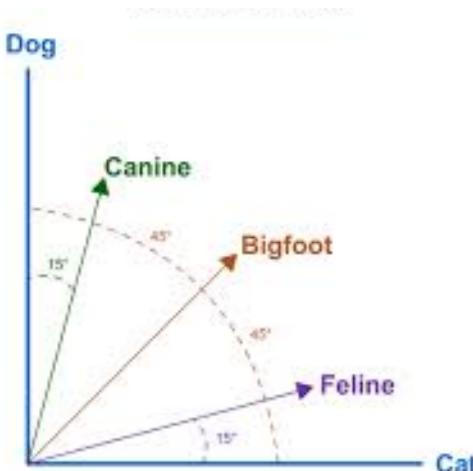


# Distributional Semantics

COMP90042

Natural Language Processing  
Lecture 10



# Lexical Databases - Problems

- Manually constructed
  - ▶ **Expensive**
  - ▶ Human annotation can be **biased** and **noisy**
- Language is **dynamic**
  - ▶ New words: slang, terminology, etc.
  - ▶ New **senses**
- The Internet provides us with massive amounts of text. Can we use that to obtain word meanings?

# Distributional Hypothesis

- “You shall know a word by the company it keeps” (Firth, 1957) *look neighbouring words I guess the meaning.*
- Document **co-occurrence** often indicative of topic (*document* as context)
  - ▶ E.g. *voting* and *politics*
- Local context reflects a word’s semantic class (*word window* as context)
  - ▶ E.g. *eat a pizza*, *eat a burger*

# Guessing Meaning from Context

- Learn unknown word from its usage
- *E.g., tezgüino*
  - (14.1) A bottle of \_\_\_\_\_ is on the table.
  - (14.2) Everybody likes \_\_\_\_\_.
  - (14.3) Don't have \_\_\_\_\_ before you drive.
  - (14.4) We make \_\_\_\_\_ out of corn.
- Look at other words in same (or similar) contexts

	(14.1)	(14.2)	(14.3)	(14.4)	...
<i>tezgüino</i>	1	1	1	1	
<i>loud</i>	0	0	0	0	
<i>motor oil</i>	1	0	0	1	
<i>tortillas</i>	0	1	0	1	
<i>choices</i>	0	1	0	0	
<i>wine</i>	1	1	1	0	

