# Deep Learning

13 Word Embeddings

Dr. Konda Reddy Mopuri
Dept. of AI, IIT Hyderabad
Jan-May 2024

ML model

IITH has been consistently ranked in the top 10 institutes in India for Engineering according to NIRF making it one of the most coveted schools for science and technology in the country.

# Why Word Embeddings?



ML model

2.3, -1.2, 0.3, 5.2, ................, 0.7, -2.4, 8.2, 9.5

IITH has been consistently ranked in the top 10 institutes in India for Engineering according to NIRF making it one of the most coveted schools for science and technology in the country.

భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

- Corpus: a collection of authentic text organized into dataset

# Terminology

- Corpus: a collection of authentic text organized into dataset
- Vocabulary (V): Set of allowed words

# Terminology

- Corpus: a collection of authentic text organized into dataset
- Vocabulary (V): Set of allowed words
- Target: Representation for every word in V

# One-hot Encoding

- Representation using discrete symbols

# One-hot Encoding

- Representation using discrete symbols
- $|V|$ words encoded as binary vectors of length $|V|$

Dictionary | Word Representation

| A | | 1 | 0 | 0 | ·········· | 0 | 0 |

| Bus | | 0 | 1 | 0 | ········· | 0 | 0 |

| Cat | | 0 | 0 | 1 | ·········· | 0 | 0 |

⋮

| Tide | | 0 | 0 | 0 | ········· | 1 | 0 |

| Zone | | 0 | 0 | 0 | ········· | 0 | 1 |

# One-hot encoding: Drawbacks

1. Space inefficient (e.g. 13M words in Google 1T corpus)

# One-hot encoding: Drawbacks

1. Space inefficient (e.g. 13M words in Google 1T corpus)
2. No notion of similarity (or, distance) between words

# One-hot encoding: Drawbacks

1. Space inefficient (e.g. 13M words in Google 1T corpus)
2. No notion of similarity (or, distance) between words
   - 'Dog' is as close to 'Cat' as it is to 'Machine'

# Notion of Meaning for words

1. What is a good notion of meaning for a word?

# Notion of Meaning for words

1. What is a good notion of meaning for a word?
2. How do we, humans, know the meaning of a word?

# Notion of Meaning for words

1. What does silla mean?

# Notion of Meaning for words

1. Let's see how this word is used in different contexts

    1. The **silla** is by the window, offering a nice view of the garden.

    2. Can you pass me that **silla** so I can join the conversation?

    3. After the event, please stack the **sillas** neatly against the wall.

    4. I found a comfortable **silla** in the corner and settled down to relax.

# Notion of Meaning for words

1. Does this context help you understand the word silla?

 

1. The silla is by the window, offering a nice view of the garden.

2. Can you pass me that silla so I can join the conversation?

3. After the event, please stack the sillas neatly against the wall.

4. I found a comfortable silla in the corner and settled down to relax.

# Notion of Meaning for words

1. Does this context help you understand the word silla?

2. { positioned near a window or against a wall or in the corner, used for conversing/events, can be used to relax }

    1. The `silla` is by the window, offering a nice view of the garden.

    2. Can you pass me that `silla` so I can join the conversation?

    3. After the event, please stack the `sillas` neatly against the wall.

    4. I found a comfortable `silla` in the corner and settled down to relax.

# Notion of Meaning for words

1. How did we do that?

# Notion of Meaning for words

① How did we do that?

② "We searched for other words that can be used in the same contexts, found some, and made a conclusion that <span style="color:red">silla</span> has to mean similar to those words."

# Notion of Meaning for words

1. Distributional Hypothesis: Words that frequently appear in similar contexts have a similar meaning

# Distributed Representations

1. Representation/meaning of a word should consider its context in the corpus

# Distributed Representations

1. Representation/meaning of a word should consider its context in the corpus
2. Use many contexts of a word to build up a representation for it

# Distributed Representations

1. **Co-occurrence matrix** is a way to can capture this!
   - size: ($\#$words $\times$ $\#$words)
   - rows: words (m), cols: context (n)
   - words and context can be of same or different size

# Distributed Representations

1. **Co-occurrence matrix** is a way to can capture this!
   - size: ($\#$words $\times$ $\#$words)
   - rows: words (m), cols: context (n)
   - words and context can be of same or different size

