



Sparse Word Embeddings Using ℓ_1 Regularized Online Learning

Fei Sun, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng

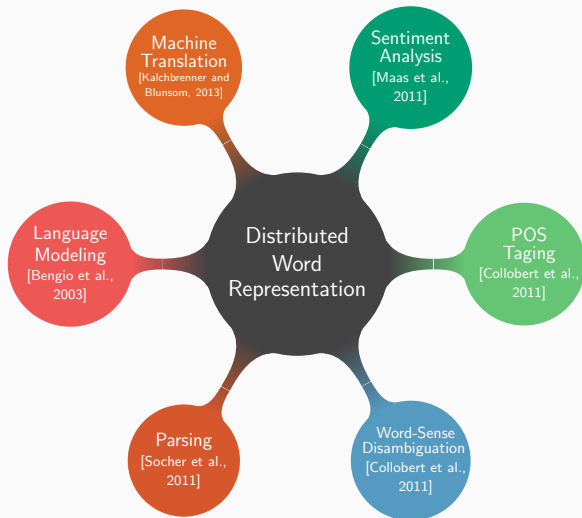
July 14, 2016

ofey.sunfei@gmail.com, {guojiafeng, lanyanyan, junxu, cxq}@ict.ac.cn

CAS Key Lab of Network Data Science and Technology

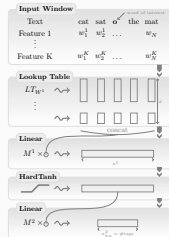
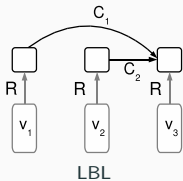
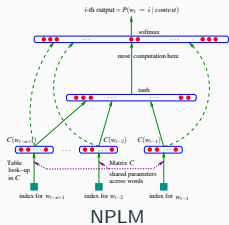
Institute of Computing Technology, Chinese Academy of Sciences

Distributed Word Representation

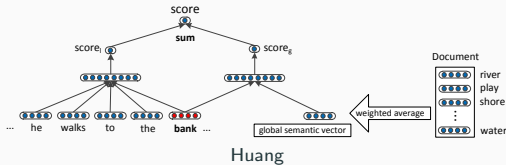
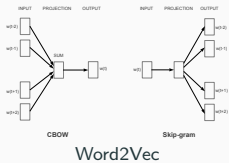


Distributed word representation is so hot in NLP community.

Models



C&W



State-Of-The-Art: CBOW and SG.

Dense Representation and Interpretability

Example¹

man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
dog	[0.0512695, ..., -0.306641, ..., 0.222656]
computer	[0.107422, ..., -0.0375977, ..., -0.0620117]

¹Vectors from GoogleNews-vectors-negative300.bin.

Dense Representation and Interpretability

Example¹

man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
dog	[0.0512695, ..., -0.306641, ..., 0.222656]
computer	[0.107422, ..., -0.0375977, ..., -0.0620117]

- Which dimension represents the gender of *man* and *woman*?

¹Vectors from GoogleNews-vectors-negative300.bin.

Dense Representation and Interpretability

Example¹

man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
dog	[0.0512695, ..., -0.306641, ..., 0.222656]
computer	[0.107422, ..., -0.0375977, ..., -0.0620117]

- Which dimension represents the gender of *man* and *woman*?
- What sort of value indicates male or female?

¹Vectors from GoogleNews-vectors-negative300.bin.

Dense Representation and Interpretability

Example¹

man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
dog	[0.0512695, ..., -0.306641, ..., 0.222656]
computer	[0.107422, ..., -0.0375977, ..., -0.0620117]

- Which dimension represents the gender of *man* and *woman*?
- What sort of value indicates male or female?
- Gender dimension(s) would be active in all the word vectors including irrelevant words like *computer*.

¹Vectors from GoogleNews-vectors-negative300.bin.

Dense Representation and Interpretability

Example¹

man	[0.326172, ..., 0.00524902, ..., 0.0209961]
woman	[0.243164, ..., -0.205078, ..., -0.0294189]
dog	[0.0512695, ..., -0.306641, ..., 0.222656]
computer	[0.107422, ..., -0.0375977, ..., -0.0620117]

- Which dimension represents the gender of *man* and *woman*?
- What sort of value indicates male or female?
- Gender dimension(s) would be active in all the word vectors including irrelevant words like *computer*.
- Difficult in interpretation and uneconomic in storage.