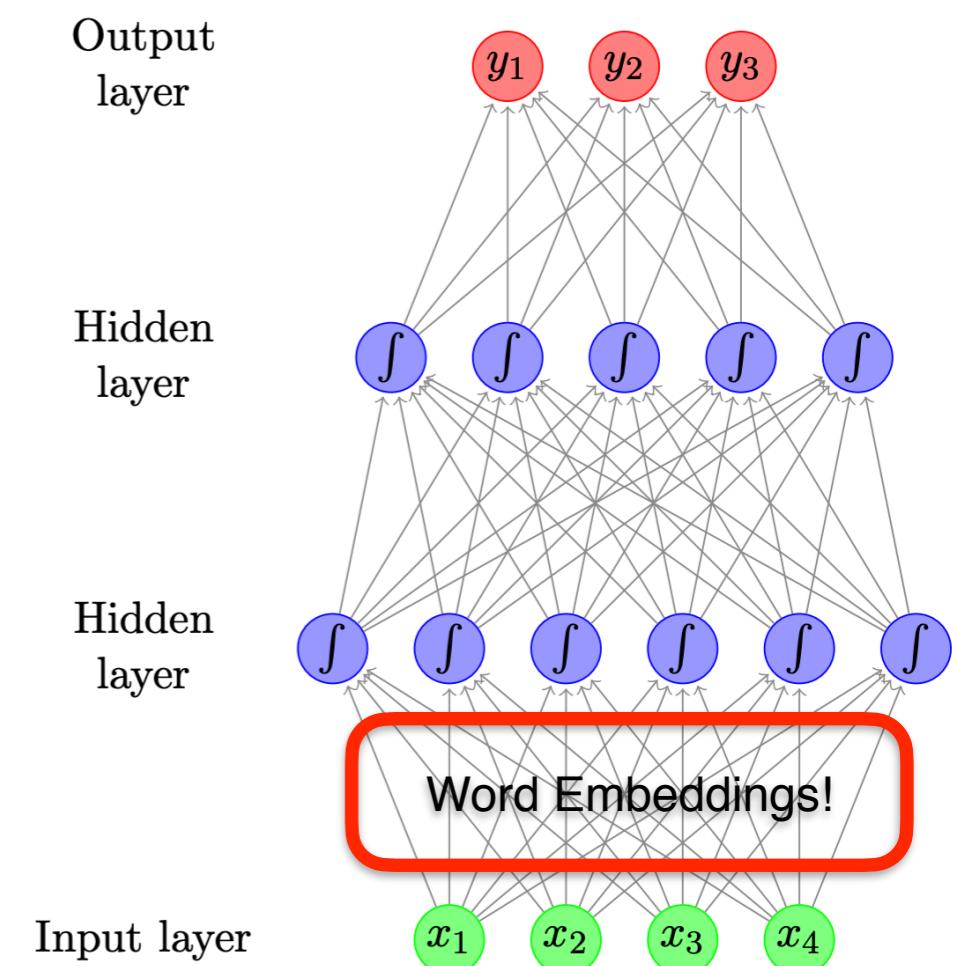
# Word Vectors

	(14.1)	(14.2)	(14.3)	(14.4)	...
<i>tezgüino</i>	1	1	1	1	
<i>loud</i>	0	0	0	0	
<i>motor oil</i>	1	0	0	1	
<i>tortillas</i>	0	1	0	1	
<i>choices</i>	0	1	0	0	
<i>wine</i>	1	1	1	0	

- Each row can be thought of a **word vector**
- It describes the **distributional properties**
  - Capture all sorts of **semantic relationships** (synonymy, analogy, etc)
    - *sort of neighbouring words that can occur around*

# Word Embeddings?

- We've seen word vectors before: word embeddings!
- Here we will learn other ways to produce word vectors
  - ▶ Count-based methods
  - ▶ More efficient neural methods designed just for learning word vectors



# Count-Based Methods

# The Vector Space Model

- Fundamental idea: represent meaning as a vector
- Consider documents as context
- One matrix, two viewpoints
  - Documents represented by their words
  - Words represented by their documents

Term matrix

	...	state	fun	heaven	...
...	0	1	0		
425	0	1	0		
426	3	0	0		
427	0	0	0		
.....					

document  
occurs and  
of words.  
Bow of Doc 425

# Manipulating the VSM

- Weighting the values (beyond frequency)
- Creating low-dimensional dense vectors
- Comparing vectors

*Sparsity problem*

for 'the'  
 $|D| = 500$   
 $df_w = 500$

## Tf-idf

means to

- Standard weighting scheme for information retrieval
- Discounts common words

penalize

	...	the	country	hell	...
...					
425		43	5	1	
426		24	1	0	
427		37	0	3	
...					
df		500	14	7	

tf matrix

total no. of documents

$$idf_w = \log \frac{|D|}{df_w}$$

document frequency.

	...	the	country	hell	...
...					
425		0	26.0	6.2	
426		0	5.2	0	
427		0	0	18.6	
...					

$tf \times idf$   
tf-idf matrix  
term frequency.

# Dimensionality Reduction

- Term-document matrices are very **sparse**
- Dimensionality reduction: create shorter, denser vectors
- More practical (less features)  
 $\equiv$  Time, Space, efficiency
- Remove noise (less overfitting)

# Singular Value Decomposition

$|D|$

$$|V| \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$A$

(term-document matrix)

$$A = U\Sigma V^T$$

$U$

(new term matrix)

$m$

Diagonal  
Matrix  $\Sigma$

(singular values)

$m$

$V^T$

(new document matrix)

