

# Deep Learning Techniques

---

1

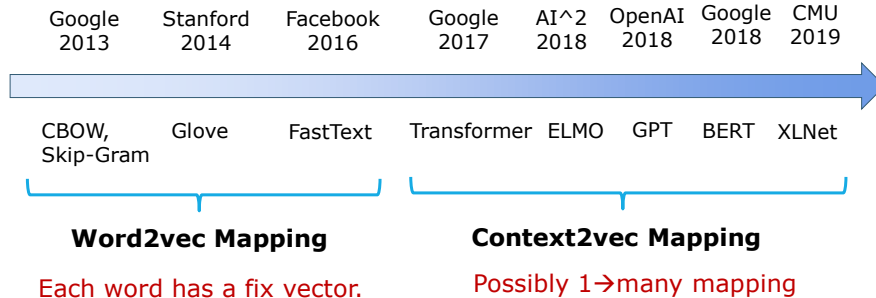
## Lectures

---

- DL1: Word Embedding
- DL2. Recurrent Neural Networks (RNN)
- DL3. Convolutional Neural Networks (CNN)
- DL4. Neural Attention Models (Transformer, etc.)
- DL5. Large Language Models (BERT, GPT, BART, etc.)

2

## Trend in Language Representation Learning



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

3

3

## Motivation

- J.R. Firth's hypothesis (1957): *"You shall know a word by the company it keeps."*
- **Word2vec mapping** (2013 – 2016)
  - Finding a k-dimensional vector (e.g., k=300) for each word based on the words co-occurring with it in many documents.
  - Widely used in text mining and NLP tasks (machine translation, question answering, text classification, information retrieval, summarization, etc.)

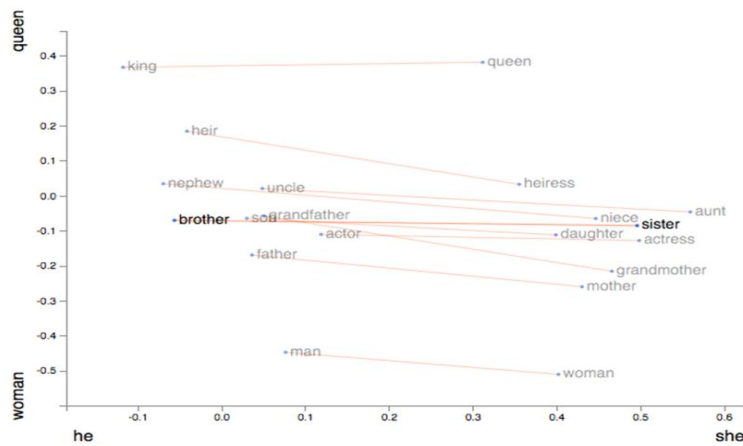
1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

4

4

## Masculine-Feminine Vectors by GloVe (Projected to 2D via PCA)



- Each word is a vector (point).
- Each relation is a line.
- Parallel lines show some "analogy" among word pairs.

<http://p.migdal.pl/2017/01/06/king-man-woman-queen-why.html>

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

5

5

## Country and Capital Vectors Projected by PCA

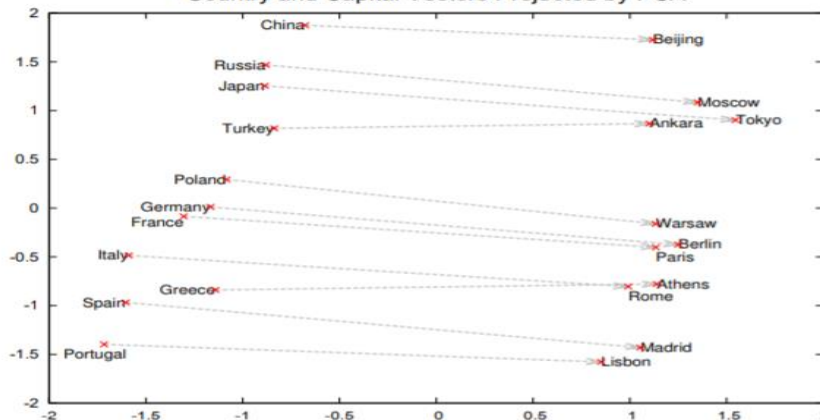


Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities. The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

6

6

## The Word Analogy Task

- Given a word pair and a word in another pair, find the word for "?"

