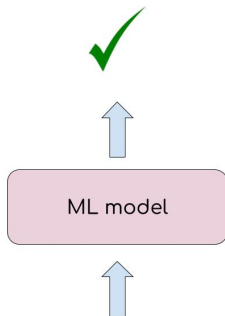


# Deep Learning

## 14 Word Embeddings

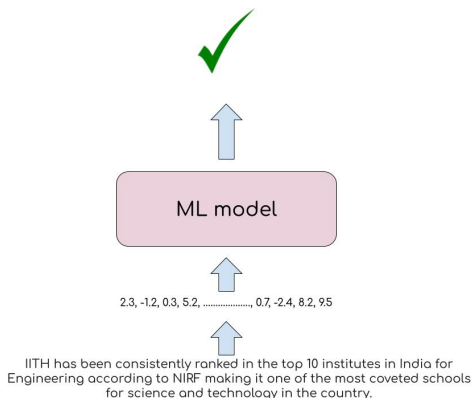
Dr. Konda Reddy Mopuri  
Dept. of AI, IIT Hyderabad  
Jan-May 2023

# Why Word Embeddings?



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?



- ① Corpus: collection of authentic text organized into dataset

# Terminology

- ① Corpus: collection of authentic text organized into dataset
- ② Vocabulary (V): Set of unique words across all the i/p streams

# Terminology

- ① Corpus: collection of authentic text organized into dataset
- ② Vocabulary ( $V$ ): Set of unique words across all the i/p streams
- ③ **Target**: Representation for every word in  $V$

# One-hot Encoding

①  $|V|$  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

# 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 **Co-occurrence matrix** 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 Representation/meaning of a word should consider its context in the corpus
- ② **Co-occurrence matrix** can capture this!
  - size: ( $\#words \times \#words$ )
  - rows: words (m), cols: context (n)
  - words and context can be of same or different size
- ③ Context can be defined as a 'h' word neighborhood

# Distributed Representations

- ① Representation/meaning of a word should consider its context in the corpus
- ② **Co-occurrence matrix** can capture this!
  - size: ( $\# \text{words} \times \# \text{words}$ )
  - rows: words ( $m$ ), cols: context ( $n$ )
  - words and context can be of same or different size
- ③ Context can be defined as a 'h' word neighborhood
- ④ Each row (column): vectorial representation of the word (context)

# Co-occurrence matrix

$$X = \begin{matrix} & \begin{matrix} I & like & enjoy & deep & learning & NLP & flying & . \end{matrix} \\ \begin{matrix} I \\ like \\ enjoy \\ deep \\ learning \\ NLP \\ flying \\ . \end{matrix} & \begin{bmatrix} 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix} \end{matrix}$$

# 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

$$\textcircled{1} X = U\Sigma V^T$$

# SVD on the Co-occurrence matrix

①  $X = U\Sigma V^T$

② 
$$\begin{bmatrix} X \\ \uparrow \dots \uparrow \\ u_1 \dots u_k \\ \downarrow \dots \downarrow \end{bmatrix}_{m \times n} = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix}_{k \times k} \cdot \begin{bmatrix} \leftarrow v_1^T \rightarrow \\ \vdots \\ \leftarrow v_k^T \rightarrow \end{bmatrix}_{k \times n}$$

# SVD on the Co-occurrence matrix

①  $X = U\Sigma V^T$

② 
$$\begin{bmatrix} X \\ \phantom{X} \end{bmatrix}_{m \times n} = \begin{bmatrix} \uparrow & \dots & \uparrow \\ u_1 & \dots & u_k \\ \downarrow & \dots & \downarrow \end{bmatrix}_{m \times k} \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix}_{k \times k} \cdot \begin{bmatrix} \leftarrow & v_1^T & \rightarrow \\ & \vdots & \\ \leftarrow & v_k^T & \rightarrow \end{bmatrix}_{k \times n}$$

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

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  $W_{\text{word}} = U_{m \times k} \cdot \Sigma_{k \times k}$