$|D|$

$$|V| \begin{bmatrix} 2.2 & 0.3 & \dots & 8.7 \\ 5.5 & -2.8 & \dots & 0.1 \\ -1.3 & 3.7 & \dots & 3.5 \\ \vdots & & \ddots & \vdots \\ 2.9 & -2.1 & \dots & -1.9 \end{bmatrix}$$

Denser Matrix

$$\begin{bmatrix} 9.1 & 0 & 0 & \dots & 0 \\ 0 & 4.4 & 0 & \dots & 0 \\ 0 & 0 & 2.3 & \dots & 0 \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0.1 \end{bmatrix}$$

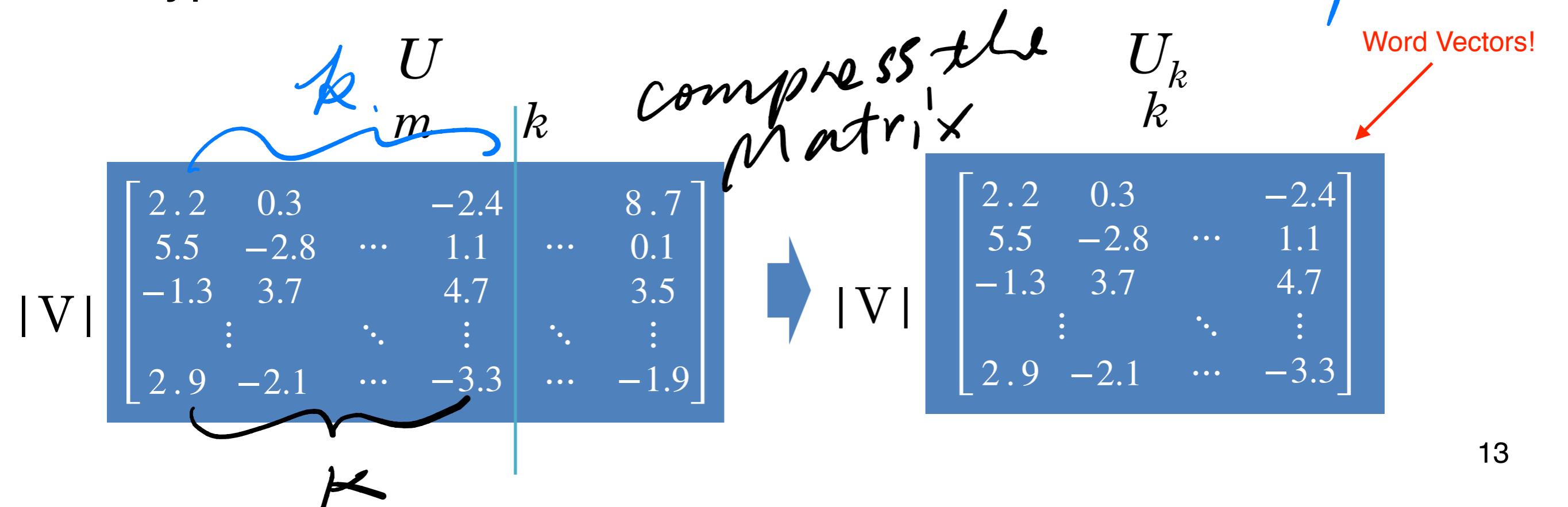
$m = \text{Rank}(A)$

$$\begin{bmatrix} -0.2 & 4.0 & \dots & -1.3 \\ -4.1 & 0.6 & \dots & -0.2 \\ 2.6 & 6.1 & \dots & 1.4 \\ \vdots & & \ddots & \vdots \\ -1.9 & -1.8 & \dots & 0.3 \end{bmatrix}$$

# Truncating – Latent Semantic Analysis

- Truncating  $U$ ,  $\Sigma$ , and  $V$  to  $k$  dimensions produces best possible  $k$  rank approximation of original matrix
- So truncated,  $U_k$  (or  $V_k^T$ ) is a new low dimensional representation of the word
- Typical values for  $k$  are 100-5000

*100 - 200 is good word representation*



# Words as Context

- Lists how often words appear with other words
  - ▶ In some predefined context (usually a window)
- The obvious problem with raw frequency:  
dominated by common words

	...	the	country	hell	...
...					
state		1973	10	1	
fun		54	2	0	
heaven		55	1	3	
.....					

occur together.

