

CS 6501 Natural Language Processing

Word Embeddings

Yangfeng Ji

October 24, 2019

Department of Computer Science
University of Virginia



ENGINEERING

1. Word Embeddings: Skip-gram
2. Word Embedding: GloVe
3. Evaluation Methods
4. Problems

Word Embeddings: Skip-gram

The objective function of a skip-gram model is defined as

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq i \leq c; i \neq 0} \log p(w_{t+i} \mid w_t) \quad (1)$$

- ▶ $\log p(w_{t+i} \mid w_t) = \mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t})$
- ▶ In practice, the vocab size could be 10K, 50K or even bigger, the computation of the log-sum-exp is prohibitively expensive

Negative Sampling

Replace

$$\log p(w_{t+i} \mid w_t) = \mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t})$$

with the following function as objective

$$\log \sigma(\mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t}) - \sum_{i=1}^k E_{w' \sim p_n(w)} \left[\log \sigma(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t}) \right] \quad (2)$$

where k is the number of negative samples

Basic Training Procedure

Example with $t = 6$, $i = 1$, and $k = 3$

... finding a better word representation ...

w_6	w_7	negative samples
better	word	larger cause window

Basic Training Procedure

Example with $t = 6$, $i = 1$, and $k = 3$

... finding a better word representation ...

w_6	w_7	negative samples
better	word	larger cause window

For a given word w_t and i

1. Treat its neighboring context word w_{t+i} as positive example
2. Randomly sample k **other** words from the vocab as negative examples
3. Optimize Equation 2 to update both v . and u .

Two Factors in Negative Sampling

$$\log \sigma(\mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t}) - \sum_{i=1}^k E_{w' \sim p_n(w)} \left[\log \sigma(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t}) \right] \quad (3)$$

Two Factors in Negative Sampling

$$\log \sigma(\mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t}) - \sum_{i=1}^k E_{w' \sim p_n(w)} \left[\log \sigma(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t}) \right] \quad (3)$$

Two factors [Mikolov et al., 2013a]

- ▶ $k = ?$
 - ▶ $5 \leq k \leq 20$ works better for small datasets
 - ▶ $2 \leq k \leq 5$ is enough for large datasets

Two Factors in Negative Sampling

$$\log \sigma(\mathbf{u}_{w_{t+i}}^\top \mathbf{v}_{w_t}) - \sum_{i=1}^k E_{w' \sim p_n(w)} \left[\log \sigma(\mathbf{u}_{w'}^\top \mathbf{v}_{w_t}) \right] \quad (3)$$

Two factors [Mikolov et al., 2013a]

- ▶ $k = ?$
 - ▶ $5 \leq k \leq 20$ works better for small datasets
 - ▶ $2 \leq k \leq 5$ is enough for large datasets
- ▶ Noisy distribution $p_n(w)$
 - ▶ $p_n(w) \propto \text{unigram-distribution}(w)^{\frac{3}{4}}$

Examples: Words and their Neighbors

The same Yelp dataset, with $k = 50$

yummy	horrible
delicious	terrible
tasty	poor
delish	awful
yum	<i>customer</i>
incredible	exceptional
superb	bad
phenomenal	astonished
fantastic	<i>pleasant</i>
<i>disappoint</i>	<i>happier</i>
awesome	<i>zero</i>

Word Embedding: GloVe

The motivation of GloVe [Pennington et al., 2014] is to find a balance between the methods based on

- ▶ global matrix factorization (e.g., LSA) and
- ▶ local context windows (e.g., Skip-gram).

Word-to-word Co-occurrence Matrix

- Define \mathbf{X} with $X_{i,j}$ denotes the frequency of word j appears in the context of word i

$$\mathbf{X} = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ X_{i,1} & \dots & X_{i,j-1} & X_{i,j} & X_{i,j+1} & \dots & X_{i,V} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (4)$$

Each row corresponds one target word, each column corresponds one context word.

Word-to-word Co-occurrence Matrix

- Define \mathbf{X} with $X_{i,j}$ denotes the frequency of word j appears in the context of word i

$$\mathbf{X} = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ X_{i,1} & \dots & X_{i,j-1} & X_{i,j} & X_{i,j+1} & \dots & X_{i,V} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (4)$$

Each row corresponds one target word, each column corresponds one context word.

