

# Word Embedding

講者：Isaac

# Outline

---

- ▶ Bag of Words
- ▶ TF-IDF vector
- ▶ N-Gram
- ▶ Skip-gram and CBOW
- ▶ GloVe
- ▶ Doc2vect

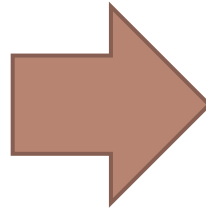


# What's word Embedding

---

- ▶ Word embedding is the collective name that words or phrases from the vocabulary are mapped to vectors of real numbers

**“hello”**



$$\begin{bmatrix} 0.3 \\ -0.2 \\ \vdots \\ 0.8 \\ -0.5 \end{bmatrix}$$

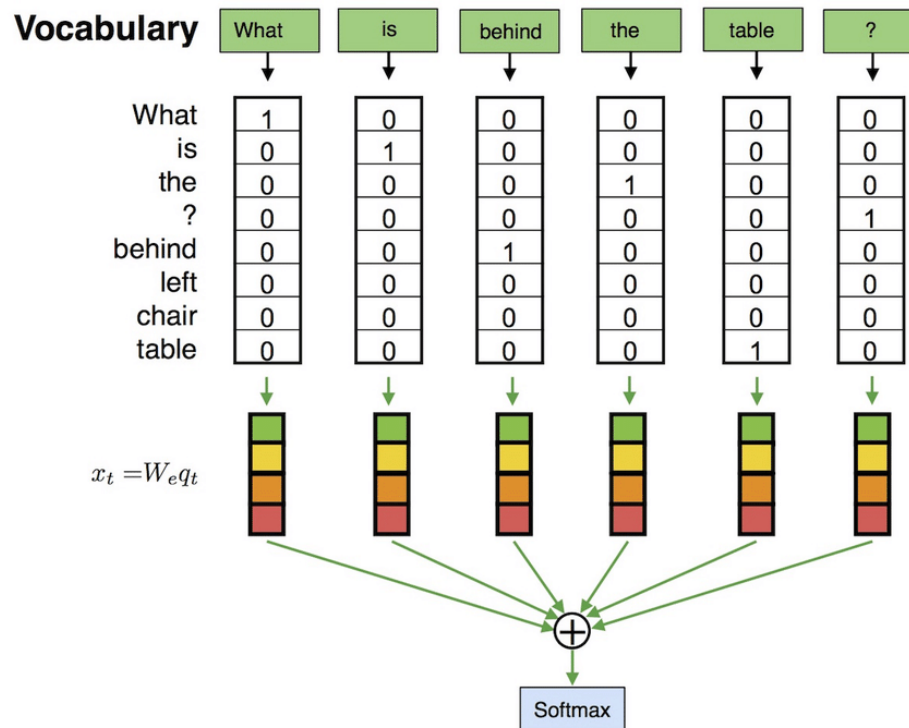
---

# Bag of Words



# What's Bag of Words

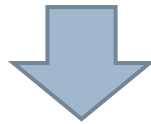
- ▶ one of the simplest methods of embedding words into numerical vectors



# What's Bag of Words

---

Document 1	High five!
Document 2	I am old.
Document 3	She is five.



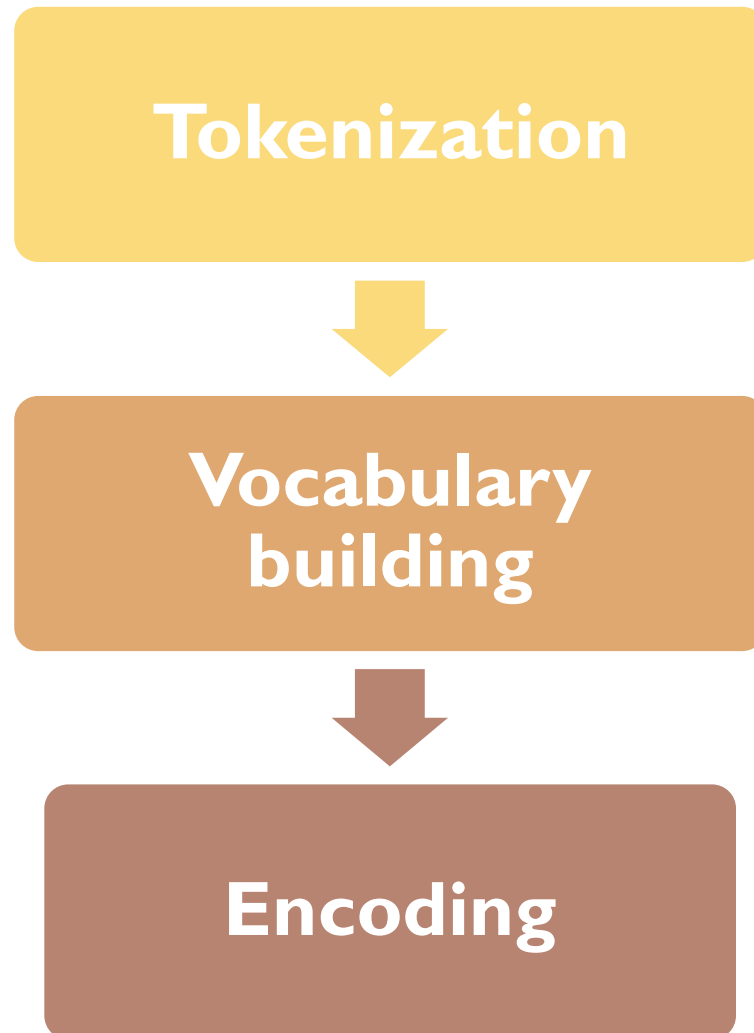
	Document 1	Document 2	Document 3
High	1	0	0
Five	1	0	1
I	0	1	0
am	0	1	0
old	0	1	0
She	0	0	1
is	0	0	1