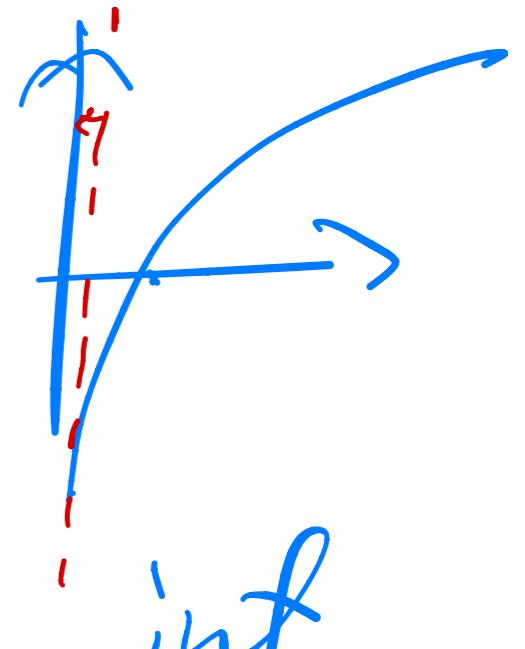
Fix

# Pointwise Mutual Information

- For two events  $x$  and  $y$ , PMI computes the discrepancy between:
    - Their joint distribution
    - Their individual distributions (assuming independence)
- largely penalise common words.*

$$PMI(x, y) = \log_2 \frac{p(x, y)}{p(x)p(y)}$$

$\log 0 \rightarrow -\inf$



# Calculating PMI

	...	the	country	hell	...	$\Sigma$
...		1973	10	1		12786
state		54	2	0		633
fun		55	1	3		627
heaven						
...						
$\Sigma$		1047519	3617	780		15871304

$x = state, y = country$

$$p(x,y) = 10/15871304 = 6.3 \times 10^{-7}$$

$$p(x) = 12786/15871304 = 8.0 \times 10^{-4}$$

$$p(y) = 3617/15871304 = 2.3 \times 10^{-4}$$

$$\begin{aligned} PMI(x,y) &= \log_2(6.3 \times 10^{-7}) / ((8.0 \times 10^{-4}) (2.3 \times 10^{-4})) \\ &= 1.78 \end{aligned}$$

# PMI Matrix

- PMI does a better job of capturing interesting semantics
  - ▶ E.g. *heaven* and *hell*
- But it is obviously biased towards rare words
- And doesn't handle zeros well

	...	the	country	hell	...
...					
state		1.22	1.78	0.63	
fun		0.37	3.79	-inf	
heaven		0.41	2.80	6.60	
.....					

never  
occur  
at same  
time -

# PMI Tricks

- Zero all negative values (PPMI)
  - ▶ Avoid  $-\infty$  and unreliable negative values
- Counter bias towards rare events
  - ▶ Smooth probabilities

Smooth,  $\rightarrow$  n-gram  
language model.  
Same idea

# SVD

*doc as context*

	...	the	country	hell	...
...					
425	0	26.0	6.2		
426	0	5.2	0		
427	0	0	18.6		
...					

*word as context*

	...	the	country	hell	...
...					
	state		1.22	1.78	0.63
	fun		0.37	3.79	0
	heaven		0.41	2.80	6.60
	.....				

*tf-idf matrix*

*document as context*

- Regardless of whether we use document or word as context, SVD can be applied to create dense vectors

*capture broader semantics of word*

*PPMI matrix*

*word as context*

*local semantics of words*

# Similarity

- Word similarity = comparison between word vectors (e.g. cosine similarity)
- Find synonyms, based on proximity in vector space
  - ▶ automatic construction of lexical resources
- Use vectors as features in classifier – more robust to different inputs (*movie* vs *film*)  
          .