- Empirical probability estimation of w_j given w_i

$$Q(w_j | w_i) = \frac{X_{ij}}{X_i} \quad (5)$$

where $X_i = \sum_j X_{i,j}$

Another way to estimate the probability of w_j given w_i is

$$P(w_j \mid w_i) = \frac{\exp(\mathbf{u}_{w_j}^\top \mathbf{v}_{w_i})}{\sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i})} \quad (6)$$

with \mathbf{u} . and \mathbf{v} . are two sets of parameters (embeddings) associated with words, similar to the Skip-gram model.

The basic idea is to learn $\{v.\}$ and $\{u.\}$, such that

$$Q(w_j | w_i) \approx P(w_j | w_i) \quad (7)$$

or

$$\log Q(w_j | w_i) \approx \log P(w_j | w_i) \quad (8)$$

The basic idea is to learn $\{\mathbf{v}.\}$ and $\{\mathbf{u}.\}$, such that

$$Q(w_j \mid w_i) \approx P(w_j \mid w_i) \quad (7)$$

or

$$\log Q(w_j \mid w_i) \approx \log P(w_j \mid w_i) \quad (8)$$

More specific

$$\log(X_{ij}) - \log(X_i) \approx \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \quad (9)$$

Starting point:

$$\log(X_{ij}) - \log(X_i) \approx \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \quad (10)$$

Starting point:

$$\log(X_{ij}) - \log(X_i) \approx \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \quad (10)$$

In order to find the best approximation, we could formulate this as a optimization problem

$$\left\{ \log(X_{ij}) - \log(X_i) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} + \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \right\}^2 \quad (11)$$

Starting point:

$$\log(X_{ij}) - \log(X_i) \approx \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} - \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \quad (10)$$

In order to find the best approximation, we could formulate this as a optimization problem

$$\left\{ \log(X_{ij}) - \log(X_i) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} + \log \sum_{w' \in \mathcal{V}} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_{w_i}) \right\}^2 \quad (11)$$

It can be further simplified as (Eq. 16 in [Pennington et al., 2014])

$$\left\{ \log(X_{ij}) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} \right\}^2 \quad (12)$$

if we only consider the **unnormalized** version of P and Q .

Objective Function

The overall objective function is defined as

$$\sum_{w_i} \sum_{w_j} (\log(X_{ij}) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i})^2 \quad (13)$$

The overall objective function is defined as

$$\sum_{w_i} \sum_{w_j} (\log(X_{ij}) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i})^2 \quad (13)$$

The objective function is further refined by discouraging high-frequency words as

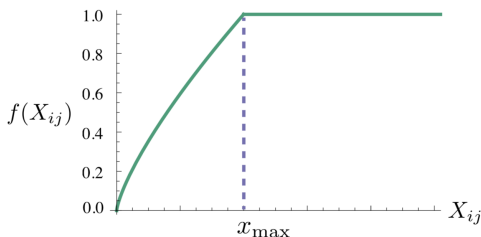
$$\sum_{w_i} \sum_{w_j} f(X_{ij}) (\log(X_{ij}) - \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i})^2 \quad (14)$$

Down-weighting

Weighting function:

$$f(x) = \begin{cases} (\frac{x}{x_{\max}})^a & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases} \quad (15)$$

where $a = 3/4$.



Skip-gram as Implicit Matrix Factorization

[Levy and Goldberg, 2014] shows that skip-gram with negative sampling can be viewed as an implicit matrix factorization over a word-word co-occurrence matrix weighted by point-wise mutual information (PMI).

$$\mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} \approx \text{PMI}(w_i, w_j) - \log k \quad (16)$$

where $\text{PMI}(w_i, w_j)$ is the mutual information of $P(w_i)$ and $P(w_j)$ with *a given window size* and k is the number of negative samples.