2. Context can be defined as a 'h' word neighborhood

# Distributed Representations

1. **Co-occurrence matrix** is a way to can capture this!
   - size: (#words $\times$ #words)
   - rows: words (m), cols: context (n)
   - words and context can be of same or different size
2. Context can be defined as a 'h' word neighborhood
3. Each row (column): vectorial representation of the word (context)

# Co-occurrence matrix

$$X = \begin{array}{c|cccccccc} & I & like & enjoy & deep & learning & NLP & flying & . \\ \hline I & 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ like & 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ enjoy & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ deep & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ learning & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ NLP & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ flying & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ . & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{array}$$

# Co-occurrence matrix

1. Very sparse

# Co-occurrence matrix

1. Very sparse
2. Very high-dimensional (grows with the vocabulary size)

# Co-occurrence matrix

1. Very sparse
2. Very high-dimensional (grows with the vocabulary size)
3. Solution:Dimensionality reduction (SVD)!

# SVD on the Co-occurrence matrix

1. $X = U\Sigma V^T$

# SVD on the Co-occurrence matrix

① $X = U\Sigma V^T$

② $\begin{bmatrix} & & \\ & X & \\ & & \end{bmatrix}_{m \times n} =$

$\begin{bmatrix} \uparrow & \cdots & \uparrow \\ u_1 & \cdots & u_k \\ \downarrow & \cdots & \downarrow \end{bmatrix}_{m \times k} \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix}_{k \times k} \cdot \begin{bmatrix} \leftarrow & v_1^T & \rightarrow \\ & \vdots & \\ \leftarrow & v_k^T & \rightarrow \end{bmatrix}_{k \times n}$

# SVD on the Co-occurrence matrix

1. $X = U\Sigma V^T$

2. $\begin{bmatrix} & & \\ & X & \\ & & \end{bmatrix}_{m \times n} =$

   $\begin{bmatrix} \uparrow & \cdots & \uparrow \\ u_1 & \cdots & u_k \\ \downarrow & \cdots & \downarrow \end{bmatrix}_{m \times k} \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix}_{k \times k} \cdot \begin{bmatrix} \leftarrow & v_1^T & \rightarrow \\ & \vdots & \\ \leftarrow & v_k^T & \rightarrow \end{bmatrix}_{k \times n}$

3. $X = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \ldots + \sigma_k u_k v_k^T$

4. $\hat{X} = \sum_{i=1}^{d < k} \sigma_i u_i v_i^T$ is a $d$-rank approximation of $X$

# SVD on the Co-occurrence matrix

1. Before the SVD, representations were the rows of $X$

# SVD on the Co-occurrence matrix

1. Before the SVD, representations were the rows of $X$
2. How do we reduce the representation size with SVD ?

# SVD on the Co-occurrence matrix

1. Before the SVD, representations were the rows of $X$
2. How do we reduce the representation size with SVD ?
3. $W_{\mathsf{word}} = U_{m \times k} \cdot \Sigma_{k \times k}$

# SVD on the Co-occurrence matrix

1. $W_{\text{word}} \in \mathbb{R}^{m \times k}$ ($k \ll |V| = m$) are considered the representation of the words

# SVD on the Co-occurrence matrix

1. $W_{\text{word}} \in \mathbb{R}^{m \times k}$ ($k \ll |V| = m$) are considered the representation of the words

2. Lesser dimensions but the same similarities! (one may verify that $XX^T = \hat{X}\hat{X}^T$)

# SVD on the Co-occurrence matrix

1. $W_{\text{word}} \in \mathbb{R}^{m \times k}$ ($k \ll |V| = m$) are considered the representation of the words

2. Lesser dimensions but the same similarities! (one may verify that $XX^T = \hat{X}\hat{X}^T$)

3. $W_{\text{context}} = V \in \mathbb{R}^{n \times k}$ are taken as the representations for the context words

# A bit more clever things...

1. Entries in the occurrence matrix can be weighted (HAL model)

# A bit more clever things...

1. Entries in the occurrence matrix can be weighted (HAL model)
2. Better associations can be used (PPMI)

# A bit more clever things...

1. Entries in the occurrence matrix can be weighted (HAL model)
2. Better associations can be used (PPMI)
3. ....

# Count-based vs prediction-based models

① Techniques we have seen so far rely on the counts (or, co-occurrence of words)