# Neural Methods

# Word Embeddings

- We've seen word embeddings used in neural networks (feedforward or recurrent)
- But these models are designed for other tasks:
  - ▶ Classification
  - ▶ Language modelling
- Word embeddings is just one part of the model

# Neural Models for Embeddings

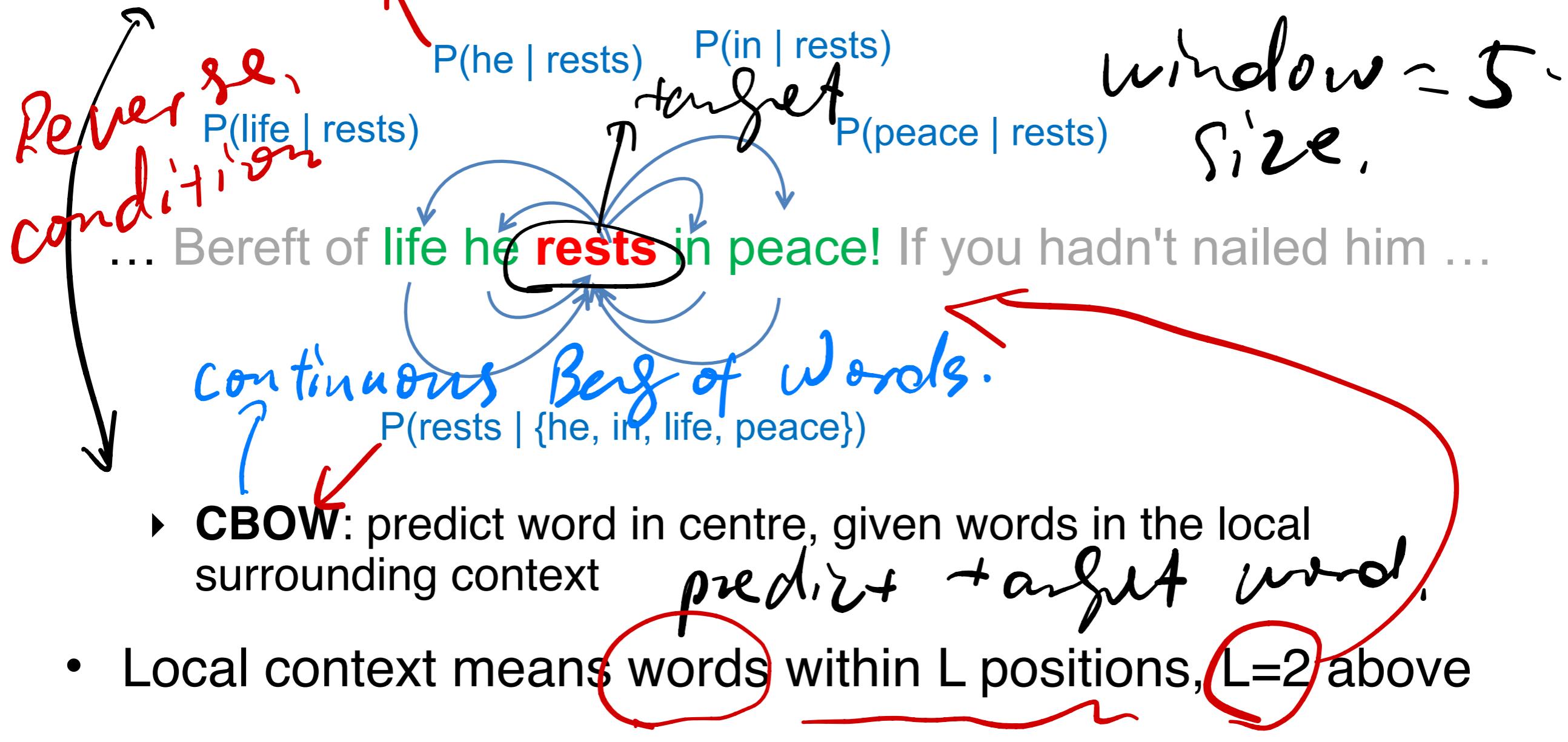
- Can we design neural networks whose goal is to learn word embeddings?
  - Desiderata:
    - ▶ Unsupervised *not labeled*,
    - ▶ Efficient *scale to learn large corpus*,
    - ▶ Useful representation