“Five”:[1,0,1]

“Document 2”:[0,0,1,1,1,0,0]

# What's Bag of Words

---



# Bag of Words - Tokenization

---

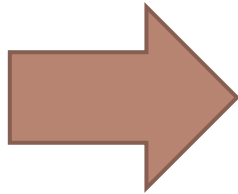
Document	Content
1	["it", "was", "the", "best", "of", "times"]
2	["it", "was", "the", "worst", "of", "times"]
3	["it", "was", "the", "age", "of", "wisdom"]
4	["it", "was", "the", "age", "of", "foolishness"]



# Bag of Words - Vocabulary building

---

“it”  
“was”  
“the”  
“best”  
“of”  
“times”  
“worst”  
“age”  
“wisdom”  
“foolishness”



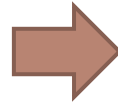
word	word ID
it	1
was	2
the	3
best	4
of	5
times	6
worst	7
Age	8
wisdom	9
foolishness	10

assuming ignore case and punctuation

# Bag of Words - Encoding

---

word	word ID
it	1
was	2
the	3
best	4
of	5
times	6
worst	7
Age	8
wisdom	9
foolishness	10



"it was the worst of times" = [1, 1, 1, 0, 1, 1, 1, 0, 0, 0]

"it was the age of wisdom" = [1, 1, 1, 0, 1, 0, 0, 1, 1, 0]

"it was the age of foolishness" = [1, 1, 1, 0, 1, 0, 0, 1, 0, 1]

# Drawback of Bag of Words

---

- ▶ ignores word order
- ▶ sparse vector problems
  - ▶ solved by filtering out stop word, stemming, lemmazation, ignoring case etc ... ..

# What's tf-idf?

---

- ▶ **tf-idf**(term frequency–inverse document frequency) is a method that reflect how important a in a collection of corpus
  - ▶ **tf** means term “frequency”
  - ▶ **idf** means term “inverse document frequency”

# Term Frequency (TF)

---

- ▶ In a single document
  - ▶ measures how frequently a term occurs in a document

$TF = (\text{Number of time the word occurs in the text}) / (\text{Total number of words in text})$

# Inverse Data Frequency (IDF)

---

- ▶ Among multiple documents

$\text{IDF} = \log(\text{Total number of documents} / \text{Number of documents with word } t \text{ in it})$

# tf-idf example

---

Document	Content
Document 1	It is going to rain today.
Document 2	Today I am not going outside.
Document 3	I am going to watch the season premiere.

# tf-idf example

---

Words / Documents	Document 1	Document 2	Document 3
going	0.16	0.16	0.12
to	0.16	0	0.12
today	0.16	0.16	0
i	0	0.16	0.12
am	0	0.16	0.12
it	0.16	0	0
is	0.16	0	0
rain	0.16	0	0

calculate TF on each document



# tf-idf example

---

Words / Documents	IDF
going	$\text{Log}(3/3)$
to	$\text{Log}(3/2)$
today	$\text{Log}(3/2)$
i	$\text{Log}(3/2)$
am	$\text{Log}(3/2)$
it	$\text{Log}(3/1)$
is	$\text{Log}(3/1)$
rain	$\text{Log}(3/1)$

calculate TF on each document

# tf-idf example

---

Words/ Documents	going	to	today	I	am	it	is	rain
Document1	0	0.07	0.07	0	0	0.17	0.17	0.17
Document2	0	0	0.07	0.07	0.07	0	0	0
Document3	0	0.05	0	0.05	0.05	0	0	0

Document1 = [0 , 0.07 , 0.07 , 0 , 0 , 0.17 , 0.17 , 0.17 ]