# 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

# 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

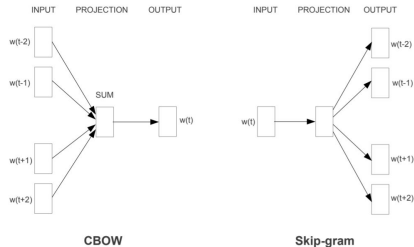


- ① Techniques we have seen so far rely on the counts (or, co-occurrence of words)
- ② Next, we see prediction based models for word embeddings

① T Mikolov et al. (2013)

- ① T Mikolov et al. (2013)
- ② Predict words from the context

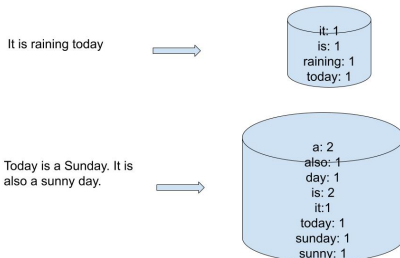
- 1 T Mikolov et al. (2013)
- 2 Predict words from the context
- 3 Two versions: Continuous Bag of Words (CBow) and Skip-gram



Caption

# Bag of Words (BoW)

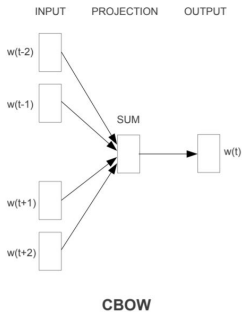
## ① Bag of Words: Collection and frequency of words



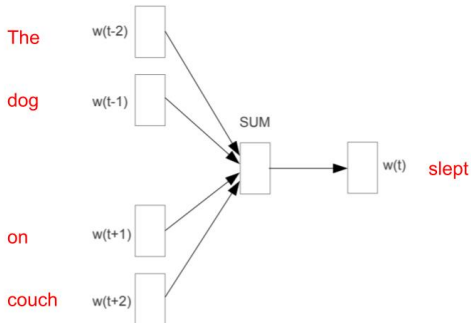


- ① Considers the embeddings of 'h' words before and 'h' words after the target word

- 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



The dog slept on couch



① Size of the vocabulary =  $m$

Vocabulary:  $m$  words,  $N$ -d real  
representation for each word

$$\left( \begin{array}{c} \mathbf{W}_{N \times m} \end{array} \right)$$



# Word Embeddings: CBoW

- ① Input layer  $W_{m \times V}$  projects the context in to  $N$ -d space

# Word Embeddings: CBoW

- ① Input layer  $W_{m \times V}$  projects the context in to  $N$ -d space
- ② Representations of all the ( $2h$ ) words in the context are summed (result is an  $V$ -d context vector)

$$\begin{matrix} & \text{context} \\ \left( \begin{matrix} W_{N \times m} \end{matrix} \right) & \left( \begin{matrix} C_{m \times 1} \end{matrix} \right) \end{matrix}$$

# Word Embeddings: CBoW

- ① Input layer  $W_{N \times m}$  projects the context in to  $N$ -d space
- ② Representations of all the ( $2h$ ) words in the context are summed (context is an  $V$ -d vector)

$$\begin{matrix} & \text{context} \\ \left( \begin{matrix} W_{N \times m} \end{matrix} \right) & \left( \begin{matrix} C_{m \times 1} \end{matrix} \right) & \Rightarrow & \left( \begin{matrix} E_{N \times 1} \end{matrix} \right)\end{matrix}$$



# Word Embeddings: CBoW

- ① Next layer has a weight matrix  $W'_{m \times N}$

# Word Embeddings: CBoW

- ① Next layer has a weight matrix  $W'_{m \times N}$
- ② Projects the accumulated embeddings onto the vocabulary

$$\begin{array}{ccc} \left( \begin{array}{c} W_{N \times m} \end{array} \right) \left( \begin{array}{c} C_{m \times 1} \end{array} \right) & \rightarrow & \left( \begin{array}{c} W'_{m \times N} \end{array} \right) \left( \begin{array}{c} E_{N \times 1} \end{array} \right) \\ \text{First layer} & & \text{Second layer} \end{array}$$