# Word2Vec

- Neural network inspired approaches seek to learn vector representations of words and their contexts
- Key idea
  - ▶ *Word embeddings should be **similar** to embeddings of **neighbouring** words*
  - ▶ *And **dissimilar** to other words that don't occur nearby*
- Using vector dot product for vector ‘comparison’
  - ▶  $u \cdot v = \sum_j u_j v_j$       *During training time*

# Word2Vec

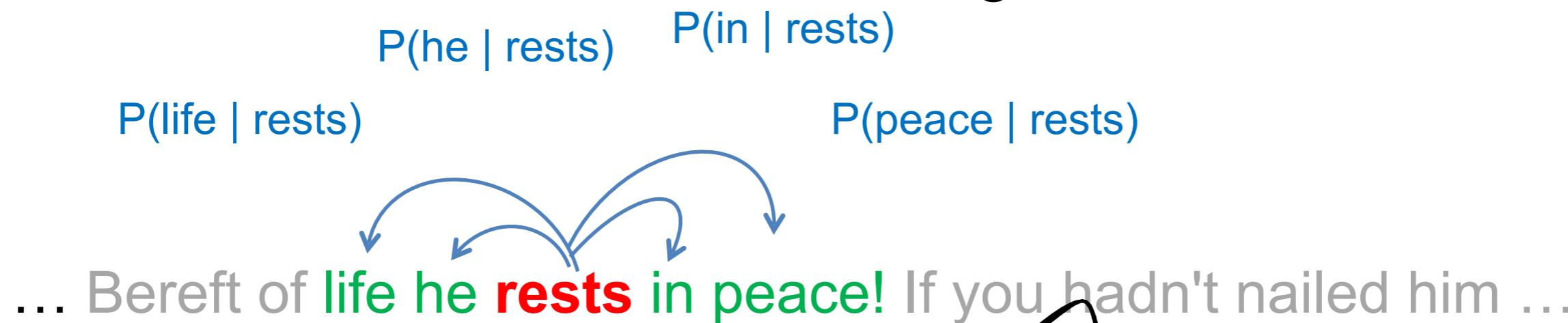
- Framed as learning a classifier
  - ▶ **Skip-gram**: predict words in local context surrounding given word



- **CBOW**: predict word in centre, given words in the local surrounding context
- Local context means words within L positions, L=2 above