---

N-Gram

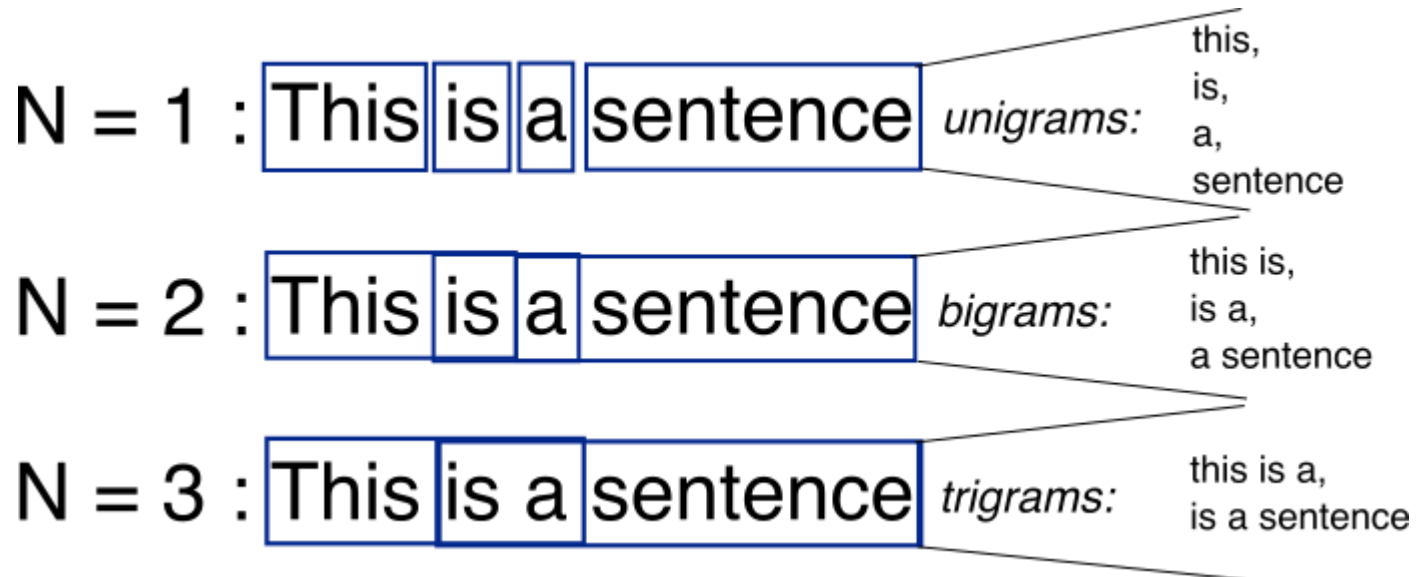
---



# What's N-Gram

---

- ▶ n-gram is a contiguous sequence of n items from a given sample of text or speech

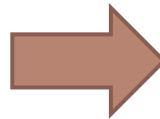


# What's N-Gram

---

- ▶ Use N-gram to form word vector
  - ▶ same as bag of word

word	word ID
this is	1
is a	2
a sentence	3
a cat	4
a dog	5



“this is a dog”

[1,1,0,0,1]

# Language model

---

- ▶ Probability a sentence occur in a text



# Language model

---

the large green \_\_ . Possible answer may be “mountain” or “tree” ?

Kate swallowed the large green \_\_ . Possible answer may be “pill” or “broccoli” ?

---

S1= “我剛吃過晚飯”

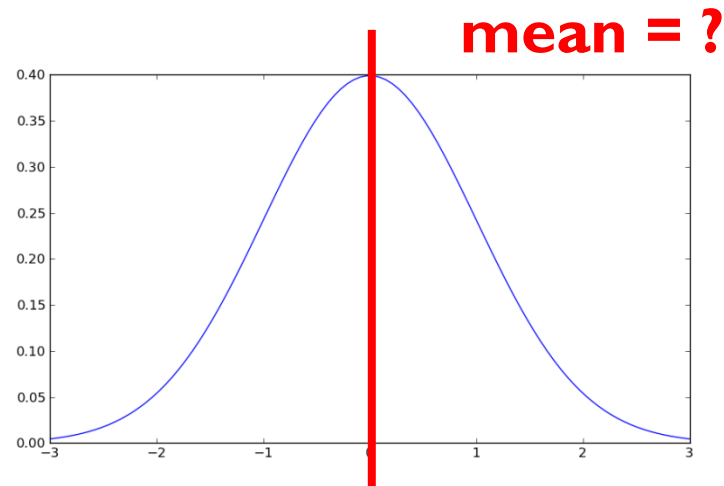
S2= “剛我過晚飯吃”

Which sentence is more reasonable?

# Likelihood Estimation

- ▶ Assume students' height is normal distribution

Student ID	1	2	3	4	5
Height (cm)	162	164	170	168	166



What's mean of the normal distribution?



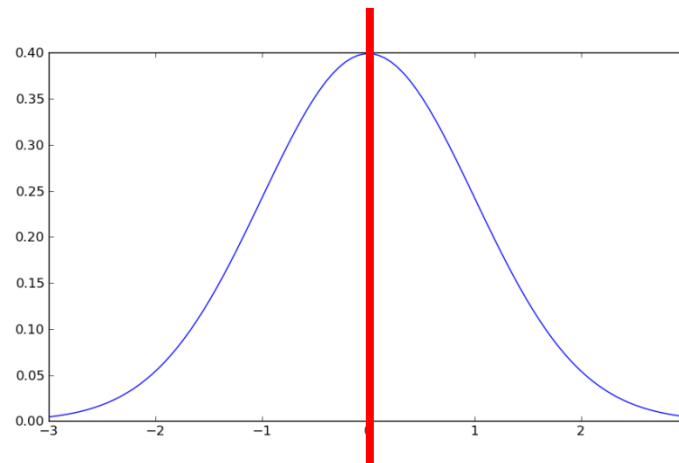
# Likelihood Estimation

- ▶ Assume students' height is normal distribution

Student ID	1	2	3	4	5
Height (cm)	162	164	170	168	166

**Intuitively:**

$$\text{mean} = (162+164+170+168+166)/5 = 166$$



**But why.....?**

# Maximum Likelihood Estimation

---

- ▶ **Maximum likelihood estimation (MLE)** is a technique used for estimating the parameters of a given distribution, using some observed data

*Given probability distribution  $f$  that have some unknown parameters  $\theta$*

*$x_1, x_2, \dots, x_n$  are observation from  $f$*

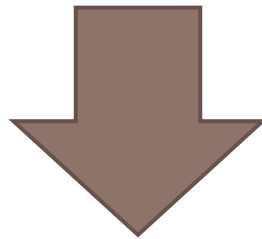
*Likelihood function:  $f(x_1, x_2, \dots, x_n | \theta) = f(x_1 | \theta) * f(x_2 | \theta) * \dots * f(x_n | \theta)$*

# Maximum Likelihood Estimation

---

Usually, we use Maximum **log** likelihood because Maximum likelihood is hard to calculate