¹Vectors from GoogleNews-vectors-negative300.bin.

Sparse Word Representation

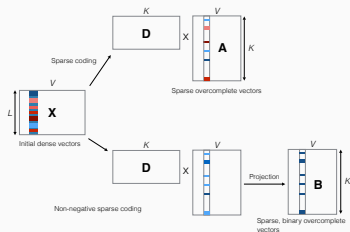
Non-Negative Sparse Embedding (NNSE) [Murphy et al., 2012]

$$\arg \min_{\substack{\mathbf{A} \in \mathbb{R}^{m \times k} \\ \mathbf{D} \in \mathbb{R}^{k \times n}}} \sum_{i=1}^m \left(\|\mathbf{X}_i - \mathbf{A}_i \times \mathbf{D}\|^2 + \lambda \|\mathbf{A}_i\|_1 \right)$$

where, $\mathbf{A}_{i,j} \geq 0, \forall 1 \leq i \leq m, \forall 1 \leq j \leq k$
 $\mathbf{D}_i \mathbf{D}_i^T \leq 1, \forall 1 \leq i \leq k$

Introduce sparse and non-negative constraints into MF.

Sparse Coding (Word2Vec) [Faruqui et al., 2015]



Convert dense vector using sparse coding in a post-processing way.

Sparse Word Representation

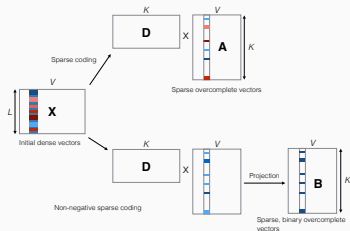
Non-Negative Sparse Embedding (NNSE) [Murphy et al., 2012]

$$\arg \min_{\substack{\mathbf{A} \in \mathbb{R}^{m \times k} \\ \mathbf{D} \in \mathbb{R}^{k \times n}}} \sum_{i=1}^m \left(\|\mathbf{X}_i - \mathbf{A}_i \times \mathbf{D}\|^2 + \lambda \|\mathbf{A}_i\|_1 \right)$$

where, $\mathbf{A}_{i,j} \geq 0, \forall 1 \leq i \leq m, \forall 1 \leq j \leq k$
 $\mathbf{D}_i \mathbf{D}_i^T \leq 1, \forall 1 \leq i \leq k$

Introduce sparse and non-negative constraints into MF.

Sparse Coding (Word2Vec) [Faruqui et al., 2015]



Convert dense vector using sparse coding in a post-processing way.

They are difficult to train on large-scale data for heavy memory usage!

Our Motivation

Good Performance


Fast

Large Scale

Interpretable

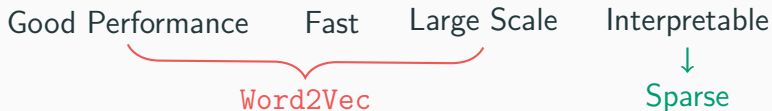
Our Motivation

Good Performance Fast Large Scale Interpretable

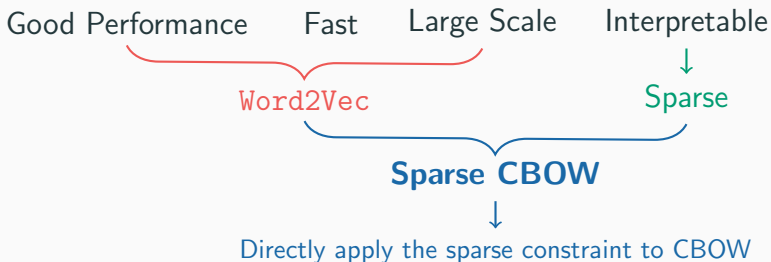


Word2Vec

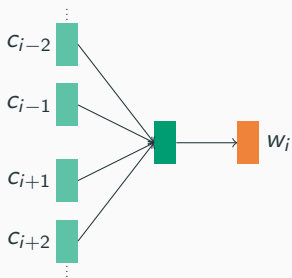
Our Motivation



Our Motivation



CBOW



CBOW

$$\mathcal{L}_{cbow} = \sum_{i=1}^N \log p(w_i | h_i)$$

$$p(w_i | h_i) = \frac{\exp(\vec{w}_i \cdot \vec{h}_i)}{\sum_{w \in \mathcal{W}} \exp(\vec{w} \cdot \vec{h}_i)}$$