Skip-gram as Implicit Matrix Factorization (II)

The definition of $\text{PMI}(w_i, w_j)$ is

$$\text{PMI}(w_i, w_j) = \log \frac{P(w_i, w_j)}{P(w_i)P(w_j)} = \log P(w_j \mid w_i) - \log P(w_j) \quad (17)$$

Skip-gram as Implicit Matrix Factorization (II)

The definition of $\text{PMI}(w_i, w_j)$ is

$$\text{PMI}(w_i, w_j) = \log \frac{P(w_i, w_j)}{P(w_i)P(w_j)} = \log P(w_j \mid w_i) - \log P(w_j) \quad (17)$$

Combine 16 and 17, we have

$$\begin{aligned} \mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} &\approx \log \frac{P(w_i, w_j)}{P(w_i)P(w_j)} - \log k \\ &= \log P(w_j \mid w_i) - \log P(w_j) - \log k \\ &= \log(X_{ij}) - \log(X_i) - \log(X_j) + \log D - \log k \end{aligned} \quad (18)$$

Similar to Eq. 8 in [Pennington et al., 2014].

Essentially,

A unified framework

$$\mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} \approx \log(X_{ij}) + g(\mathbf{X}) \quad (19)$$

A unified framework

$$\mathbf{u}_{w_j}^\top \mathbf{v}_{w_i} \approx \log(X_{ij}) + g(\mathbf{X}) \quad (19)$$

Which one matters?

- ▶ $g(\mathbf{X})$, or
- ▶ Implicit/explicit optimization, or
- ▶ Other tricks (down-sampling, hyper-parameters, etc.)

Evaluation Methods

- ▶ Intrinsic Evaluation
 - ▶ Word similarity
 - ▶ Word analogy
 - ▶ Word intrusion
- ▶ Extrinsic Evaluation

Let w_i and w_j be two words, and \mathbf{v}_{w_i} and \mathbf{v}_{w_j} be the corresponding word embeddings, word similarity can be obtained by computing their cosine similarity between \mathbf{v}_{w_i} and \mathbf{v}_{w_j} as

$$\cos(\mathbf{v}_{w_i}, \mathbf{v}_{w_j}) = \frac{\mathbf{v}_{w_i}^\top \mathbf{v}_{w_j}}{\|\mathbf{v}_{w_i}\|_2 \cdot \|\mathbf{v}_{w_j}\|_2} \quad (20)$$

Examples

Word ₁	Word ₂	Similarity score [0,10]
love	sex	6.77
stock	jaguar	0.92
money	cash	9.15
development	issue	3.97
lad	brother	4.46

Figure: Sample word pairs along with their human similarity judgment from WS-353 [Faruqui et al., 2016].

Available word similarity datasets

Dataset	Word pairs	Reference
RG	65	Rubenstein and Goodenough (1965)
MC	30	Miller and Charles (1991)
WS-353	353	Finkelstein et al. (2002)
YP-130	130	Yang and Powers (2006)
MTurk-287	287	Radinsky et al. (2011)
MTurk-771	771	Halawi et al. (2012)
MEN	3000	Bruni et al. (2012)
RW	2034	Luong et al. (2013)
Verb	144	Baker et al. (2014)
SimLex	999	Hill et al. (2014)

Figure: Word similarity datasets [Faruqui et al., 2016].

the **basis** for other intrinsic evaluations

Word Analogy

- ▶ It is sometimes referred as *linguistic regularity* [Mikolov et al., 2013b]
- ▶ The basic setup

$$w_a : w_b = w_c : ?$$

where $w_{a,b,c}$ are words and w_a, w_b are related under a certain linguistic relation

- ▶ Calculation: $(\mathbf{v}_{w_a} - \mathbf{v}_{w_b})^\top (\mathbf{v}_{w_c} - \mathbf{v}_{w_d})$

Word Analogy

- ▶ It is sometimes referred as *linguistic regularity* [Mikolov et al., 2013b]
- ▶ The basic setup

$$w_a : w_b = w_c : ?$$

where $w_{a,b,c}$ are words and w_a, w_b are related under a certain linguistic relation

- ▶ Calculation: $(\mathbf{v}_{w_a} - \mathbf{v}_{w_b})^\top (\mathbf{v}_{w_c} - \mathbf{v}_{w_d})$
- ▶ Example
 - ▶ Semantic love : like
 - ▶ Syntactic quick : quickly
 - ▶ Gender king : man
 - ▶ Others Beijing : China

