie Friedrich, and Alexis Palmer. Semantically Enriched Models for Modal Sense Classification. In *Proceedings of the EMNLP Workshop LSDSem: Linking Models of Lexical, Sentential and Discourse-level Semantics*, pages 44–53, Lisbon, Portugal, 2015.
- [7] Ana Marasovic, Mengfei Zhou, Alexis Palmer, and Anette Frank. Modal sense classification at large: Paraphrase-driven sense projection, semantically enriched classification models and cross-genre evaluations. *Linguistic Issues in Language Technology, Special issue on Modality in Natural Language Understanding*, 14(3):1–58, 2016.
- [8] Ana Marasovic and Anette Frank. Multilingual modal sense classification using a convolutional neural network. In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 111–120, Berlin, Germany, 2016.
- [9] Ignacio Iacobacci, Mohammad Taher Pilehvar, and Roberto Navigli. Embeddings for word sense disambiguation: An evaluation study. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, ACL '16, pages 897–907, 2016.
- [10] Kaveh Taghipour and Hwee Tou Ng. Semi-supervised word sense disambiguation using word embeddings in general and specific domains. In *Proceedings of the 2015 Annual Conference of the North American Chapter of the ACL*, NAACL '15, pages 314–323, 2015.
- [11] Dayu Yuan, Julian Richardson, Ryan Doherty, Colin Evans, and Eric Altendorf. Semi-supervised word sense disambiguation with neural models. In *Proceedings of the 26th International Conference on Computational Linguistics*, COLING '16, pages 1374–1385, 2016.
- [12] Mikael Kågebäck and Hans Salomonsson. Word sense disambiguation using a bidirectional lstm. In *Proceedings of the 5th Workshop on Cognitive Aspects of the Lexicon*, pages 51–56, 2016.
- [13] Alessandro Raganato, Claudio Delli Bovi, and Roberto Navigli. Neural sequence learning models for word sense disambiguation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, EMNLP '17, pages 1156–1167, 2017.
- [14] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):27:1–27:27, May 2011.
- [15] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119. 2013.