$$\vec{h}_i = \frac{1}{2l} \sum_{\substack{j=i-l \\ j \neq i}}^{i+l} \vec{c}_j$$

$$\mathcal{L}_{cbow}^{ns} = \sum_{i=1}^N \left(\log \sigma(\vec{w}_i \cdot \vec{h}_i) + k \cdot \mathbf{E}_{\vec{w} \sim P_{\vec{w}}} \log \sigma(-\vec{w} \cdot \vec{h}_i) \right)$$

$$\mathcal{L}_{\text{s-cbow}}^{\text{ns}} = \mathcal{L}_{\text{cbow}}^{\text{ns}} - \lambda \sum_{w \in W} \|\vec{w}\|_1$$

$$\mathcal{L}_{\text{s-cbow}}^{\text{ns}} = \mathcal{L}_{\text{cbow}}^{\text{ns}} - \lambda \sum_{w \in W} \|\vec{w}\|_1$$

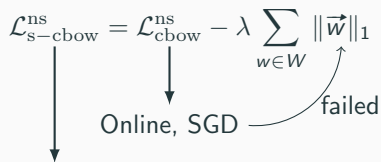
↓
Online, SGD

Sparse CBOW

$$\mathcal{L}_{\text{s-cbow}}^{\text{ns}} = \mathcal{L}_{\text{cbow}}^{\text{ns}} - \lambda \sum_{w \in W} \|\vec{w}\|_1$$

↓
Online, SGD

failed

$$\mathcal{L}_{s\text{-cbow}}^{\text{ns}} = \mathcal{L}_{\text{cbow}}^{\text{ns}} - \lambda \sum_{w \in W} \|\vec{w}\|_1$$


Regularized Dual Averaging (RAD) [Xiao, 2009]

Truncating using online average subgradients

RDA algorithm for Sparse CBOW

- 1: **procedure** SPARSECBOW(\mathcal{C})
- 2: **Initialize:** $\vec{w}, \forall w \in W, \vec{c}, \forall c \in \mathcal{C}, \vec{g}_{\vec{w}}^0 = \vec{0}, \forall w \in W$
- 3: **for** $i = 1, 2, 3, \dots$ **do**
- 4: $t \leftarrow$ update time of word w_i
- 5:
$$\vec{h}_i = \frac{1}{2l} \sum_{\substack{j=i-l \\ j \neq i}}^{i+l} \vec{c}_j$$
- 6:
$$\vec{g}_{\vec{w}_i}^t = \left[\mathbb{1}_{h_i}(w_i) - \sigma(\vec{w}_i^t \cdot \vec{h}_i) \right] \vec{h}_i$$
- 7:
$$\vec{g}_{\vec{w}_i}^t = \frac{t-1}{t} \vec{g}_{\vec{w}_i}^{t-1} + \frac{1}{t} \vec{g}_{\vec{w}_i}^t \quad \triangleright \text{Keeping track of the online average subgradients}$$
- 8: Update \vec{w}_i element-wise according to
- 9:
$$\vec{w}_{ij}^{t+1} = \begin{cases} 0 & \text{if } |\vec{g}_{\vec{w}_{ij}}^t| \leq \frac{\lambda}{\#(w_i)}, \\ \eta t (\vec{g}_{\vec{w}_{ij}}^t - \frac{\lambda}{\#(w_i)} \text{sgn}(\vec{g}_{\vec{w}_{ij}}^t)) & \text{otherwise,} \end{cases} \quad \triangleright \text{Truncating}$$

where, $j = 1, 2, \dots, d$
- 10: **for** $k = -l, \dots, -1, 1, \dots, l$ **do**
- 11: update \vec{c}_{i+k} according to
- 12:
$$\vec{c}_{i+k} := \vec{c}_{i+k} + \frac{\alpha}{2l} \left[\mathbb{1}_{h_i}(w_i) - \sigma(\vec{w}_i^t \cdot \vec{h}_i) \right] \vec{w}_i^t$$
- 13: **end for**
- 14: **end for**
- 15: **end procedure**