# Count-based vs prediction-based models

1. Techniques we have seen so far rely on the counts (or, co-occurrence of words)
2. Next, we see prediction based models for word embeddings

# Word2Vec

1. T Mikolov et al. (2013)

# Word2Vec

1. T Mikolov et al. (2013)
2. Two versions: Predict words from the contexts (or contexts from words)

# Word2Vec

1. T Mikolov et al. (2013)
2. Two versions: Predict words from the contexts (or contexts from words)
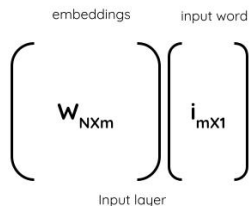3. Continuous Bag of Words (CBoW) and Skip-gram

# Word2Vec



CBOW                    Skip-gram

# Word Embeddings: Skip-gram

Skip-gram

# Word Embeddings: Skip-gram



embeddings      input word

$$W_{NXm} \quad i_{mX1}$$

Input layer

# Word Embeddings: Skip-gram

1. Start: huge corpus and random initialization of the word embeddings

# Word Embeddings: Skip-gram

1. Start: huge corpus and random initialization of the word embeddings
2. Process the text with a sliding window (one word at a time)

# Word Embeddings: Skip-gram

1. Start: huge corpus and random initialization of the word embeddings
2. Process the text with a sliding window (one word at a time)

1. At each step, there is a central word and context words (other words in the window)

# Word Embeddings: Skip-gram

1. Start: huge corpus and random initialization of the word embeddings
2. Process the text with a sliding window (one word at a time)

1. At each step, there is a central word and context words (other words in the window)
2. Given the central word, compute the probabilities for the context words

# Word Embeddings: Skip-gram

1. Start: huge corpus and random initialization of the word embeddings
2. Process the text with a sliding window (one word at a time)

 

1. At each step, there is a central word and context words (other words in the window)
2. Given the central word, compute the probabilities for the context words
3. Modify the word embeddings to increase these probabilities
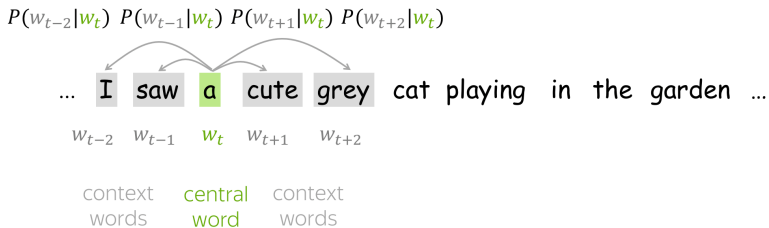
# Word Embeddings: Skip-gram

$P(w_{t-2}|w_t)$ $P(w_{t-1}|w_t)$ $P(w_{t+1}|w_t)$ $P(w_{t+2}|w_t)$

... I saw a cute grey cat playing in the garden ...

$w_{t-2}$ $w_{t-1}$ $w_t$ $w_{t+1}$ $w_{t+2}$

context words     central word     context words

Figure from Lena Voita

# Word Embeddings: Skip-gram



$P(w_{t-2}|w_t)$  $P(w_{t-1}|w_t)$  $P(w_{t+1}|w_t)$  $P(w_{t+2}|w_t)$

... I  saw  a  cute  grey  cat  playing  in  the  garden  ...

$w_{t-2}$  $w_{t-1}$  $w_t$  $w_{t+1}$  $w_{t+2}$

context words    central word    context words

Figure from Lena Voita

# Word Embeddings: Skip-gram

$$P(w_{t-2}|w_t) \quad P(w_{t-1}|w_t) \quad P(w_{t+1}|w_t) \quad P(w_{t+2}|w_t)$$

... I saw  a  cute  grey  cat  playing  in  the  garden ...

$$w_{t-2} \quad w_{t-1} \quad w_t \quad w_{t+1} \quad w_{t+2}$$

context words    central word    context words

Figure from Lena Voita

# Word Embeddings: Skip-gram

$$P(w_{t-2}|w_t) \quad P(w_{t-1}|w_t) \quad P(w_{t+1}|w_t) \quad P(w_{t+2}|w_t)$$

... I saw a cute grey cat playing in the garden ...

$$w_{t-2} \quad w_{t-1} \quad w_t \quad w_{t+1} \quad w_{t+2}$$

| context words | central word | context words |
|---|---|---|

Figure from Lena Voita

# Word Embeddings: Skip-gram

$$P(w_{t-2}|w_t)\ P(w_{t-1}|w_t)\ P(w_{t+1}|w_t)\ P(w_{t+2}|w_t)$$

... I saw a cute grey $\boxed{\text{cat}}$ $\boxed{\text{playing}}$ $\boxed{\text{in}}$ $\boxed{\text{the}}$ $\boxed{\text{garden}}$ ...

