

# Compound Embedding Features for Semi-supervised Learning

Mo Yu<sup>1</sup>, Tiejun Zhao<sup>1</sup>, Daxiang Dong<sup>2</sup>, Hao Tian<sup>2</sup> and Dianhai Yu<sup>2</sup>

Harbin Institute of Technology, Harbin, China

Baidu Inc., Beijing, China

{yumo, tjzhao}@mtlab.hit.edu.cn

{dongdaxiang, tianhao, yudianhai}@baidu.com

## Abstract

There has been a recent trend in discriminative methods of NLP to use representations of lexical items learned from unlabeled data as features, in order to overcome the problem of data sparsity. In this paper, we investigated the usage of word representations learned by neural language models, i.e. word embeddings. We built compound features of continuous word embeddings based on clustering to overcome the disadvantages of the direct usage of them, such as inadequacy with dealing word ambiguity and rare-words, large amount of computation and the problem of linear non-separability. Experiments showed that the compound features improved the performances on several NLP tasks while ran faster, suggesting the potential of embeddings.

## 1 Introduction

Supervised learning methods have achieved great successes in the field of *Natural Language Processing (NLP)*. However, these methods are usually limited by the problem of data sparsity in practice since it is impossible to obtain sufficient labeled data for any NLP tasks. In these situations semi-supervised learning methods will help, which could make use of both the labeled data and the easy-to-obtain unlabeled data.

A semi-supervised framework widely applied in NLP is to first learn word representations, which are feature vectors of lexical items, from unlabeled data and then plug them into a supervised system. These methods are very effective in utilizing large-scale unlabeled data and have successfully im-

proved performances of state-of-the-art supervised systems on a variety of tasks (Koo et al., 2008; Huang and Yates, 2009; Täckström et al., 2012).

With the development of *neural language models (NLM)* (Bengio et al., 2003; Mnih and Hinton, 2009), researchers become interested in word representations learned by those models (also called embeddings) recently. Word embeddings are dense, low dimensional and real-valued vectors. Embedding for each word is composed by some latent features, which are expected to capture useful syntactic and semantic properties. Word embeddings usually served as the first layer in deep learning systems for NLP (Collobert and Weston, 2008; Socher et al., 2011a, 2011b) and helped these systems perform comparably with state-of-the-art models based on handcrafted features. They have also been directly added as features to state-of-the-art models of chunking and NER, and achieved significant improvements (Turian et al. 2010).

Although the direct usage of continuous embeddings was proved an effective method for enhancing the state-of-the-art supervised models, it has some disadvantages, which made them outperformed by simpler Brown cluster features (Turian et al, 2010) and made them computationally complicated. Firstly, embeddings of rare words are insufficiently trained since they have not been updated many times and are close to their random initial values. As shown in (Turian et al, 2010), this is a main reason for the more errors made by models with embedding features than those with Brown cluster features. Secondly, in NLMs, each word has its unique representation, so it is difficult to represent different senses for an ambiguous word.

Thirdly, word embeddings are unsuitable for linear models in some tasks as will be proven in

Section 4.4. This is possibly because in these tasks, the target labels are correlated with combinations of different dimensions of word embeddings, or discriminative information may be coded in different intervals in a same dimension. So treating embeddings directly as inputs to a linear model could not make full use of them. We will have a deep insight on the reasons of the problem in the future.

In addition, since embeddings are dense vectors, it will introduce large amount of computations when they are directly used as inputs, making the method impractical. Such computations are constant for each data sample. For example, using a 64 dimension embedding as features in chunking, which has a original template of 20 entries, will introduce 3 times extra features.

In this paper, we first introduce the idea of clustering embeddings to overcome the last two disadvantages of the direct usage of word embeddings discussed above. The high-dimensional cluster features make samples from different classes better separated by linear models. And models with these features still run fast because the clusters are sparse and discrete which are easier to incorporate into discriminative methods.