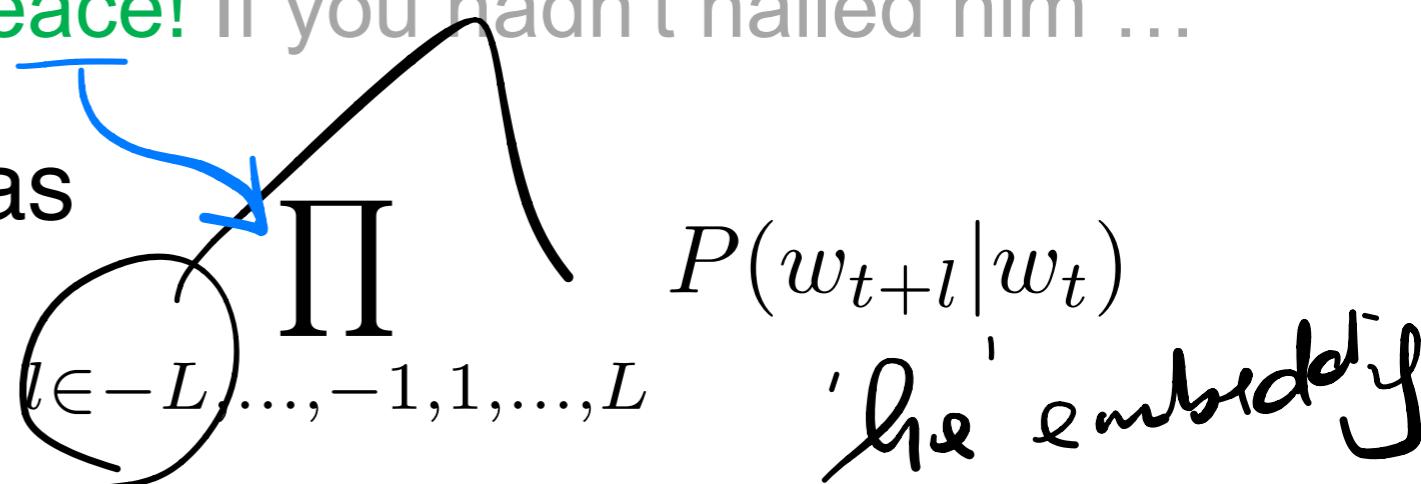
# Skip-gram Model

- Generates each word in context given centre word



- Total probability defined as

Where subscript denotes position in running text



- Using a **logistic regression model**

$$P(w_k | w_j) = \frac{\exp(c_{w_k} \cdot v_{w_j})}{\sum_{w' \in V} \exp(c_{w'} \cdot v_{w_j})}$$

all vocabulary

# Embedding parameterisation

- Two parameter matrices, with  $d$ -dimensional embedding for all words

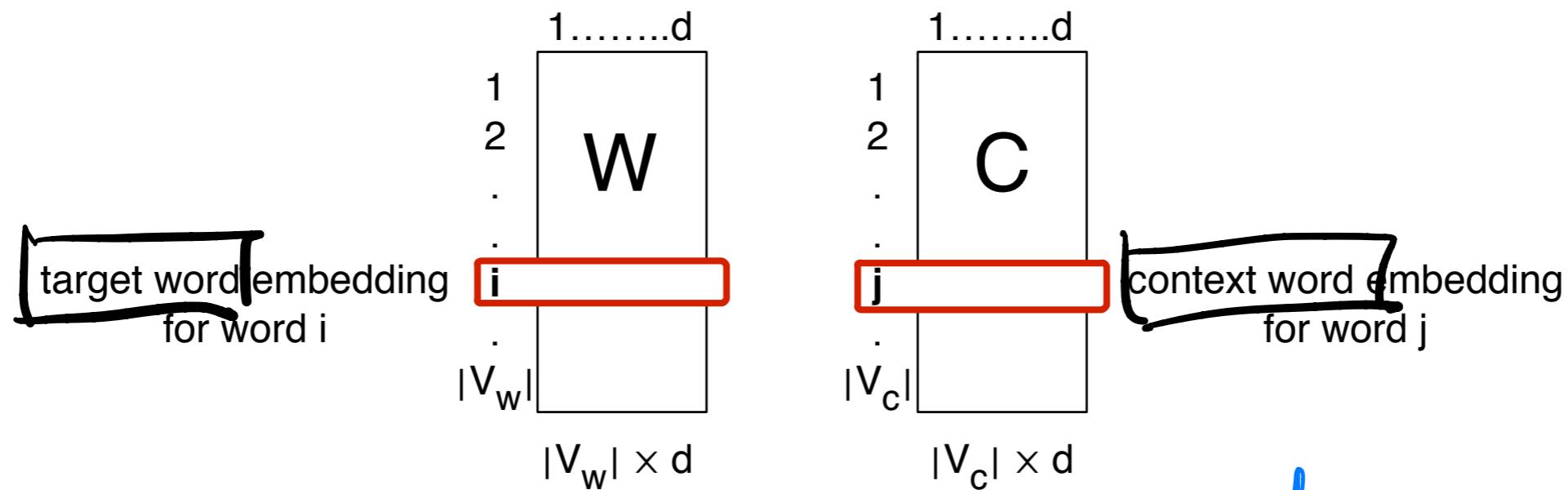


Fig 19.17, JM3

*target word parameters*

- Words are numbered, e.g., by sorting vocabulary and using word location as its index