$\quad\quad\quad\quad\quad\quad\quad w_{t-2}\quad\quad w_{t-1}\quad\quad w_t\quad w_{t+1}\quad\quad w_{t+2}$

context words    central word    context words

Figure from Lena Voita

# Word Embeddings: Skip-gram

- For each position in $t = 1, 2, \ldots T$ in the corpus, Skip-gram predicts the context words in $m-$sized window ($\theta$ is the variables to be optimized)

$$\text{Likelihood} \quad L(\theta) = \prod_{t=1}^{T} \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t, \theta)$$

# Word Embeddings: Skip-gram

- The loss is mean NLL

$$\text{Loss} \quad J(\theta) = -\frac{1}{T} \prod_{t=1}^{T} \prod_{-m \le j \le m, j \ne 0} \log P(w_{t+j}|w_t, \theta)$$

# Word Embeddings: Skip-gram

- What are the parameters $(\theta)$ to be learned?



Figure from Lena Voita

# Word Embeddings: Skip-gram

- How to compute $P(w_{t+j}|w_t, \theta)$?

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

Dot product: measures similarity of $o$ and $c$
Larger dot product = larger probability

Normalize over entire vocabulary
to get probability distribution

Figure from Lena Voita

# Word Embeddings: Skip-gram

$P(u_l|v_a)$  $P(u_{saw}|v_a)$  $P(u_{cute}|v_a)$  $P(u_{grey}|v_a)$

... I  saw  a  cute  grey  cat  playing  in  the  garden  ...

$w_{t-2}$  $w_{t-1}$  $w_t$  $w_{t+1}$  $w_{t+2}$



Figure from Lena Voita

# Word Embeddings: Skip-gram



Figure from Lena Voita

# Word Embeddings: Skip-gram



$$P(u_a|v_{grey}) \quad P(u_{cute}|v_{grey}) \quad P(u_{cat}|v_{grey}) \quad P(u_{playing}|v_{grey})$$

… I saw a cute grey cat playing in the garden …

$w_{t-2}$ $w_{t-1}$ $w_t$ $w_{t+1}$ $w_{t+2}$

grey

a
cute
cat
playing

$v$ $u$

Figure from Lena Voita

# Word Embeddings: Skip-gram



$P(u_{cute}|v_{cat})$   $P(u_{grey}|v_{cat})$   $P(u_{playing}|v_{cat})$   $P(u_{in}|v_{cat})$

... I  saw  a  cute  grey  cat  playing  in  the  garden  ...

$w_{t-2}$   $w_{t-1}$   $w_t$   $w_{t+1}$   $w_{t+2}$

in
cute
grey
playing

cat

$v$        $u$

Figure from Lena Voita

# Word Embeddings: Skip-gram



$P(u_{grey}|v_{playing})$  $P(u_{cat}|v_{playing})$   $P(u_{in}|v_{playing})$  $P(u_{the}|v_{playing})$

... I  saw  a  cute  grey  cat  playing  in  the  garden  ...

$w_{t-2}$  $w_{t-1}$  $w_t$  $w_{t+1}$  $w_{t+2}$

playing

$v$     $u$

Figure from Lena Voita

# Word Embeddings: Skip-gram

$$P(u_{grey}|v_{in}) \ P(u_{cat}|v_{in}) \quad P(u_{in}|v_{in}) \ P(u_{the}|v_{in})$$

... I saw a cute grey cat playing in the garden ...