# Word Embeddings: CBoW

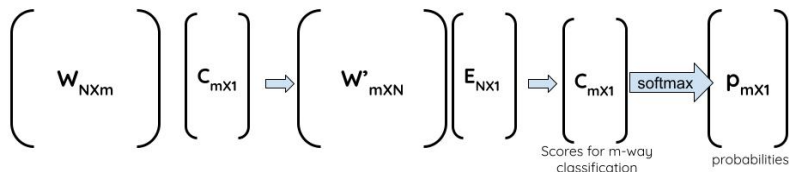
- ① Next layer has a weight matrix  $W'_{V \times N}$
- ② Projects the accumulated embeddings onto the vocabulary

$$\begin{pmatrix} W_{N \times m} \end{pmatrix} \begin{pmatrix} C_{m \times 1} \end{pmatrix} \Rightarrow \begin{pmatrix} W'_{m \times N} \end{pmatrix} \begin{pmatrix} E_{N \times 1} \end{pmatrix} \Rightarrow \begin{pmatrix} C_{m \times 1} \end{pmatrix}$$

Scores for m-way classification

# Word Embeddings: CBoW

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



# Word Embeddings: CBoW

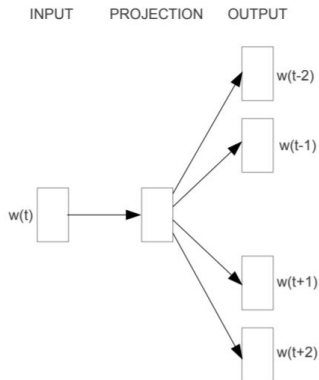
①  $W_{N \times m}$  is the  $W_{\text{context}}$

# Word Embeddings: CBoW

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

- ① Softmax at the o/p is very expensive  $\hat{y}_w = \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})}$

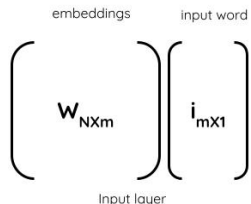
# Word Embeddings: Skip-gram



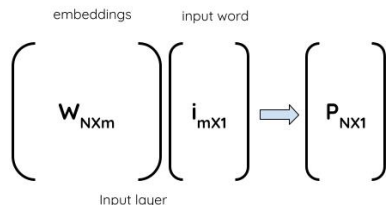
**Skip-gram**



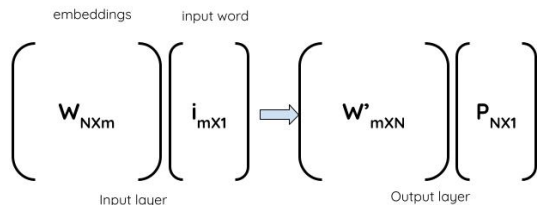
# Word Embeddings: Skip-gram



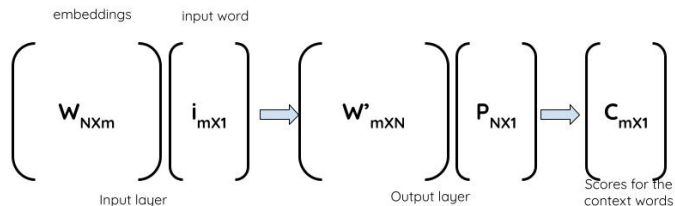
# Word Embeddings: Skip-gram



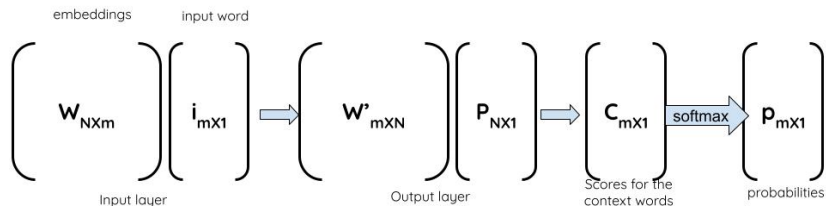
# Word Embeddings: Skip-gram



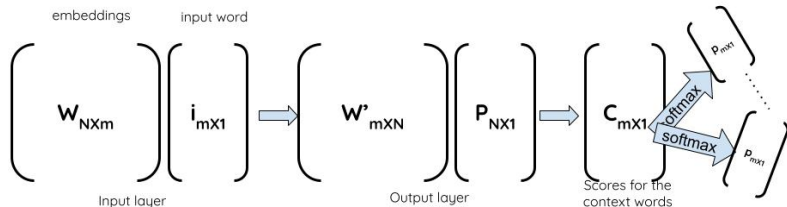
# Word Embeddings: Skip-gram



# Word Embeddings: Skip-gram



# Word Embeddings: Skip-gram



# Word Embeddings: Skip-gram

① Loss is the summation of  $\mathbb{L}(\theta) = - \sum_{i=1}^{2h} \log(\hat{y}_{w_i})$