*Likelihood function:*  $f(x_1, x_2, \dots, x_n | \theta) = f(x_1 | \theta) * f(x_2 | \theta) * \dots * f(x_n | \theta)$



*Log Likelihood function:*  $f(x_1, x_2, \dots, x_n | \theta) = \sum_{i=1}^n \log(f(x_i | \theta))$

# Likelihood Estimation

---

Assume apple weights is normal distribution

We got three apples and their weights are 9, 9.5, 11 respectively

**Our goal is to estimate the mean/std of total apples on the tree**

$$P(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{(x - \mu)^2}{2\sigma^2} \right)$$

# Likelihood Estimation

---

We want to maximize likelihood(the following equation )

$$P(9, 9.5, 11; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9 - \mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9.5 - \mu)^2}{2\sigma^2}\right) \\ \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(11 - \mu)^2}{2\sigma^2}\right)$$

**But it is hard to derivative on this equation !**

# Likelihood Estimation

---

**maximize likelihood equivalent to maximize log likelihood**

$$P(9, 9.5, 11; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9 - \mu)^2}{2\sigma^2}\right) \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(9.5 - \mu)^2}{2\sigma^2}\right) \\ \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(11 - \mu)^2}{2\sigma^2}\right)$$

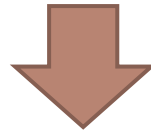


$$\ln(P(x; \mu, \sigma)) = \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(9 - \mu)^2}{2\sigma^2} + \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(9.5 - \mu)^2}{2\sigma^2} \\ + \ln\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(11 - \mu)^2}{2\sigma^2}$$

# Likelihood Estimation

---

$$\frac{\partial \ln(P(x; \mu, \sigma))}{\partial \mu} = \frac{1}{\sigma^2} [9 + 9.5 + 11 - 3\mu] .$$



$$\mu = \frac{9 + 9.5 + 11}{3} = 9.833$$

# Likelihood Estimation Table

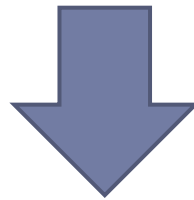
---

Distribution	Estimated parameters
$\frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$	$\mu = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$ $\sigma = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
$\lambda e^{-\lambda x}$	$\lambda = \frac{1}{\bar{x}}$
$\frac{e^{-\lambda} \lambda^k}{k!}$	$\lambda = \bar{x}$
$\vdots$	$\vdots$



- 
- ▶ Assume we have  $m$  word sequence, the probability of this sentence is

$$P(w_1, w_2, \dots, w_m) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots P(w_m|w_1, \dots, w_{m-1})$$



Assume Markov Property

$$P(w_i|w_1, \dots, w_{i-1}) = P(w_i|w_{i-n+1}, \dots, w_{i-1})$$

---

$$P(w_i|w_1, \dots, w_{i-1}) = P(w_i|w_{i-n+1}, \dots, w_{i-1})$$

unigram model (n=1)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i)$$

$$P(w_i) = \frac{C(w_i)}{M}$$

bigram model (n=2)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i|w_{i-1})$$

$$P(w_i|w_{i-1}) = \frac{C(w_{i-1}w_i)}{C(w_{i-1})}$$

trigram model (n=3)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i|w_{i-2}w_{i-1})$$

$$P(w_i|w_{i-n+1}, \dots, w_{i-1}) = \frac{C(w_{i-n+1}, \dots, w_i)}{C(w_{i-n+1}, \dots, w_{i-1})}$$

# Example

---

- ▶ Assume we have a corpus

Sentence	Content
1	<s1> <s2> yes no no no no yes </s1> </s2>
2	<s1> <s2> no no no yes yes yes no </s1> </s2>

We want to calculate probability of the following sentence:

<s1> <s2> yes no no yes </s1> </s2>

# Example

---

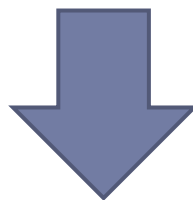
Use trigram as example

$$P(\text{yes} | \langle s1 \rangle \langle s2 \rangle) = \frac{1}{2}, \quad P(\text{no} | \langle s2 \rangle \text{yes}) = 1$$

$$P(\text{no} | \text{yes no}) = \frac{1}{2}, \quad P(\text{yes} | \text{no no}) = \frac{2}{5}$$

$$P(\langle /s2 \rangle | \text{no yes}) = \frac{1}{2}, \quad P(\langle /s1 \rangle | \text{yes} \langle /s2 \rangle) = 1$$

$\langle s1 \rangle \langle s2 \rangle \text{yes no no yes} \langle /s1 \rangle \langle /s2 \rangle$



$$\frac{1}{2} \times 1 \times \frac{1}{2} \times \frac{2}{5} \times \frac{1}{2} \times 1 = 0.05$$

---

lots of |

lots of love

lots of fish

lots of discharge

lots of lollies

Press Enter to search.

search engine



text generation

---

# Skip-gram and CBOW



# Skip-gram and CBOW

---

- ▶ NN-based algorithm
- ▶ Core concept
  - ▶ the meaning of a word can be inferred by the company it keeps

I like math.

I like programming.

Today is Friday.

Today is a good day.