$$\mathcal{R}(\text{man}, \text{woman}) \approx \mathcal{R}(\text{king}, ?)$$

- Given word embeddings, system finds the answer as

$$w_{?}^{*} = \operatorname{argmin}_{w_{?}} \|(w_{\text{man}} - w_{\text{woman}}) + (w_{\text{king}} - w_{?})\|$$

$$w_{\text{man}} - w_{\text{woman}} \approx w_{\text{king}} - w_{?}$$

$$w_{?} \approx w_{\text{king}} - w_{\text{man}} + w_{\text{woman}}$$

1/18/2024

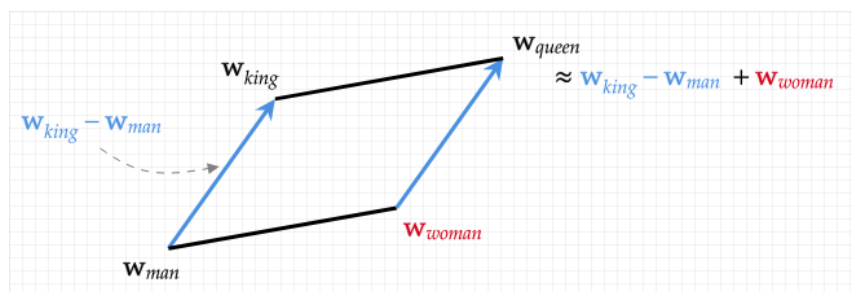
@Yiming Yang, lecture on Deep Learning for Text Mining

7

7

## Word Analogy: A Geometrical View

- We can draw a parallelogram as



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

8

8

## Neural Word Embedding Methods

- CBOW and Skip-Gram
- GloVe (Global Vectors for Word Representation)

## Training Data Generation

- Consider a sliding window of words over a sequence as  
 "we have  $w_{i-2}$   $w_{i-1}$   $w_i$   $w_{i+1}$   $w_{i+2}$  in NSH 1305 ..."
- The word in the middle is called the **target word** ( $w_i$ ), and the surrounding words together are called the **context**.
- Apply the sliding window over a large corpus of text we obtain many  $\langle \text{word}, \text{context} \rangle$  pairs as the training set.

## CBOW vs. Skip-Gram

- **CBOW** (predicting each target word given its context)

$$\max_{\theta} \sum_i \log P_{\theta}(w_i | \underbrace{w_{<i}, w_{>i}}_{\text{context } c_i}) \quad c_i = \{w_j : j \in i \pm k\}$$

- **Skip-Gram** (predicting the context given a target word)

$$\max_{\theta} \sum_i \log P_{\theta}(c_i | w_i) = \sum_i \sum_{w_j \in c_i} \log P_{\theta}(w_j | w_i)$$

- Both methods ignore the word order in the context window.

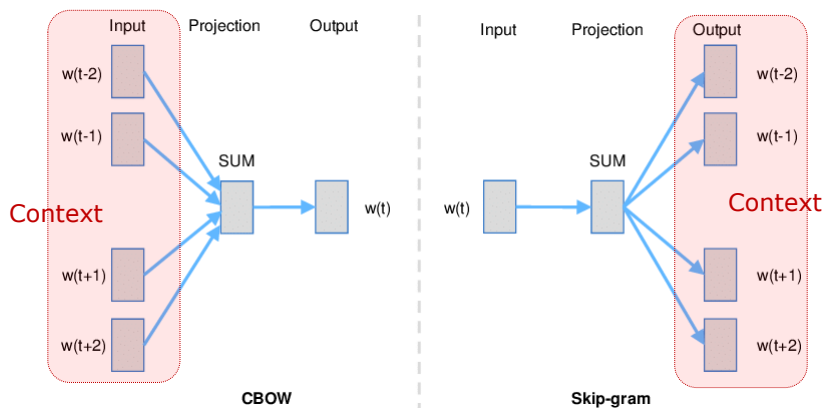
1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

11

11

## CBOW and Skip-Gram for word2vec



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

12

12

## CBOW Architecture

- Input:** Each context word is represented as **one-hot vector** (all the elements are 0 except one), whose dimension is the vocabulary size  $V$ .
- Hidden Layer:**  $h \in \mathbb{R}^N$  (with  $N \ll V$ ) is the context embedding calculated as

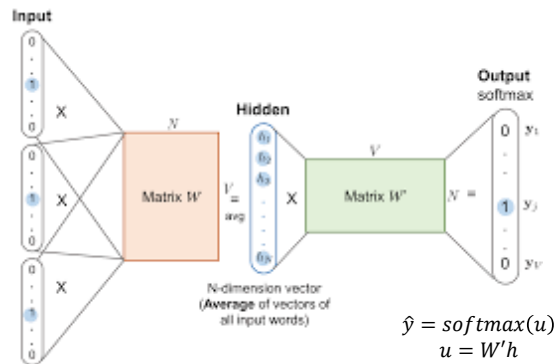
$$h = \frac{1}{|c|} \sum_{j \in c} W^T x_j$$

where  $W \in \mathbb{R}^{V \times D}$  is a matrix of learnable parameters.