Word Analogy: Examples

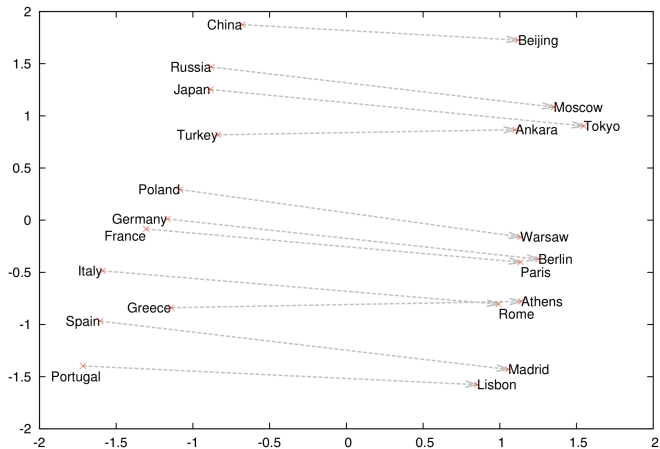


Figure: Word analogy examples.

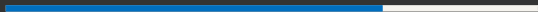
From [Faruqui et al., 2014]

naval, industrial, technological, marine, **identity**

- ▶ constructed from word embeddings
- ▶ evaluated by human annotators

- ▶ Implicit assumption: there is a consistent, global ranking of word embedding quality, and that higher quality embeddings will necessarily improve results on *any* downstream task.
- ▶ Unfortunately, this assumption does not hold in general [Schnabel et al., 2015].
- ▶ Examples
 - ▶ empirical results show that it may not be able give much help to syntactic parsing [Andreas and Klein, 2014]
 - ▶ adding surface-form features always help ([Ji and Eisenstein, 2014a] and many other works)

Problems



$$\boldsymbol{v}_{\text{man}} - \boldsymbol{v}_{\text{woman}} \approx \boldsymbol{v}_{\text{computer programmer}} - \boldsymbol{v}_{\text{homemaker}} \quad (21)$$

$$\boldsymbol{v}_{\text{father}} - \boldsymbol{v}_{\text{mother}} \approx \boldsymbol{v}_{\text{doctor}} - \boldsymbol{v}_{\text{nurse}} \quad (22)$$

[Bolukbasi et al., 2016]

Word embeddings like this not only reflect such stereotypes but also
amplify them

Three steps [Bolukbasi et al., 2016]

1. find gender neutral words with biases in the original embeddings;
2. identify the gender-specific space V and its orthogonal complement V^\perp
3. project embeddings of the gender neutral words to the subspace V^\perp

Example



Can we have an interpretability of each dimension?

Solution: post-processing on word embeddings

- ▶ reconstructing with sparsity constraint [Faruqui et al., 2015]
- ▶ rotating word embedding space using factor analysis [Park et al., 2017]

Interpretability is *derived* from the sparsity constraint as

$$\operatorname{argmin}_{\mathbf{D}, \mathbf{A}} \sum_{i=1}^V \|\mathbf{x}_i - \mathbf{D}\mathbf{a}_i\|_2^2 + \lambda \|\mathbf{a}_i\|_1 + \tau \|\mathbf{D}\|_2^2 \quad (23)$$

where \mathbf{x}_i and \mathbf{a}_i are the original and sparse embeddings of word i , \mathbf{D} is the transformation matrix.

Example

X	combat, guard, honor, bow, trim, naval 'll, could, faced, lacking, seriously, scored see, n't, recommended, depending, part due, positive, equal, focus, respect, better sergeant, comments, critics, she, videos
A	fracture, breathing, wound, tissue, relief relationships, connections, identity, relations files, bills, titles, collections, poems, songs naval, industrial, technological, marine stadium, belt, championship, toll, ride, coach

Figure: Top-ranked words per-dimension before and after reconstruction.
Each line shows words from a different dimension.