$w_{t-2} \qquad w_{t-1} \qquad w_t \qquad w_{t+1} \qquad w_{t+2}$



$v$ $\qquad$ $u$

Figure from Lena Voita

# Word Embeddings: Skip-gram

- Train using Gradient Descent

# Word Embeddings: Skip-gram

- Train using Gradient Descent
- For one word at a time, i.e., (a center word, one of the context words)

# Word Embeddings: Skip-gram

- Train using Gradient Descent
- For one word at a time, i.e., (a center word, one of the context words)
- $J_{t,j}(\theta) = -\log P(cute|cat) = -\log \dfrac{\exp u_{cute}^T v_{cat}}{\sum\limits_{w \in Voc} \exp u_w^T v_{cat}} =$

  $-u_{cute}^T v_{cat} + \log \sum\limits_{w \in Voc} \exp u_w^T v_{cat}$

# Word Embeddings: Skip-gram

**1**. Take dot product of $v_{cat}$ with **all** $u$     **2**. exp     **3**. sum all



**4**. get loss (for this one step)

$$J_{t,j}(\theta) = \underbrace{-u_{cute}^T v_{cat}}_{1} + \underbrace{\log \sum_{w \in V} \exp(u_w^T v_{cat})}_{2}$$

**5**. evaluate the gradient, make an update

$$v_{cat} := v_{cat} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{cat}}$$

$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \ \forall \ w \in V$$

Figure from Lena Voita

# Word Embeddings: Skip-gram

- Training is slow (for each central word, all the context words need to be updated)

# Word Embeddings: Skip-gram

 భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

- Training is slow (for each central word, all the context words need to be updated)
- Negative sampling: not all the context words are considered, but a random sample of them

footer_navigation">Dr. Konda Reddy Mopuri   dl - 13/ Word Embeddings   44

# Word Embeddings: Skip-gram

- Training is slow (for each central word, all the context words need to be updated)

- Negative sampling: not all the context words are considered, but a random sample of them

- Training over a large corpus leads to sufficient updates for each vector

# Word Embeddings: Skip-gram



Dot product of $v_{cat}$:
- with $u_{cute}$ - increase,
- with **all other** $u$ - decrease

Dot product of $v_{cat}$:
- with $u_{cute}$ - increase,
- with <u>a subset of other</u> $u$ - decrease

Negative samples: randomly selected K words

decrease
$u_{w1}^T v_{cat}$
$u_{w3}^T v_{cat}$
$u_{cute}^T v_{cat}$ increase
$u_{wn}^T v_{cat}$ decrease

decrease
$u_{w_{i1}}^T v_{cat}$
$u_{w_{i2}}^T v_{cat}$
$u_{cute}^T v_{cat}$ increase
$u_{w_{ik-1}}^T v_{cat}$
$u_{w_{ik}}^T v_{cat}$ decrease

Parameters to be updated:
- $v_{cat}$
- $u_w$ for all $w$ in the vocabulary

|V| + 1 vectors

Parameters to be updated:
- $v_{cat}$
- $u_{cute}$ and $u_w$ for $w$ in K negative examples

K + 2 vectors

Figure from Lena Voita

# Word Embeddings: Skip-gram

embeddings    input word

$$W_{NXm} \quad i_{mX1} \Rightarrow P_{NX1}$$

Input layer

Can be viewed as a Neural Network

# Word Embeddings: Skip-gram



Can be viewed as a Neural Network

# Word Embeddings: Skip-gram



Can be viewed as a Neural Network

# Word Embeddings: Skip-gram



Can be viewed as a Neural Network

# Word Embeddings: Skip-gram

① $W_{N \times m}$ is the $W_{\text{word}}$ (used for representing the words)

# Word Embeddings: Skip-gram

1. $W_{N \times m}$ is the $W_{\text{word}}$ (used for representing the words)
2. $W'_{m \times N}$ is the $W_{\text{context}}$ (may be ignored after the training)

# Word Embeddings: Skip-gram

1. $W_{N \times m}$ is the $W_{\text{word}}$ (used for representing the words)
2. $W'_{m \times N}$ is the $W_{\text{context}}$ (may be ignored after the training)
3. Some showed averaging word and context vectors may be more beneficial

# Bag of Words (BoW)

1. Bag of Words: Collection and frequency of words

# CBoW

1. Considers the embeddings of 'h' words before and 'h' words after the target word

# CBoW

1. Considers the embeddings of 'h' words before and 'h' words after the target word
2. Adds them (order is lost) for predicting the target word

# CBoW

The dog slept on couch

1. Size of the vocabulary $= m$

# CBoW

1. Size of the vocabulary $= m$
2. Dimension of the embeddings $= N$

# Word Embeddings: CBoW

① Input layer $W_{N \times m}$ (embeddings for the context words) projects the context (sum of 1-hot vectors of all the context vectors) into $N$-dim space