- Output Layer:** the predicted probabilistic distribution of candidate words in a vector

$$\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_V) = \text{softmax}(W'h)$$

where  $W' \in \mathbb{R}^{N \times V}$  is another matrix of learnable parameters.



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

13

13

## CBOL with a Compact Input

- Input Layer (with a merged vector)

$$x = \frac{1}{|c|} (x_{w_1} + x_{w_2} + \dots + x_{w_{|c|}})$$

- Hidden layer

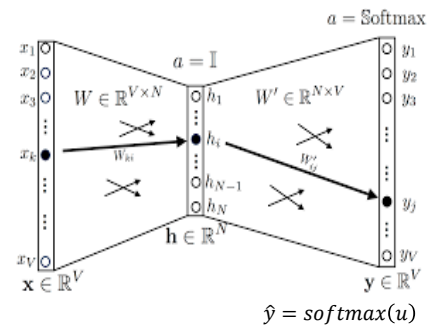
$$h = W^T x$$

$$W = \frac{1}{|c|} (W^T x_{w_1} + W^T x_{w_2} + \dots + W^T x_{w_{|c|}})$$

- Model Parameters

$$\theta = (W, W')$$

- Matrix  $W \in \mathbb{R}^{V \times N}$  contains all the word embeddings (each row is the embedding of a word).**



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

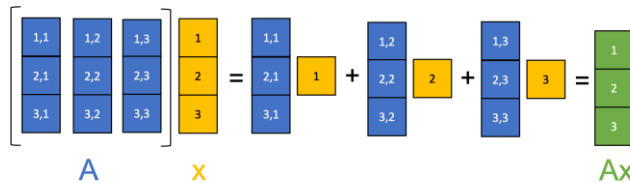
14

14

## Word-embedding matrix $W$

- Hidden layer

$$\mathbf{h} = W^T \mathbf{x} = \frac{1}{|c|} W^T (\mathbf{x}_{w_1} + \mathbf{x}_{w_2} + \dots + \mathbf{x}_{w_{|c|}})$$



In CBOW,  $A = W^T$  and  $x$  is the sum of a few one-hot vectors.

$Ax$  picks a few columns of  $A$  to sum up (and average over) for **context embedding**.  
Each column of  $W^T$  (each row of  $W$ ) is the **embedding of a word**.

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

15

15

## Training Process

- Denote Set:  $\mathcal{D} = \{(x_i, y_i)\}$  context-word pairs.
- Model Initialization: Randomly initialize  $W$  and  $W'$
- Model Update: For each pair  $(x, y) \in \mathcal{D}$ 
  - Forward Propagation**: Fix  $W$  and  $W'$ , compute  $h$ ,  $u$  and  $\hat{y}$  given  $x$ ;
  - Backpropagation**: Fix  $h$ ,  $u$  and  $\hat{y}$ , update  $W$  and  $W'$  with mini-batch gradients.

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

16

16



## Model Parameter Optimization

- Iterative update

$$W_{kl}^{(new)} := W_{kl}^{(old)} - \eta \nabla_{W_{kl}} \left( \frac{1}{|B|} \sum_{y_i \in B} l_{\theta}(\hat{y}_i, y_i) \right)$$

$$W'_{ij}{}^{(new)} := W'_{ij}{}^{(old)} - \eta \nabla_{W'_{ij}} \left( \frac{1}{|B|} \sum_{y_i \in B} l_{\theta}(\hat{y}_i, y_i) \right)$$

where  $l_{\theta}(\hat{y}_i, y_i)$  is the loss on each training pair;

$B$  is a randomly sampled mini-batch from the training set.

07/23/2017

@Yiming Yang, lecture on Deep Learning for Text Mining

17

17

## Loss Function (on a single training pair for simplicity)

- Cross entropy loss

$$l(\hat{y}, y) = -\sum_{j=1}^V y_j \log \hat{y}_j = -\log \hat{y}_{j^*}$$

where  $j^*$  is the index of the target word in training pair  $(x_i, y_i)$ .