- *like* is the context of target *I*
- *math* is the context of target *like*
- *programming* is also the context of target *like*

# Skip-Gram

---

Transform many sentences into training pairs

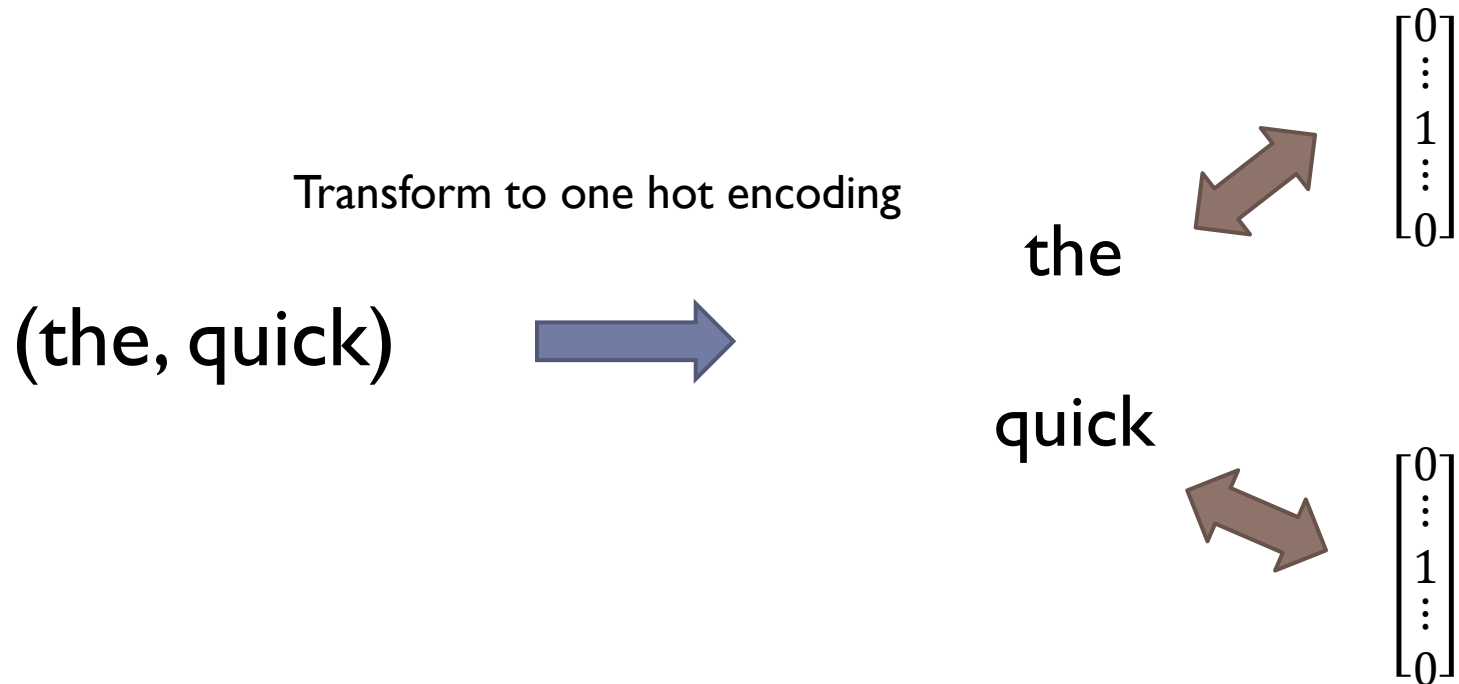
Source Text	Training Samples
<div>The quick brown fox jumps over the lazy dog.</div> <div><div>The quick brown</div> →</div>	(the, quick) (the, brown)
<div>The quick brown fox jumps over the lazy dog.</div> <div><div>The quick brown fox</div> →</div>	(quick, the) (quick, brown) (quick, fox)
<div>The quick brown fox jumps over the lazy dog.</div> <div><div>The quick brown fox jumps</div> →</div>	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
<div>The quick brown fox jumps over the lazy dog.</div> <div><div>The quick brown fox jumps over</div> →</div>	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Window size = 2