# Word Embeddings: CBoW

1. Input layer $W_{N \times m}$ (embeddings for the context words) projects the context (sum of 1-hot vectors of all the context vectors) into $N$-dim space

2. Representations of all the $(2h)$ words in the context are summed (result is an $N$-dim context vector)

# Word Embeddings: CBoW

① Input layer $W_{N \times m}$ (embeddings for the context words) projects the context (sum of 1-hot vectors of all the context vectors) into $N$-dim space

context

$$\begin{pmatrix} W_{NXm} \end{pmatrix} \begin{pmatrix} C_{mX1} \end{pmatrix} \Rightarrow \begin{pmatrix} E_{NX1} \end{pmatrix}$$

# Word Embeddings: CBoW

1. Input layer $W_{N \times m}$ (embeddings for the context words) projects the context (sum of 1-hot vectors of all the context vectors) into $N$-dim space

2. Representations of all the $(2h)$ words in the context are summed (result is an $N$-dim context vector)

context

$$\left( \quad W_{N \times m} \quad \right) \left( C_{m \times 1} \right) \Longrightarrow \left( E_{N \times 1} \right)$$

# Word Embeddings: CBoW

1. Next layer has a weight matrix $W'_{m \times N}$ (embeddings for the center words)

# Word Embeddings: CBoW

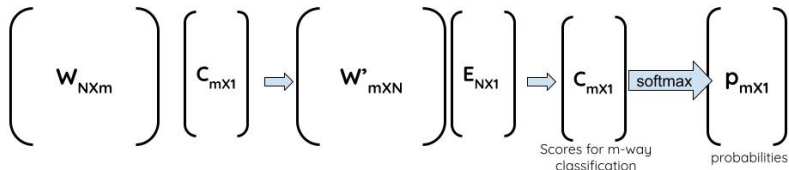1. Next layer has a weight matrix $W'_{m \times N}$ (embeddings for the center words)

2. Projects the accumulated embeddings onto the vocabulary

# Word Embeddings: CBoW

1. Next layer has a weight matrix $W'_{m \times N}$ (embeddings for the center words)

2. Projects the accumulated embeddings onto the vocabulary



Scores for m-way classification

# Word Embeddings: CBoW

① $m$- way classification $\rightarrow$ (after a softmax) maximizes the probability for the target word



Scores for m-way classification

probabilities

# Word Embeddings: CBoW

1. $W_{N \times m}$ is the $W_{\text{context}}$

# Word Embeddings: CBoW

1. $W_{N \times m}$ is the $W_{\text{context}}$
2. $W'_{m \times N}$ is the $W_{\text{words}}$

# Glove

1. Glove - Global Vectors

# Glove

1. Glove - Global Vectors
2. Combines the score-based and predict-based approaches

# Glove

1. $X_{ij}$ in the co-occurrence matrix encodes the global info. about words $i$ and $j$

# Glove

1. $X_{ij}$ in the co-occurrence matrix encodes the global info. about words $i$ and $j$
   - $p(j/i) = \frac{X_{ij}}{X_i}$
2. Glove attempts to learn representations that are faithful to the co-occurrence info.

# Glove

1. Glove attempts to learn representations that are faithful to the co-occurrence info.

# Glove

1. Glove attempts to learn representations that are faithful to the co-occurrence info.
2. Try to enforce $v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$
   - $v_i$ - central representation of word $i$, $c_j$ - context representation of word $j$

# Glove

1. Glove attempts to learn representations that are faithful to the co-occurrence info.
2. Try to enforce $v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$
   - $v_i$ - central representation of word $i$, $c_j$ - context representation of word $j$
3. Similarly, $v_j^T c_i = \log P(i/j) = \log X_{ij} - \log X_j$ (aim is to learn such embeddings $v_i$ and $c_i$)

# Glove

1. To realize the exchange symmetry of
   $$v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$$

# Glove

① To realize the exchange symmetry of
$v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$

- we may capture the $\log X_i$ as a bias $b_i$ of the word $w_i$

# Glove

① To realize the exchange symmetry of
$v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$

- we may capture the $\log X_i$ as a bias $b_i$ of the word $w_i$
- And, an additional term $\tilde{b_j}$

# Glove

1. To realize the exchange symmetry of
   $$v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$$
   - we may capture the $\log X_i$ as a bias $b_i$ of the word $w_i$
   - And, an additional term $\tilde{b}_j$