- Example

$$y = (0 \quad \color{red}{1} \quad 0 \quad 0 \quad 0) \quad \color{red}{j^*=2}$$

$$\hat{y} = (0.1 \quad 0.5 \quad 0.2 \quad 0.1 \quad 0.1)$$

$$L(\hat{y}, y) = -\log \hat{y}_2 = -\log 0.5$$

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

18

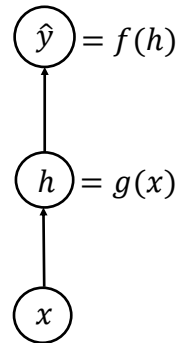
18

[L.P. Morency: CMU 11-777]

# Gradient Computation

## □ Chain rule

$$\frac{df}{dx} = \frac{df}{dh} \frac{dh}{dx}$$



07/23/2017

@Yiming Yang, lecture on Deep Learning for Text Mining

19

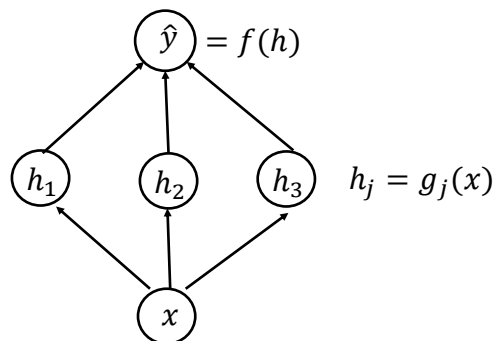
19

[L.P. Morency: CMU 11-777]

# Gradient Computation (cont'd)

## □ Chain rule

$$\frac{df}{dx} = \sum_j \frac{\partial f}{\partial h_j} \frac{dh_j}{dx}$$



07/23/2017

@Yiming Yang, lecture on Deep Learning for Text Mining

20

20

[L.P. Morency: CMU 11-777]

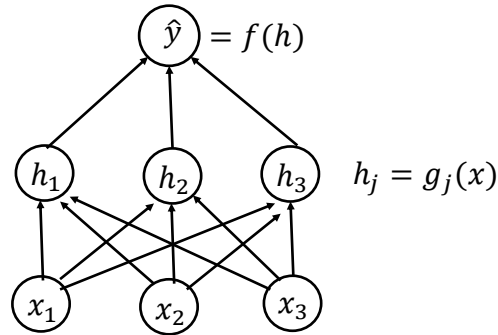
## Gradient Computation (cont'd)

### □ Chain rule

$$\frac{\partial f}{\partial x_1} = \sum_j \frac{\partial f}{\partial h_j} \frac{\partial h_j}{\partial x_1}$$

$$\frac{\partial f}{\partial x_2} = \sum_j \frac{\partial f}{\partial h_j} \frac{\partial h_j}{\partial x_2}$$

$$\frac{\partial f}{\partial x_3} = \sum_j \frac{\partial f}{\partial h_j} \frac{\partial h_j}{\partial x_3}$$



07/23/2017

@Yiming Yang, lecture on Deep Learning for Text Mining

21

21

## Gradient Computation (cont'd)

### □ the gradient (scalar-by-vector)

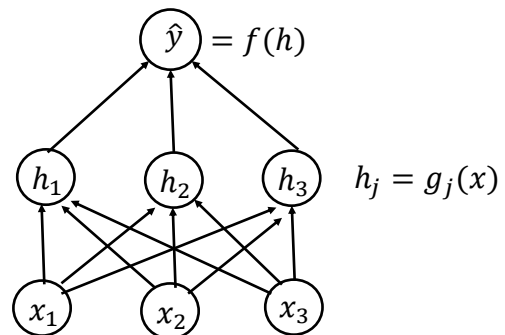
$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_3} \end{bmatrix}$$

$$= \nabla_h f \left( \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right)$$

Gradient of  $f$  w.r.t.  $\mathbf{h}$ 

$$\begin{bmatrix} \frac{\partial f}{\partial h_1} & \frac{\partial f}{\partial h_2} & \frac{\partial f}{\partial h_3} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \frac{\partial h_1}{\partial x_3} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \frac{\partial h_2}{\partial x_3} \\ \frac{\partial h_3}{\partial x_1} & \frac{\partial h_3}{\partial x_2} & \frac{\partial h_3}{\partial x_3} \end{bmatrix}$$