Problem

- ▶ Word embeddings from either Word2vec or GloVe encode not just semantic information

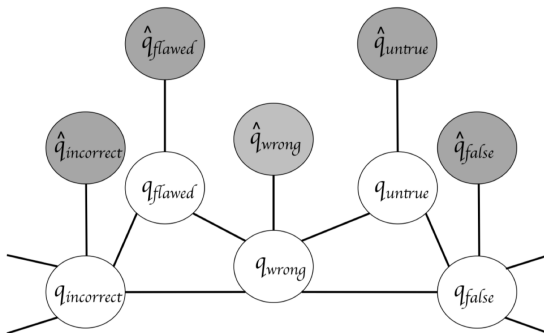
- ▶ Word embeddings from either Word2vec or GloVe encode not just semantic information
- ▶ In some applications, we want to emphasize one particular aspect of linguistic information
 - ▶ Semantic information [Faruqui et al., 2014]
 - ▶ Discourse information [Ji and Eisenstein, 2014b]

- ▶ Word embeddings from either Word2vec or GloVe encode not just semantic information
- ▶ In some applications, we want to emphasize one particular aspect of linguistic information
 - ▶ Semantic information [Faruqui et al., 2014]
 - ▶ Discourse information [Ji and Eisenstein, 2014b]
- ▶ Solutions
 - ▶ retrofitting word embeddings [Faruqui et al., 2014]
 - ▶ learning from supervision information [Ji and Eisenstein, 2014b]

Retrofitting

Retrofitting with WordNet [Miller, 1995]

- ▶ $\Omega = (V, E)$ be a semantic graph over words, where V is the node set with each element as a word, and E is the edge set with each edge representing a semantic relation between two words.



- ▶ The goal is to learn word embeddings $\{\tilde{v}\}$ such that \tilde{v}_i and \tilde{v}_j are close enough if $(i, j) \in E$.
- ▶ In addition, $\{\tilde{v}\}$ should also satisfy the constraint from original word embeddings, such that \tilde{v}_i and \tilde{v}_i are close enough for every word in \mathcal{V} .

$$\Psi(\tilde{\mathbf{V}}) = \sum_{i=1}^{|\mathcal{V}|} \left[\alpha_i \|\mathbf{v}_i - \tilde{\mathbf{v}}_i\|^2 + \sum_{(i,j) \in E} \beta_{ij} \|\tilde{\mathbf{v}}_i - \tilde{\mathbf{v}}_j\|^2 \right] \quad (24)$$

Learning from Supervision Signal

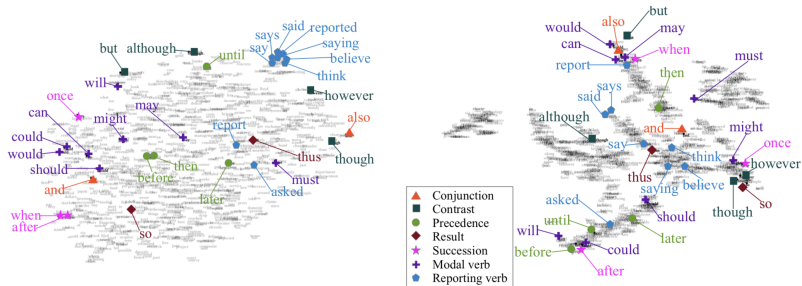


Figure: (Left) Word embeddings learned with supervision signal; (Right) Unsupervised word embeddings.



Andreas, J. and Klein, D. (2014).

How much do word embeddings encode about syntax?

In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 822–827.



Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016).

Man is to computer programmer as woman is to homemaker? debiasing word embeddings.

In *Advances in Neural Information Processing Systems*, pages 4349–4357.



Faruqui, M., Dodge, J., Jauhar, S. K., Dyer, C., Hovy, E., and Smith, N. A. (2014).

Retrofitting word vectors to semantic lexicons.

arXiv preprint arXiv:1411.4166.



Faruqui, M., Tsvetkov, Y., Rastogi, P., and Dyer, C. (2016).

Problems with evaluation of word embeddings using word similarity tasks.

arXiv preprint arXiv:1605.02276.



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

Sparse overcomplete word vector representations.

arXiv preprint arXiv:1506.02004.



Ji, Y. and Eisenstein, J. (2014a).

One vector is not enough: Entity-augmented distributional semantics for discourse relations.

arXiv preprint arXiv:1411.6699.



Ji, Y. and Eisenstein, J. (2014b).

Representation learning for text-level discourse parsing.

In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 13–24.



Levy, O. and Goldberg, Y. (2014).

Neural word embedding as implicit matrix factorization.

In *Advances in neural information processing systems*, pages 2177–2185.