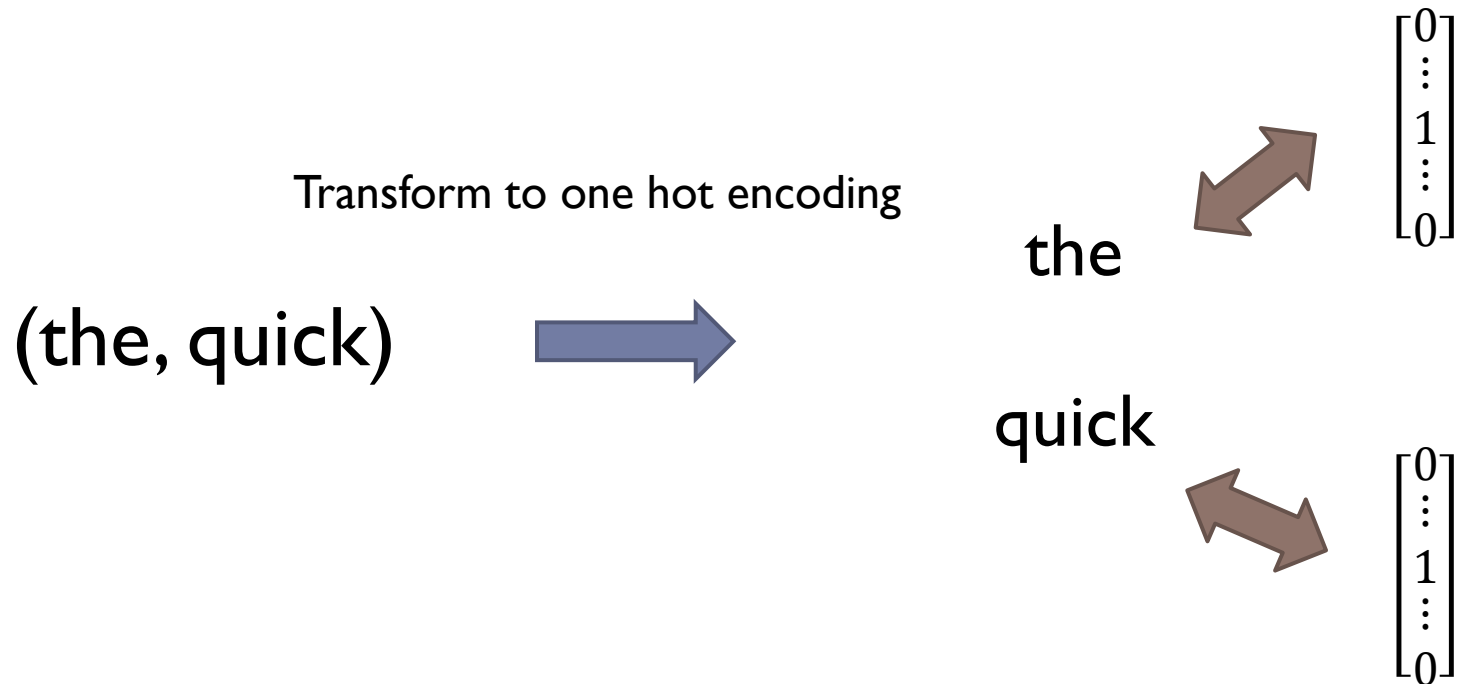
# Skip-Gram

---

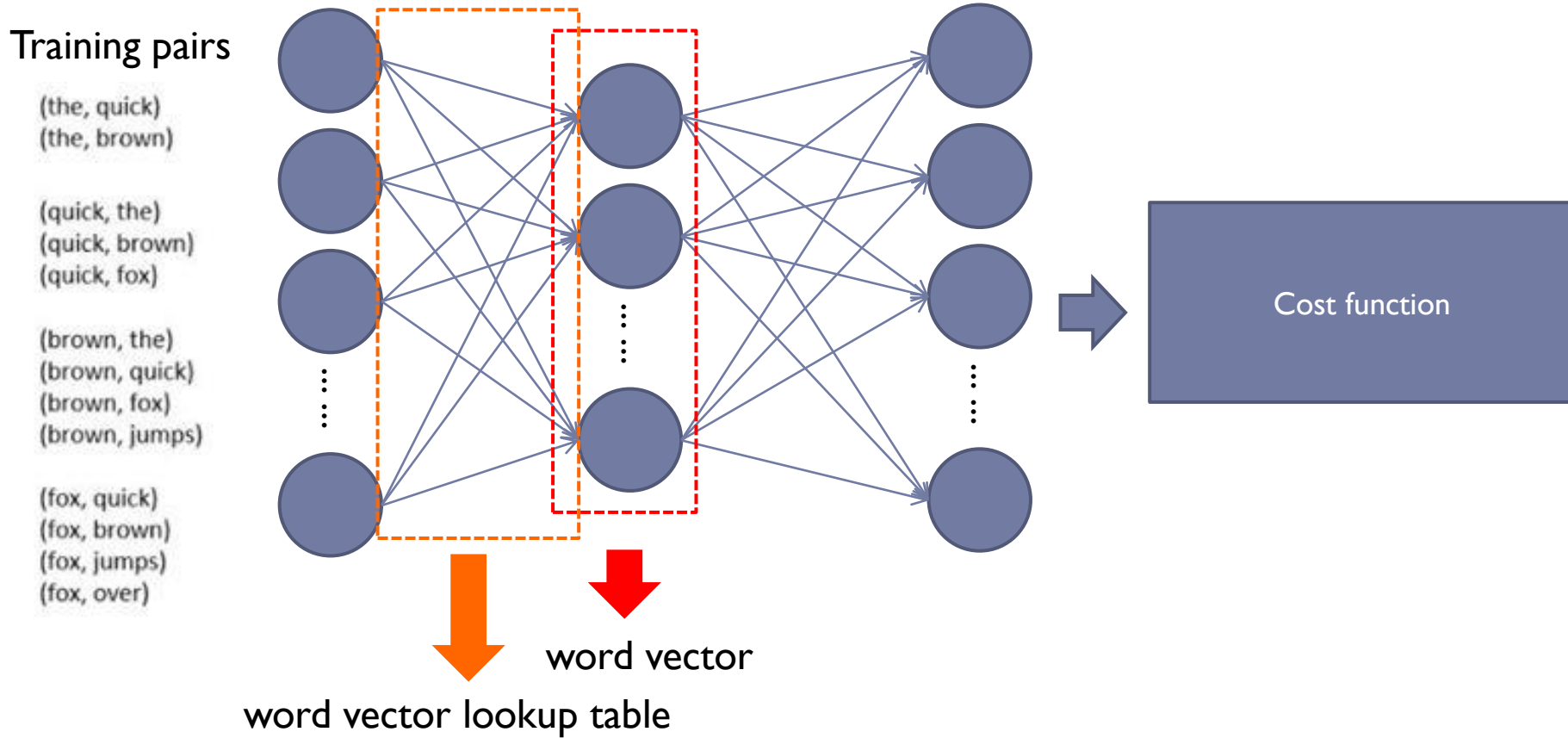


# Skip-Gram

---

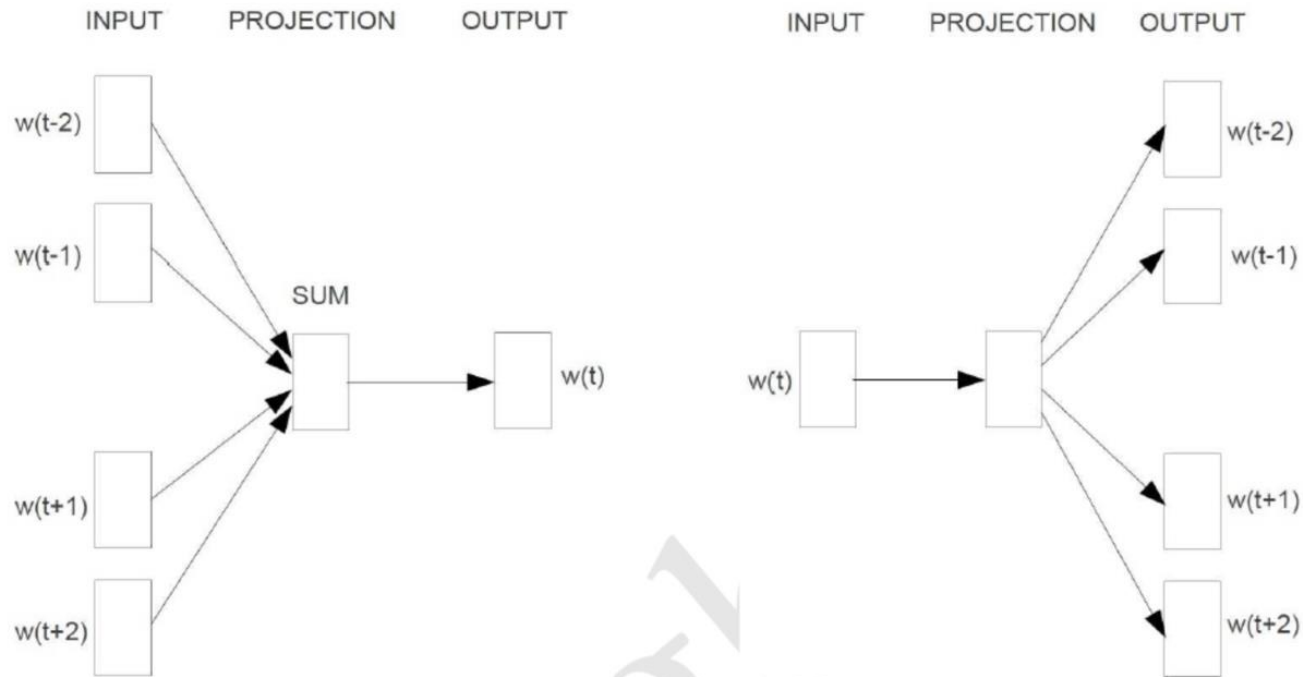


# Skip-Gram



# Skip-Gram V.S. CBOW

---



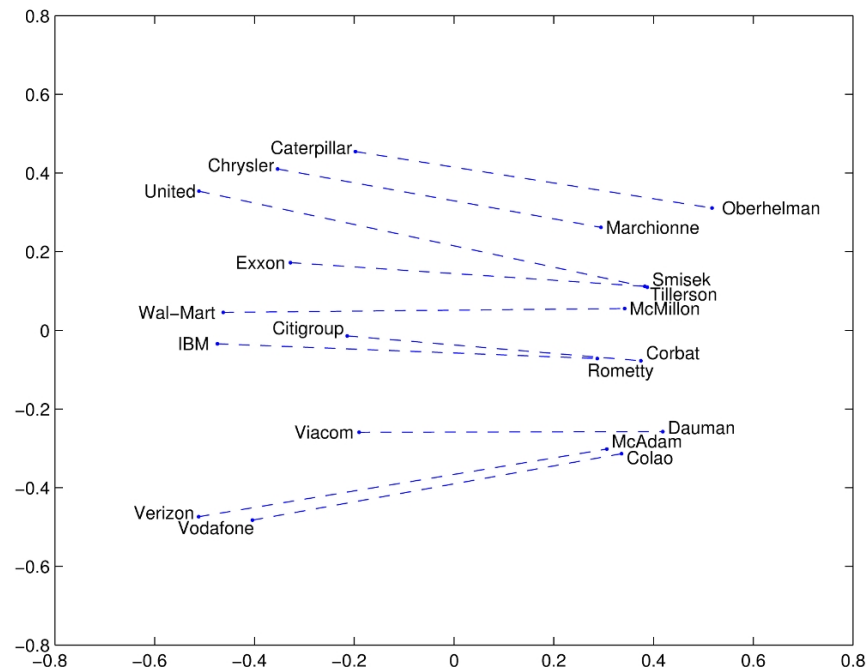
---

GloVe



# What's GloVe

- ▶ GloVe(Global Vectors for Word Representation) is a count-based and overall statistics word representation method



# Construct Co-occurrence Matrix

---

$X =$

	the	cat	sat	on	mat
the	0	1	0	1	1
cat	1	0	1	0	0
sat	0	1	0	1	0
on	1	0	1	0	0
mat	1	0	0	0	0