Evaluation

Baseline

- Dense representation models
 - GloVe [Pennington et al., 2014]
 - CBOW and SG [Mikolov et al., 2013]
- Sparse representation models
 - Sparse Coding (SC) [Faruqui et al., 2015]
 - Positive Pointwise Mutual Information (PPMI) [Bullinaria and Levy, 2007]
 - NNSE [Murphy et al., 2012]

Tasks

- Interpretability
 - Word Intrusion
- Expressive Power
 - Word Analogy
 - Word Similarity

Experimental Settings

Corpus: Wikipedia 2010 (1B words)

Parameters Setting:

window	negative	iteration	λ	learning rate	noise distribution
10	10	20	grid search	0.05	$\propto \#(w)^{0.75}$

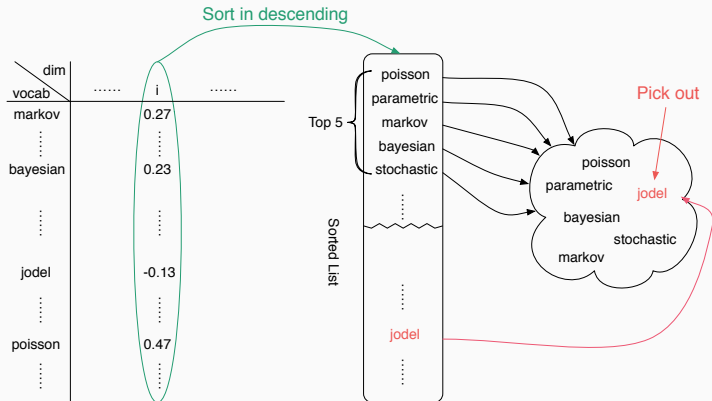
Baseline:

Model	Setting
GloVe, CBOW, SG	same setting with released tools
SC	Embeddings of CBOW as input
NNSE	PPMI matrix, 4,000 words, SPAMS ¹

¹<http://spams-devel.gforge.inria.fr>

Interpretability

Word Intrusion [Chang et al., 2009]



1. Sort dimension i in descending order.
2. A set: {top 5 words, 1 word (bottom 50% in i & top 10% in $j, j \neq i$)}
3. Pick out the intruder word

More interpretable, more easy to pick out.

Evaluation Metric

traditional metric: human assignment, **subjective, costly**.

Definition

$$\text{IntraDist}_i = \sum_{w_j \in \text{top}_k(i)} \sum_{\substack{w_k \in \text{top}_k(i) \\ w_k \neq w_j}} \frac{\text{dist}(w_j, w_k)}{k(k-1)}$$

$$\text{InterDist}_i = \sum_{w_j \in \text{top}_k(i)} \frac{\text{dist}(w_j, w_{b_i})}{k}$$

$$\text{DistRatio} = \frac{1}{d} \sum_{i=1}^d \frac{\text{InterDist}_i}{\text{IntraDist}_i}$$

IntraDist_i: Average distance between top k words

InterDist_i: Average distance between top k words and bottom word

Intuition: The intruder word should be dissimilar to the top words while those top words should be similar to each other.

Word Intrusion Results

Table 1: 300 dimension, running 10 times.

Model	Sparsity	DistRatio
GloVe	0%	1.07
CBOW	0%	1.09
SG	0%	1.12
NNSE (PPMI)	89.15%	1.55
SC (CBOW)	88.34%	1.24
Sparse CBOW	90.06%	1.39

- Sparse is better than dense.
- Sparse CBOW **VS.** SC (CBOW): information loss in a separate sparse coding step.
- Non-negative constraint might also be a good choice.

Case Study

Model	Top 5 Words
CBOW	beat, finish, wedding, prize, read
	rainfall, footballer, breakfast, weekdays, angeles
	landfall, interview, asked, apology, dinner
	becomes, died, feels, resigned, strained
	best, safest, iucn, capita, tallest
Sparse CBOW	poisson, parametric, markov, bayesian, stochastic statistical learning
	ntfs, gzip, myfile, filenames, subdirectories file system
	hugely, enormously, immensely, wildly, tremendously adverb for degree
	earthquake, quake, uprooted, levees, spectacularly disasters
	bosons, accretion, higgs, neutrinos, quarks particles

The dimensions of Sparse CBOW reveal some clear and consistent semantic meanings.