Second, we propose the compound features based on clustering. Compound features, which are conjunctive features of neighboring words, have been widely used in NLP models for improving the performances because of their discriminative strength. Compound features of embeddings could also help a model better predict labels associated with rare-words and ambiguous words. For example when a rare-word has inaccurate embedding, embeddings of its context words are possibly accurate when they are not rare. So compound features made up of embeddings of nearby words can help to better describe the property of this rare-word. Compound features are difficult to build on dense embeddings. However they are easy to induce from the sparse embedding clusters proposed in this paper. Experiments on chunking and NER showed that based on the same embeddings, the compound features managed to achieve better performances.

In addition, although Brown clustering was considered better in (Turian et al 2010), our results and comparisons show that our compound features from embedding clustering is at least comparable with those from Brown clustering. Since embeddings can greatly benefit from the improvement and developing of deep learning in the future, we

believe that our proposed method has a large space of performance growth and will benefit more applications in NLP.

In the rest of this paper, Section 2 introduces how word embeddings were obtained and how the clusters of word embeddings were induced. In Section 3, we describe the evaluation tasks and the compound features for these tasks. Section 4 gives experimental results and analysis. Section 5 gives the conclusions.

## 2 Clustering of Word Embeddings

### 2.1 Learning Word Embeddings

Word embeddings in this paper were trained by NLMs (Bengio et al., 2003). The model predicts the scores of probabilities of words given their context information in the sentences.

When the NLM is used to predict the score of a word, it first converts current word and its context words (e.g.  $n-1$  words before it as in  $n$ -gram models) into their embeddings. Then these embeddings are put together and propagate forward on the network to compute the score of current word. After minimizing the loss on training data, embeddings are learned and can be further used as smoothing representations for words.

### 2.2 Clustering of embeddings

In order to get compound features of embeddings, we first induce discrete clusters from the embeddings. Concretely, the  $k$ -means clustering algorithm is used. Every word is treated as a single sample. A cluster is represented as the mean of the embeddings of words assigned to it. Similarities between words and clusters are measured by Euclidean distance. Since different numbers of  $k$ s contain information of different granularity, we combine clustering results achieved by different  $k$ s as features to better utilize the embeddings.

## 3 Supervised evaluation tasks

We built compound features based on clusters of embeddings for the same chunking and NER tasks in (Turian et al., 2010). The Brown cluster features were also used for comparison, which shared the same feature template used by clusters of embeddings. To compare with the work of (Turian et al,

2010), which aimed at the same purpose with ours but using embedding directly, we used the same word embeddings (CW 50) and Brown clusters (1000 clusters) they provided.

Moreover, we wish to find out whether our method extends well to languages other than English. So we conducted experiments on Chinese NER, where large amount of training data exists, making the accuracies difficult to improve much.

### 3.1 Chunking

The first application is chunking. The data in the CoNLL2000 shared task was used. The features used are shown in Table 1. The feature  $y_{-1}/y_0/c_{-1}/c_1$  in the last row is an example of compound features based on embedding clusters.

Words	$w_{i,i \in \{-2,2\}}, w_{i-1}/w_{i,i \in \{0,1\}}$
POS	$p_{i,i \in \{-2,2\}}, p_{i-1}/p_{i,i \in \{-1,2\}}$
Cluster	$c_{i,i \in \{-2,2\}}, c_{i-1}/c_{i,i \in \{0,1\}}, c_{-1}/c_1$
Transition	$y_{-1}/y_0/\{w_0, p_0, c_0, c_{-1}/c_1\}$

Table 1: Chunking features. Cluster features are suitable for both Brown clusters and embedding clusters. Symbol  $y_i$  is the tag predicted on word  $w_i$ .

The embeddings in (Turian et al, 2010) are trained on RCV data. Since the CoNLL2000 data is a part of the WSJ corpus, we believe that word representations trained on the same domain may better help to improve the results. So in the experiments, we also used embeddings and Brown clusters trained on unlabeled WSJ data from (Nivre et al, 2007) for comparison.

### 3.2 Named entity recognition

We evaluated our method on the tasks of NER for both English, for comparison with (Turian et al, 2010), and Chinese, for showing the extensibility of our method. For English, the standard evaluation benchmark in the CoNLL03 shared task was used. Table 2 shows the features used in the model.