# Word Embeddings: Skip-gram

- ① Loss is the summation of  $\mathbb{L}(\theta) = - \sum_{i=1}^{2h} \log(\hat{y}_{w_i})$
- ②  $W_{N \times m}$  is the  $W_{\text{word}}$



# Word Embeddings: Skip-gram

- ① Loss is the summation of  $\mathbb{L}(\theta) = - \sum_{i=1}^{2h} \log(\hat{y}_{w_i})$
- ②  $W_{N \times m}$  is the  $W_{\text{word}}$
- ③  $W'_{m \times N}$  is the  $W_{\text{context}}$

# Skip-gram: Issues

- ① Expensive softmax operation at the o/p (same as that of CBoW)

# Skip-gram: Issues

- ① Expensive softmax operation at the o/p (same as that of CBoW)
- ② Negative sampling: subset of incorrect words participate (instead of all)

# Skip-gram: Issues

- ① Expensive softmax operation at the o/p (same as that of CBoW)
- ② Negative sampling: subset of incorrect words participate (instead of all)
- ③ Other solutions: Contrastive estimation, and hierarchical softmax

## ① Glove - Global Vectors

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

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

- ①  $X_{ij}$  in the cooccurrence matrix encodes the global info. about words  $i$  and  $j$
- ② Glove attempts to learn representations that are faithful to the cooccurrence info.



- ①  $X_{ij}$  in the cooccurrence matrix encodes the global info. about words  $i$  and  $j$
- ② Glove attempts to learn representations that are faithful to the cooccurrence info.
- ③  $v_i^T v_j = \log P(j/i) = \log X_{ij} - \log X_i$

- ①  $X_{ij}$  in the cooccurrence matrix encodes the global info. about words  $i$  and  $j$
- ② Glove attempts to learn representations that are faithful to the cooccurrence info.
- ③  $v_i^T v_j = \log P(j/i) = \log X_{ij} - \log X_i$
- ④  $v_j^T v_i = \log P(i/j) = \log X_{ij} - \log X_j$

① (add the two equations)  $v_i^T v_j = \log X_{ij} - \frac{1}{2} \log X_i - \frac{1}{2} \log X_j$

- ① (add the two equations)  $v_i^T v_j = \log X_{ij} - \frac{1}{2} \log X_i - \frac{1}{2} \log X_j$
- ② Since  $\log X_i$  and  $\log X_j$  depend on the words  $i$  and  $j$ , they can be considered as the word specific biases

- ① (add the two equations)  $v_i^T v_j = \log X_{ij} - \frac{1}{2} \log X_i - \frac{1}{2} \log X_j$
- ② Since  $\log X_i$  and  $\log X_j$  depend on the words  $i$  and  $j$ , they can be considered as the word specific biases
- ③  $v_i^T v_j + b_i + b_j = \log X_{ij}$

- ① (add the two equations)  $v_i^T v_j = \log X_{ij} - \frac{1}{2} \log X_i - \frac{1}{2} \log X_j$
- ② Since  $\log X_i$  and  $\log X_j$  depend on the words  $i$  and  $j$ , they can be considered as the word specific biases
- ③  $v_i^T v_j + b_i + b_j = \log X_{ij}$
- ④  $\sum_{i,j} (v_i^T v_j + b_i + b_j - \log X_{ij})^2$

# Evaluating the embeddings

- ① Semantic relatedness (compute the correlation with humans' similarity)

# Evaluating the embeddings

- ① Semantic relatedness (compute the correlation with humans' similarity)
- ② Synonym detection



# Evaluating the embeddings

- ① Semantic relatedness (compute the correlation with humans' similarity)
- ② Synonym detection
- ③ Analogy

# Comparison among different models

① Difficult to judge!

# Comparison among different models

- ① Difficult to judge!
- ② Some studies favor the predict-based, some the cooccurrence based!!