Expressive Power

Expressive Power

TASK	Word Analogy	Word Similarity
TESTSET	Google [Mikolov et al., 2013]	Rare Word (RW) [Luong et al., 2013] WordSim-353 (WS-353) [Finkelstein et al., 2002] SimLex-999 (SL-999) [Hill et al., 2015]
EXAMPLE	Beijing: China ~ Paris: ? big: bigger ~ deep:?	(tiger cat 7.35)
SOLUTION	3COSMUL [Levy and Goldberg, 2014]	cos
METRIC	Percentage	Spearman Rank Correlation

Results

Table 2: Precision (%) for analogy and spearman correlation for similarity.

Model	Dim	Sparsity	Sem	Syn	Total	WS-353	SL-999	RW
GloVe	300	0%	79.31	61.48	69.57	59.18	32.35	34.13
CBOW	300	0%	79.38	68.80	73.60	67.21	38.82	45.19
SG	300	0%	77.79	67.32	72.09	70.74	36.07	45.55
PPMI(W-C)	40000	86.55%	74.02	38.99	53.02	62.35	24.10	30.45
PPMI(W-C)	388723	99.61%	58.55	31.19	43.60	58.99	23.01	27.98
NNSE (PPMI) ²	300	89.15%	29.89	27.68	28.56	68.61	27.60	41.82
SC (CBOW)	300	88.34%	28.99	28.43	28.68	59.85	30.44	38.75
SC (CBOW)	3000	95.85%	74.71	61.24	67.35	68.22	39.12	44.75
Sparse CBOW	300	90.06%	73.24	67.48	70.10	68.29	44.47	42.30

- Sparse CBOW is a competitive model using much less memory.
- Similar performance on analogy if we reduce its sparsity level $< 85\%$.

²The input matrix of NNSE is the 40000-dimensional representations of PPMI in fourth row.

- A sparse word representation model.
- A new evaluation metric for word intrusion task.
- Improvement in word vector interpretability.
- Similar performance with less memory usage.

Thanks!

Q & A

References I



Bullinaria, J. A. and Levy, J. P. (2007).

Extracting semantic representations from word co-occurrence statistics: A computational study.

Behavior Research Methods, 39(3):510–526.



Chang, J., Gerrish, S., Wang, C., Boyd-graber, J. L., and Blei, D. M. (2009).

Reading tea leaves: How humans interpret topic models.

In *NIPS*, pages 288–296. Curran Associates, Inc.



Faruqui, M., Tsvetkov, Y., Yogatama, D., Dyer, C., and Smith, N. A. (2015).

Sparse overcomplete word vector representations.

In *Proceedings of ACL*, pages 1491–1500.



Finkelstein, L., Gabrilovich, E., Matias, Y., Rivlin, E., and Gadi Wolfman, Z. S., and Ruppin, E. (2002).

Placing search in context: The concept revisited.

ACM Trans. Inf. Syst., 20(1):116–131.

References II



Hill, F., Reichart, R., and Korhonen, A. (2015).
Simlex-999: Evaluating semantic models with (genuine) similarity estimation.
Computational Linguistics, pages 665–695.



Levy, O. and Goldberg, Y. (2014).
Linguistic regularities in sparse and explicit word representations.
In *Proceedings of CoNLL*, pages 171–180.



Luong, M.-T., Socher, R., and Manning, C. D. (2013).
Better word representations with recursive neural networks for morphology.
In *Proceedings of CoNLL*, pages 104–113.



Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013).
Efficient estimation of word representations in vector space.
In *Proceedings of ICLR*.



Murphy, B., Talukdar, P., and Mitchell, T. (2012).
Learning effective and interpretable semantic models using non-negative sparse embedding.
In *Proceedings of COLING*, pages 1933–1950.

References III



Pennington, J., Socher, R., and Manning, C. D. (2014).

Glove: Global vectors for word representation.

In *Proceedings of EMNLP*, pages 1532–1543.



Xiao, L. (2009).

Dual averaging method for regularized stochastic learning and online optimization.

In *Proceedings of NIPS*, pages 2116–2124.