Words	$w_{i,i \in \{-2,2\}}, w_{i-1}/w_{i,i \in \{0,1\}}$
Pre/suffix	$w_{0,i \in \{2,4\}}^{ji}, w_{0,i \in \{1,4\}}^{-i-1}$
Orthography	$Hyp(w_0), Cap(w_0)$
POS	$p_{i,i \in \{-2,2\}}, p_{i-1}/p_{i,i \in \{-1,2\}}$
Chunking	$b_{i,i \in \{-2,2\}}, b_{i-1}/b_{i,i \in \{-1,2\}}$
Cluster	$c_{i,i \in \{-2,2\}}, c_{i-1}/c_{i,i \in \{0,1\}}, c_{-1}/c_1$
Transition	$y_{-1}/y_0/\{w_0, p_0, c_0, c_{-1}/c_1\}$

Table 2: NER features. Hyp indicates if word contains hyphen and Cap indicates if first letter is capitalized.

For Chinese, we used data from People’s Daily (Jan.-Jun. 1998). The data was converted following the style of Penn CTB (Xue et al, 2005). Data from April was chosen as test set (1,309,616 words in 55,177 sentences), others for training (6,119,063 words in 255,951 sentences). The Chinese word representations were trained on Chinese Wikipedia until March 2011. The features used in Chinese NER are similar to those in English, except for the orthography, pre/suffixes, and chunking features.

## 4 Experiments

### 4.1 Experimental settings

We did little pre-processing work for the training of word representations on WSJ data. The datasets were tokenized and capital words were kept. For training of Chinese Wikipedia, we retained the bodies of all articles and replaced words with frequencies lower than 10 as an “UK\_WORD” token. On each dataset, we induced embeddings with 64 dimensions based on 7-gram models and 1000 Brown clusters. Training of NLMs is very slow when the dataset has a large vocabulary. In this paper, the method in (Schwenk, 2007) was used to accelerate the training processes. The NLMs were trained for 5 epochs on each dataset.

For clustering of embeddings we choose  $k=500$  and 2500 since such combination performed best on development set as shown in the next section. We chose the Sofia-ml toolkit (Sculley 2010) for clustering of embeddings in order to save time.

In the experiments CRF models were used and were optimized by ASGD (implemented by Léon Bottou) for shortening the periods of experiments. For comparison we re-implemented the direct usage of embeddings in (Turian et al, 2010) with CRFsuite (Okazaki, 2007) since their features contain continuous values.

### 4.2 Performances

Table 3 shows the chunking results. The results reported in (Turian et al. 2010) were denoted as “direct”. Based on the same word representations, our compound features got better performances in all cases. The embedding features trained on unlabeled WSJ data yield further improvements, showing that word representations from similar domains can better help the supervised tasks.

System	Direct	Compound
Baseline	93.75	
+Embedding (RCV)	94.10	94.19
+Brown (RCV)	94.11	94.24
+Brown&Emb (RCV)	94.35	<b>94.42</b>
+Embedding (WSJ)	94.20	94.37
+Brown (WSJ)	94.25	94.36
+Brown&Emb (WSJ)	94.43	<b>94.58</b>

Table 3: F1-scores of chunking

In the experiments of NER, first we evaluated how the numbers of clusters  $k$  will affect the performances on development set (Figure 1). The results showed that both the cluster features (excluding all compound embedding features) and compound features could achieve better results than direct usage of the same embeddings. It is also found that the performances did not vary much when  $k$  was between 500 and 3000. When  $k=2500$ , the result was a little higher than others. We finally chose combination of  $k=500$  and 2500, which achieved best results on development set.

The performances of NER on test set are shown in Table 4. Our baseline is a bit lower than that in (Turian et al, 2010), since the first-order CRF could not utilize more context information of NE tags. Despite of this, the same conclusions as in chunking held.

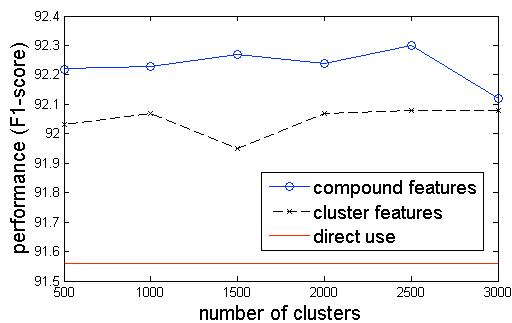


Figure 1: Relation between numbers of clusters  $k$  and performances on development set.