Jacobian matrix of size  $|\mathbf{h}| \times |\mathbf{x}|$ 

07/23/2017

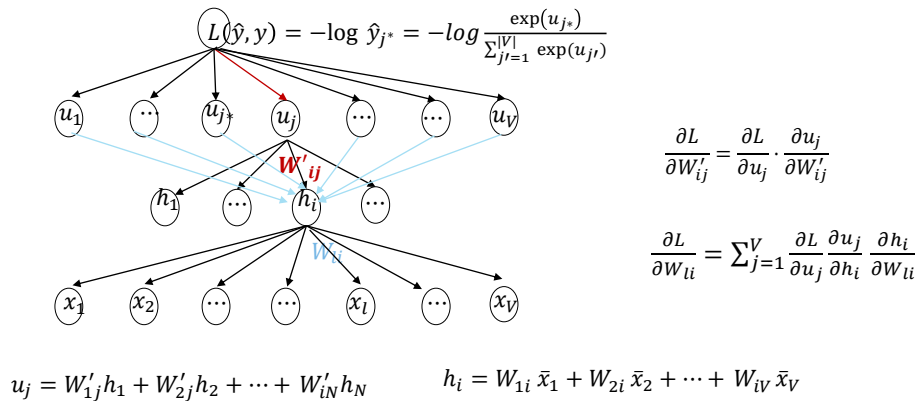
@Yiming Yang, lecture on Deep Learning for Text Mining

22

@Yiming Yang, lecture on Deep Learning for Text Mining

22

## Packpropagation in CBOL



1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

23

23

## Backpropagation [Xin Rong, arXiv 2016]

- Partial derivative with respect to each element of  $W'$

$$\frac{\partial L}{\partial W'_{ij}} = \frac{\partial L}{\partial u_j} \cdot \frac{\partial u_j}{\partial W'_{ij}}$$

$$\begin{aligned} \frac{\partial L}{\partial u_j} &= \frac{\partial}{\partial u_j} \left( -\log \frac{\exp(u_{j^*})}{\sum_{j'=1}^V \exp(u_{j'})} \right) = \frac{\partial}{\partial u_j} (-u_{j^*}) + \frac{\partial}{\partial u_j} (\log(\sum_{j'=1}^V \exp(u_{j'}))) \\ &= -\delta_j + \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})} = -\delta_j + \hat{y}_j \end{aligned} \quad (1)$$

( $\delta_j$  is the Kronecker delta function,  $\delta_j = 1$  if  $j=j^*$ ; otherwise  $\delta_j = 0$ )

$$\frac{\partial u_j}{\partial W'_{ij}} = \frac{\partial}{\partial W'_{ij}} (W'_{1j} h_1 + W'_{2j} h_2 + \dots + W'_{ij} h_i + \dots) = h_i \quad (2)$$

- Combine (1) and (2), we have

$$\frac{\partial L}{\partial W'_{ij}} = (\hat{y}_j - \delta_j) h_i \quad (3)$$

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

24

24

## Backpropagation (cont'd)

- Partial derivative w.r.t. the elements of  $W$

$$\frac{\partial L}{\partial W_{li}} = \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial W_{li}} = \sum_{j=1}^V \frac{\partial L}{\partial u_j} \frac{\partial u_j}{\partial h_i} \frac{\partial h_i}{\partial W_{li}}$$

$$\frac{\partial L}{\partial u_j} = (\hat{y}_j - \delta_j) \quad (1)$$

$$\frac{\partial u_j}{\partial h_i} = \frac{\partial}{\partial h_i} (W'_{1j}h_1 + W'_{2j}h_2 + \dots + W'_{ij}h_i + \dots) = W'_{ij} \quad (2)$$

$$\frac{\partial h_i}{\partial W_{li}} = \frac{\partial}{\partial W_{li}} (W_{1i}x_1 + W_{2i}x_2 + \dots + W_{li}x_l + \dots) = x_l \quad (3)$$

- Finally,  $\frac{\partial L}{\partial W_{li}} = \sum_{j=1}^V (\hat{y}_j - \delta_j) \cdot W'_{ij}x_l$

## Neural Word Embedding Methods

- ✓ CBOW and Skip-Gram
- GloVe (Global Vectors for Word Representation)

## GloVe [Jeffrey Pennington et. al., EMNLP 2014]

- Input matrix  $X \in R^{V \times V}$ 
  - $X_{ij}$  is the **global sum** of the weight of word  $j$  in **local context** window of  $\pm 10$  from word  $i$ ;
  - if word  $j$  is  $k$  ( $\leq 10$ ) words apart from word  $i$ , then it has the weight of  $\frac{1}{k}$ .
  - Taking a log scale of each element as  $X_{ij} \xrightarrow{\text{logscale}} \log X_{ij}$
- Matrix  $W \in R^{V \times D}$  (word embeddings) and bias vector  $b \in R^V$ 
  - Matrix  $W$  and vector  $b = (b_1, b_2, \dots, b_V)$  are the model parameters to be optimized given  $X$ .
  - Each row of  $W$  is the embedding of a word, denoted as  $w_i$ .

