

# Introduction to learning of text representations

Prakhar Gupta

Department of Computer and Communication Sciences  
EPFL

Advanced Topics in Machine Learning

## 1 Learning of word representations

- GloVe
- Word2Vec

## 2 Second Main Section

- Another Subsection

## 1 Learning of word representations

- GloVe
- Word2Vec

## 2 Second Main Section

- Another Subsection

# Matrix Factorization Methods

- Models that can be described as optimization problems of the form

$$\min_{U,V} F(UV^T) \quad (1)$$

# GloVe: Global Vectors for Word Representation

## Some important definitions and notation

# GloVe: Global Vectors for Word Representation

## Some important definitions and notation

- Context window is defined as the neighbourhood of a word and its length can be chosen as desired.

## Some important definitions and notation

- Context window is defined as the neighbourhood of a word and its length can be chosen as desired.
- Word co-occurrence matrix is denoted by  $X$



# GloVe: Global Vectors for Word Representation

## Some important definitions and notation

- Context window is defined as the neighbourhood of a word and its length can be chosen as desired.
- Word co-occurrence matrix is denoted by  $X$ 
  - where  $X_{ij}$  = Number of times word  $j$  appears in context of word  $i$

# GloVe: Global Vectors for Word Representation

# GloVe: Global Vectors for Word Representation

- The aim is to capture  $X_{ij}$  using the source embeddings  $u_i$  and target embeddings  $v_j$

# GloVe: Global Vectors for Word Representation

- The aim is to capture  $X_{ij}$  using the source embeddings  $u_i$  and target embeddings  $v_j$
- The GloVe problem is thus formulated as

$$\min_{U,V} \sum_{i,j \in W} f(X_{ij})(u_i^T v_j + b_i + c_j + \log(X_{ij})) \quad (2)$$

where  $f(X_{ij})$  is the weight assigned to the source-target pair,  $b_i$  and  $c_j$  are the biases associated with  $u_i$  and  $v_j$  respectively and  $W$  is the vocabulary.

# GloVe: Global Vectors for Word Representation

- The aim is to capture  $X_{ij}$  using the source embeddings  $u_i$  and target embeddings  $v_j$
- The GloVe problem is thus formulated as

$$\min_{U,V} \sum_{i,j \in W} f(X_{ij})(u_i^T v_j + b_i + c_j + \log(X_{ij})) \quad (2)$$

where  $f(X_{ij})$  is the weight assigned to the source-target pair,  $b_i$  and  $c_j$  are the biases associated with  $u_i$  and  $v_j$  respectively and  $W$  is the vocabulary.

- $f(X_{ij})$  is often chosen to be  $(\frac{X_{ij}}{Y})^\alpha$  where  $Y = \max_{kl} X_{kl}$ .

# GloVe: Global Vectors for Word Representation

- The aim is to capture  $X_{ij}$  using the source embeddings  $u_i$  and target embeddings  $v_j$
- The GloVe problem is thus formulated as

$$\min_{U,V} \sum_{i,j \in W} f(X_{ij})(u_i^T v_j + b_i + c_j + \log(X_{ij})) \quad (2)$$

where  $f(X_{ij})$  is the weight assigned to the source-target pair,  $b_i$  and  $c_j$  are the biases associated with  $u_i$  and  $v_j$  respectively and  $W$  is the vocabulary.

- $f(X_{ij})$  is often chosen to be  $(\frac{X_{ij}}{Y})^\alpha$  where  $Y = \max_{kl} X_{kl}$ .
- Empirically  $\alpha = \frac{3}{4}$  gives the best performance.

## 1 Learning of word representations

- GloVe
- Word2Vec

## 2 Second Main Section

- Another Subsection

# Word2Vec: CBOW and Skipgram models



# Word2Vec: CBOW and Skipgram models

- Uses two different architectures

# Word2Vec: CBOW and Skipgram models

- Uses two different architectures
  - 1 Continuous bag-of-words (CBOW)

# Word2Vec: CBOW and Skipgram models

- Uses two different architectures
  - 1 Continuous bag-of-words (CBOW)
  - 2 Continuous Skipgram

# Word2Vec:CBOW

## Intuition

For the CBOW architecture, the task is to

## Intuition

For the CBOW architecture, the task is to

- predict the word  $w$  given the context  $\mathcal{C}(w)$

## Intuition

For the CBOW architecture, the task is to

- predict the word  $w$  given the context  $\mathcal{C}(w)$

## Formulation

Given a sequence of training words  $w_1, \dots, w_n$ , the maximum likelihood formulation for the CBOW architecture can be written as

$$\sum_{i=1}^n \sum_{w_j \in \mathcal{C}(w_i)} \log p(w_j | w_i) \quad (3)$$

# Word2Vec:Skipgram



## Intuition

For the Skipgram architecture, the task is to

## Intuition

For the Skipgram architecture, the task is to

- predict the context  $\mathcal{C}(w)$  given the word  $w$

## Intuition

For the Skipgram architecture, the task is to

- predict the context  $\mathcal{C}(w)$  given the word  $w$

## Formulation

Given a sequence of training words  $w_1, \dots, w_n$ , the maximum likelihood formulation for the CBOW architecture can be written as

$$\sum_{i=1}^n \sum_{w_j \in \mathcal{C}(w_i)} \log p(w_i | w_j) \quad (4)$$

## 1 Learning of word representations

- GloVe
- Word2Vec

## 2 Second Main Section

- Another Subsection

## Block Title

You can also highlight sections of your presentation in a block, with it's own title

## Theorem

*There are separate environments for theorems, examples, definitions and proofs.*



## Example

Here is an example of an example block.

# Summary

- The **first main message** of your talk in one or two lines.
- The **second main message** of your talk in one or two lines.
- Perhaps a **third message**, but not more than that.
- Outlook
  - Something you haven't solved.
  - Something else you haven't solved.

# For Further Reading I

-  A. Author.  
*Handbook of Everything*.  
Some Press, 1990.
-  S. Someone.  
On this and that.  
*Journal of This and That*, 2(1):50–100, 2000.