System	Direct	Compound
Baseline	83.78	
+Embedding	87.38	88.46
+Brown	88.14	88.23
+Brown&Embedding	88.85	89.06

Table 4: F1-scores of English NER on test data

Performances on Chinese NER are shown in Table 5. Word representation features could greatly improve the performances. One of the reasons is that in Chinese, person names and location names are usually single words. So both Brown clustering and NLM can learn good representations for them.

System	Direct	Compound
Baseline	88.24	
+Embedding	89.98	90.37
+Brown	90.24	90.55
+Brown&Embedding	90.66	<b>90.96</b>

Table 5: F1-scores of Chinese NER on test data

Above results gave evidences that although the clustering of embeddings may loss information, the derived compound features did have better performances. The compound features can also improve the performances of Brown clusters, but not as much as they did on embeddings. And the combination of embedding-clusters and Brown-clusters could further improve the performances, since they made use of different type of context information.

The compound features also reduced the time cost of using embedding features. For example, the time for tagging one sentence in English NER was reduced to 1.6 ms from 5.6 ms as shown in Table 6.

Embedding	Time (ms)
Baseline	1.2
Embeddings (direct)	5.6
Embeddings (compound)	1.6

Table 6: Running time of different features

### 4.3 Analysis on rare-words and ambiguous words

Our compound embedding features greatly outperformed the direct usage of same embeddings on English NER. One of the reasons for the improvements is that these features worked better on rare-words and ambiguous words as we hypothesized.

To show the compound features have stronger abilities to handle rare words, we counted the numbers of errors made on words with different frequencies on unlabeled data. Here the word frequencies are from the results of Brown clustering provided by (Turian et al. 2010). We compared our compound embedding features with direct usage of embeddings as well as Brown clusters, which is believed to work better on rare words. Figure 2(a) shows that the compound features indeed resulted in fewer errors than the two baseline methods in most cases. Errors of embeddings occurred on words with frequencies lower than 2K and those in the range of 16 to 256 were reduced by 10.55% and 24.44%, respectively.

Our compound features also reduced the errors caused by ambiguous words, as shown in Figure 2(b), where the numbers of senses for a word are

measured by the numbers of different POS tags it has in Penn Treebank. 12.1% of the errors on ambiguous words were reduced, comparing to 8.4% of the errors on unambiguous ones.

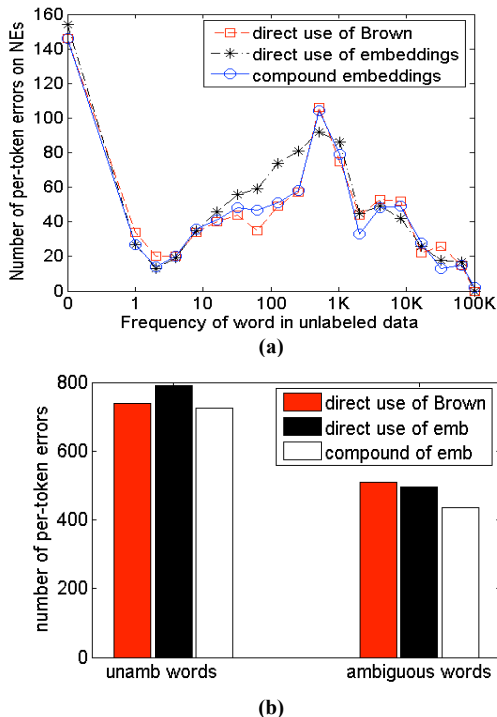


Figure 2: Errors incurred on words with different frequencies (a) and ambiguous words (b) in NER.

#### 4.4 Analysis on linear separability of embeddings in NER

Another reason for the good performances of compound features on NER is that they made linear models better separate *named entities (NEs)* and *non-NEs*, which are more difficult to be linearly separated when embeddings are directly used as features. Here we designed an experiment to prove this. Based on training data of CoNLL2003, a classification task was built to tell whether a word belongs to NE or not. Linear SVM and a non-linear model *Multilayer Perceptron (MLP)* were used to build the classifiers. As shown in Table 7, when embeddings were directly used as features, MLP performed much better than linear SVM. And the linear model was under-fitting on this task since it had similar accuracies on both training set and development set. Above observations showed that linear models could not separate NEs and non-NEs well in the space of embeddings.