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

27

27

## GloVe [Jeffrey Pennington et. al., EMNLP 2014]

### Objective

$$\min_{W, b} \sum_{i,j=1}^m f(X_{ij}) (\underbrace{w_i^T w_j + b_i + b_j}_{\hat{X}_{ij}} - X_{ij})^2$$

where  $f(x) = \begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$

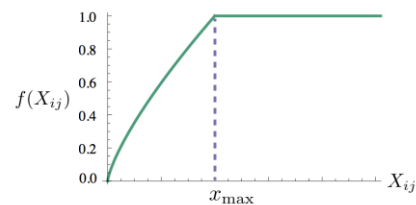


Figure 1: Weighting function  $f$  with  $\alpha = 3/4$ .

### Algorithm

- stochastic gradient descent

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

28

28

## Differences of GloVe from CBOW

### □ GloVe have the following differences

- **Weight** a local contextual word by  $1/K$  (sensitive to word position)
- Giving more weights via  $f(X_{ij})$  to larger cells (common word pairs) in  $X$
- Adding a bias term per word in addition to its embedding vector
- Despite the name, GloVe (**G**lobal **V**ectors for word representation) is still based on the **local** context around each target word
- In fact, all the MF methods (SVD, NMF, PMF) can also be applied to the GloVe's input matrix  $X$  for producing the word-embedding matrix  $W$ .

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

29

29

## Evaluation Results

- Word analogy task [Jeffrey Pennington et. al., EMNLP 2014]
- Other downstream tasks (not included here)
  - Name Entity Recognition (NER)
  - Sentimental Classification
  - Language Modeling
  - Neural Machine Translation

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

30

30

# The Word Analogy Task

- Input question likes "a is to b as c is to \_\_\_?"
- System's answer by finding  $w_d$  that is closest to  $w_b - w_a + w_c$  according to cosine similarity
- Dataset contains 19,544 questions, divided into two subsets
  - **Semantic**: "Athens is to Greece as Berlin is to \_\_\_?"
  - **Syntactic**: "dance is to dancing as fly is to \_\_\_?"

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

31

31

## Word Analogy Results [Jeffrey Pennington et. al., EMNLP 2014]

- **Baseline Methods**
  - HPCA: PMI version of LSA (PCA) [10]
  - vLBL, ivLBL: log-bilinear model [9]
  - SG: skip gram (another variant of w2v)
  - CBOW: continuous bag-of-words
  - SVD-S: take SVD of  $\sqrt{X_{trunc}}$
  - SVD-L: take SVD of  $\log(1 + X_{trunc})$
- **Metric: Accuracy**
- **Size: number of tokens in training set**

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW <sup>†</sup>	300	6B	63.6	<u>67.4</u>	65.7
SG <sup>†</sup>	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	<u>67.0</u>	<u>71.7</u>
CBOW	1000	6B	<u>57.3</u>	<u>68.9</u>	<u>63.7</u>
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

32

32



## Summary of word2vec methods

- CBOW (and SkipGram) treats local context as **a set of words** (ignoring word order) and learns the word embeddings with a one-hidden-layer neural network.
- GloVe use a matrix to aggregates **the global co-occurrence counts** (weighted by proximity and learns word embeddings via gradient descent.
- Both methods produces **a fixed embedding** for each word that cannot differentiate word meanings under different contexts (more discussions in the DL5 Lecture).

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

33

33

## References

1. Distributed representations of words and phrases and their compositionality. Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, Jeffrey Dean, NIPS 2013
2. Efficient Estimation of Word Representations in Vector Space. Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean ICLR Workshop 2013
3. Xin Rong: word2vec Parameter Learning Explained. [CoRR abs/1411.2738](https://arxiv.org/abs/1411.2738) (2016)
4. GloVe: Global Vectors for Word Representation. Jeffrey Pennington, Richard Socher, Christopher D. Manning. EMNLP 2014

1/18/2024

@Yiming Yang, lecture on Deep Learning for Text Mining

34

34