2. $v_i^T c_j + b_i + \tilde{b}_j = \log X_{ij}$

# **Glove**

1. To realize the exchange symmetry of
   $v_i^T c_j = \log P(j/i) = \log X_{ij} - \log X_i$
   - we may capture the $\log X_i$ as a bias $b_i$ of the word $w_i$
   - And, an additional term $\tilde{b}_j$

2. $v_i^T c_j + b_i + \tilde{b}_j = \log X_{ij}$

3. Since $\log X_i$ and $\log X_j$ depend on the words $i$ and $j$, they can be considered as the word specific biases (learnable)

# Glove

1. Learning objective becomes

$$\underset{v_i, c_j, b_i, \tilde{b_j}}{\text{argmin}} \; J() = \sum_{i,j} \left( v_i^T c_j + b_i + \tilde{b_j} - \log X_{ij} \right)^2$$

# Glove

1. Learning objective becomes
$$\underset{v_i, c_j, b_i, \tilde{b_j}}{\mathrm{argmin}} \; J() = \sum_{i,j} \left( v_i^T c_j + b_i + \tilde{b_j} - \log X_{ij} \right)^2$$

# Glove

1. Learning objective becomes
   $$\underset{v_i, c_j, b_i, \tilde{b_j}}{\operatorname{argmin}} J() = \sum_{i,j} \left( v_i^T c_j + b_i + \tilde{b_j} - \log X_{ij} \right)^2$$

2. Much of the entries in the co-occurrence matrix are zeros (noisy or less informative)

# Glove

1. Learning objective becomes
   $$\underset{v_i, c_j, b_i, \tilde{b_j}}{\operatorname{argmin}} J() = \sum_{i,j} \left( v_i^T c_j + b_i + \tilde{b_j} - \log X_{ij} \right)^2$$

2. Much of the entries in the co-occurrence matrix are zeros (noisy or less informative)

3. Suggests to apply a weight

# Glove



context vector

word vector

bias terms (also learned)

$$J(\theta) = \sum_{w, c \, \in V} f(\mathsf{N}(w, c)) \cdot (u_c^T v_w + b_c + \overline{b_w} - \log \mathsf{N}(w, c))^2$$

Weighting function to:
- penalize rare events
- not to over-weight frequent events

$f(x)$

1

0

$x_{max}$

x

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

$\alpha = 0.75, x_{max} = 100$

Figure from Lena Voita

# Evaluating the embeddings

1. Intrinsic - studying the internal properties (how well they capture the meaning: word similarity, analogy, etc.)

# Evaluating the embeddings

1. Intrinsic - studying the internal properties (how well they capture the meaning: word similarity, analogy, etc.)
2. Extrinsic - studying how they perform a task

# Analysing the embeddings

1. Walking the semantic space

# Analysing the embeddings

1. Walking the semantic space
2. Structure - (form clusters) nearest neighbors have a similar meaning, Linear structure

# Glove

semantic: $v(\textbf{king}) - v(\textbf{man}) + v(\textbf{woman}) \approx v(\textbf{queen})$

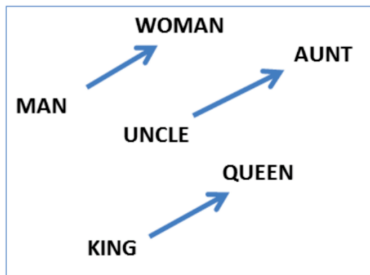syntactic: $v(\textbf{kings}) - v(\textbf{king}) + v(\textbf{queen}) \approx v(\textbf{queens})$
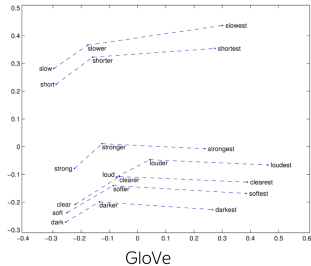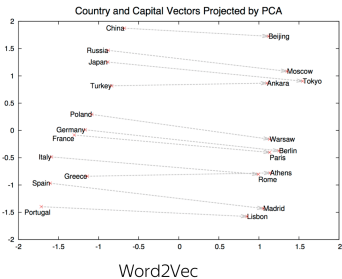


Figure from Lena Voita

# Glove



Figure from Lena Voita