When clusters of embeddings were used as features, the accuracies of linear models increased

even when there were only one or two non-zero features for each sample. At the same time the performances of MLP decreased because of the loss of information during clustering. The gaps between accuracies of linear models and non-linear ones decreased in the spaces of clusters, showing that cluster features are more suitable for linear models. At last, the compound features made the linear model out-perform all non-linear ones, since extra context information could be utilized.

Embeddings	Models	Accuracy
direct	linear	94.38
direct	MLP	96.87
cluster 1000	linear	95.31
cluster 1000	MLP	95.32
cluster 500+2500	linear	96.10
cluster 500+2500	MLP	96.02
compound	linear	<b>97.30</b>

Table 7: Performances of linear and non-linear models on development set with different embedding features.

## 5 Conclusion and perspectives

In this paper, we first introduced the idea of clustering embeddings and then proposed the compound features based on clustering, in order to overcome the disadvantages of the direct usage of continuous embeddings. Experiments showed that the compound features built on the same original word representation features (either embeddings or Brown clusters) achieve better performances on the same tasks. Further analyses showed that the compound features reduced errors on rare-words and ambiguous words and could be better utilized by linear models.

The usage of word embeddings also has some limitations, e.g. they are weak in capturing structural information of languages, which is necessary in NLP. In the future, we will research on task-specific representations for sub-structures, such as phrases and sub-trees, based on word embeddings.

## Acknowledgments

We would like to thank Dr. Hua Wu, Haifeng Wang, Jie Zhou and Rui Zhang for many discussions and thank the anonymous reviewers for their valuable comments and helpful suggestions. This work was supported by National Natural Science Foundation of China (61173073), and the Key Project of the National High Technology Research and Development Program of China (2011AA01A207).

## References

- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language models. *Journal of Machine Learning Research*, 3:1137–1155.
- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM.
- Finkel, J., Grenager, T., and Manning, C. (2005). Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 363–370. Association for Computational Linguistics.
- Huang, F. and Yates, A. (2009). Distributional representations for handling sparsity in supervised sequence labeling. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 495–503. Association for Computational Linguistics.
- Koo, T., Carreras, X., and Collins, M. (2008). Simple semi-supervised dependency parsing.
- Mnih, A., & Hinton, G. E. (2009). A scalable hierarchical distributed language model. *NIPS*. pages 1081–1088.
- Nivre, J., Hall, J., K<sup>u</sup>bler, S., McDonald, R., Nilsson, J., Riedel, S., and Yuret, D. (2007). The CoNLL 2007 shared task on dependency parsing. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL*, pages 915–932.
- Okazaki, N. (2007). Crfsuite: a fast implementation of conditional random fields (crfs). URL <http://www.chokkan.org/software/crfsuite>.
- Schwenk, H. (2007). Continuous space language models. *Computer Speech & Language*, 21(3):492–518.
- Sculley, D. (2010). Web-scale k-means clustering. In *Proceedings of the 19th international conference on World Wide Web*, pages 1177–1178. ACM.
- Socher, R., Huang, E., Pennington, J., Ng, A., and Manning, C. (2011a). Dynamic pooling and unfolding recursive auto-encoders for paraphrase detection. *Advances in Neural Information Processing Systems*, 24:801–809.
- Socher, R., Pennington, J., Huang, E., Ng, A., and Manning, C. (2011b). Semi-supervised recursive auto-encoders for predicting sentiment distributions. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 151–161. Association for Computational Linguistics.
- Täckström, O., McDonald, R., and Uszkoreit, J. (2012). Cross-lingual word clusters for direct transfer of linguistic structure.
- Turian, J., Ratinov, L., and Bengio, Y. (2010). Word representations: a simple and general method for semi-supervised learning. *Urbana*, 51:61801.
- Xue, N., Xia, F., Chiou, F., and Palmer, M. (2005). The penn chinese treebank: Phrase structure annotation of a large corpus. *Natural Language Engineering*, 11(2):207.