---

$$w_i^T \tilde{w}_j + b_i + \tilde{b}_j = \log(X_{ij})$$

$w_i, w_j$ : word vector

$b_i, b_j$ : bias term

$X_{ij}$ : entity  $i, j$  in co – occurrence matrix

relationship between word vector and co-occurrence matrix



# Define Cost

---

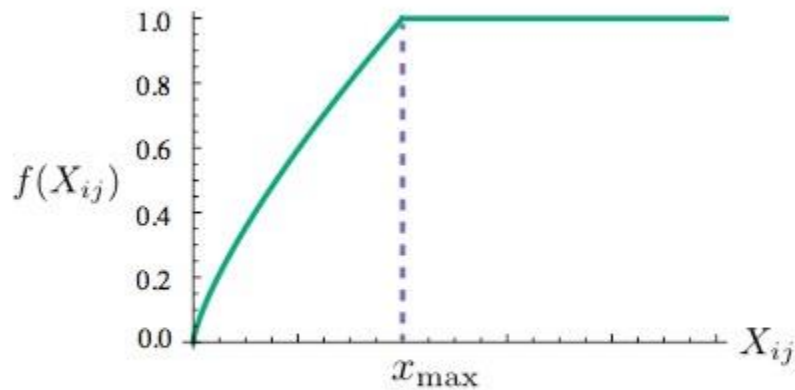
$$J = \sum_{i,j=1}^V \boxed{f(X_{ij})} (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log(X_{ij}))^2$$

avoid  $X_{ij}$  is zero

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

# Define Cost

---



$$f(x) = \begin{cases} \left(\frac{x}{x_{\max}}\right)^{\alpha}, & \text{if } x < x_{\max} \\ 1 & , \text{if } x \geq x_{\max} \end{cases}$$

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

---

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
$P(k \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
$P(k \text{ice})/P(k \text{steam})$	8.9	$8.5 \times 10^{-2}$	1.36	0.96

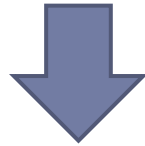
If word  $k$  is very similar to ice but irrelevant to steam (e.g.  $k=\text{solid}$ )  
 $\rightarrow P(k|\text{ice})/P(k|\text{steam})$  will be very high ( $>1$ )

If word  $k$  is very similar to steam but irrelevant to ice (e.g.  $k=\text{gas}$ )  
 $\rightarrow P(k|\text{ice})/P(k|\text{steam})$  will be very small ( $<1$ )

If word  $k$  is related or unrelated to either words  
 $\rightarrow P(k|\text{ice})/P(k|\text{steam})$  will be close to 1

---

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$



$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

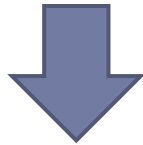
Word vectors are linear systems  
 $\text{vec}(\text{king}) - \text{vec}(\text{male}) + \text{vec}(\text{queen}) = \text{vec}(\text{female})$



$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$

Make LHS scalar

$$F((w_i - w_j)^T \tilde{w}_k)$$



homomorphism property  
 $F(A-B) = F(A)/F(B)$

$$F((w_i - w_j)^T \tilde{w}_k) = \frac{F(w_i^T \tilde{w}_k)}{F(w_j^T \tilde{w}_k)}$$

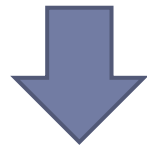
$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$



$$F(w_i^T \tilde{w}_k) = P_{ik} = \frac{X_{ik}}{X_i}$$

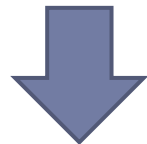
---

$$F(w_i^T \tilde{w}_k) = P_{ik} = \frac{X_{ik}}{X_i}$$



assume  $F$  is exp

$$w_i^T \tilde{w}_k = \log(P_{ik}) = \log(X_{ik}) - \log(X_i)$$



assume  $b_i + \tilde{b}_j = X_i$

$$w_i^T \tilde{w}_j + b_i + \tilde{b}_j = \log(X_{ij})$$

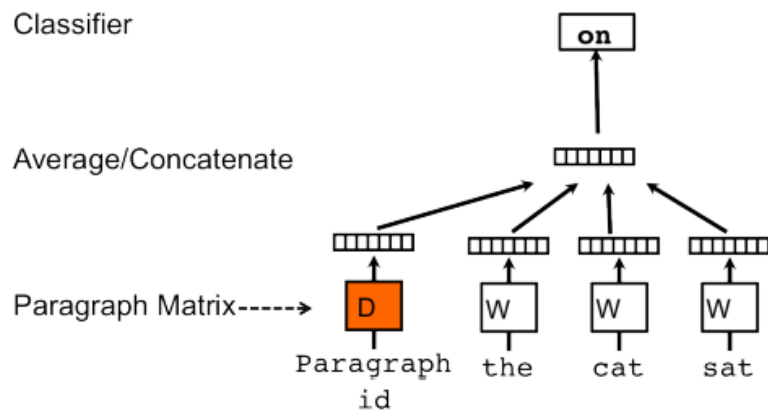
---

Doc2vect

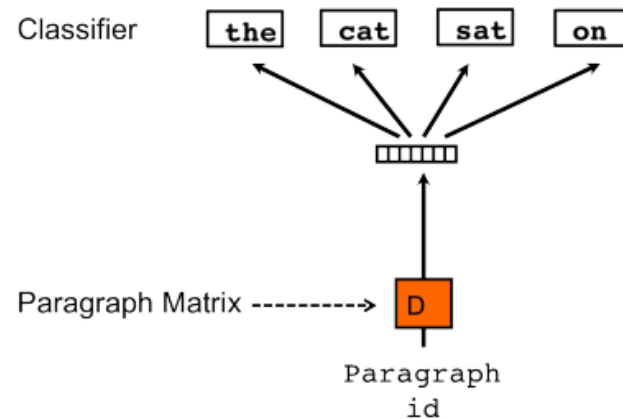


# Doc2vect

---



Distributed Memory version of Paragraph Vector  
(PV-DM)



Words version of Paragraph Vector  
(PV-DBOW)