# Skip-gram model

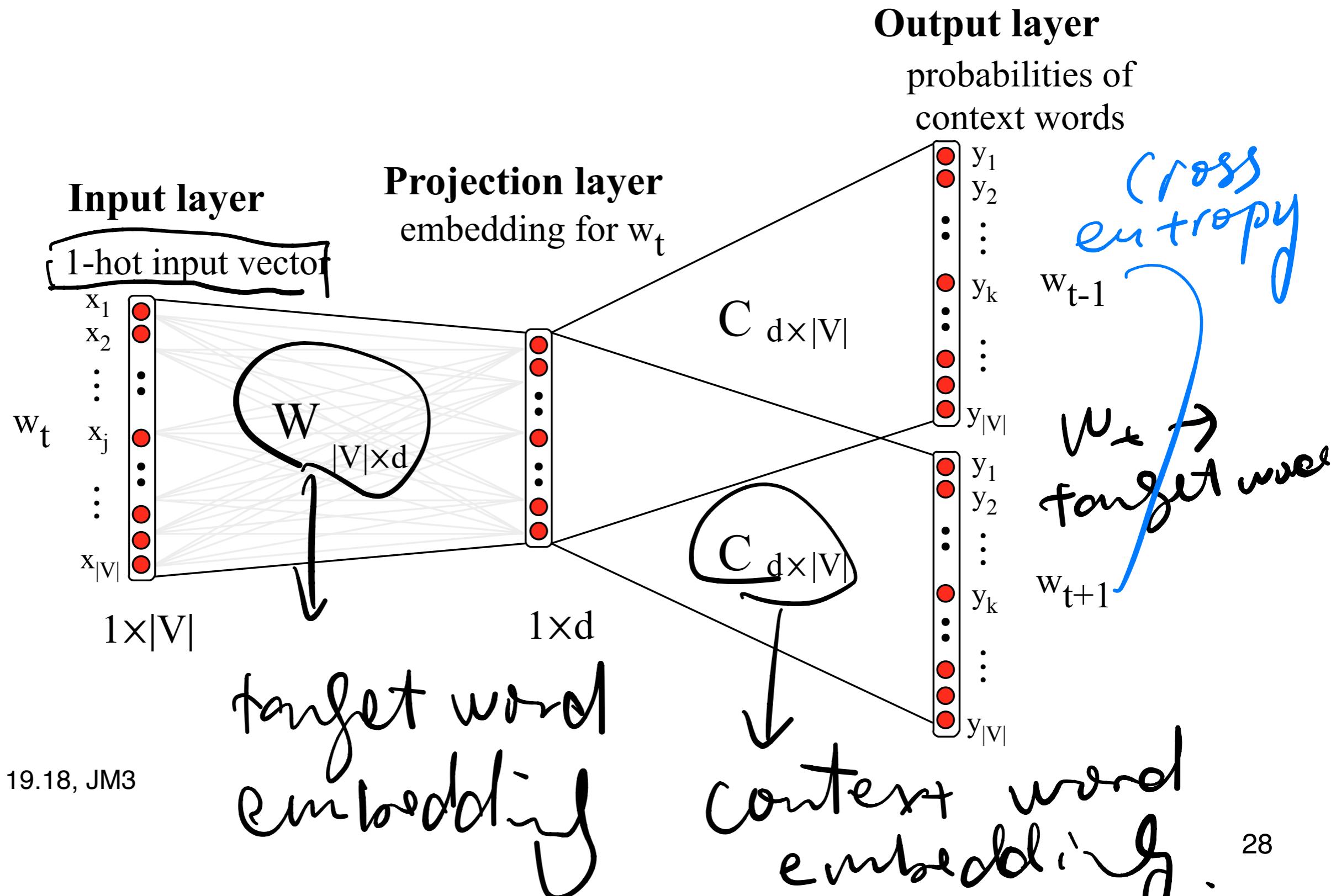


Fig 19.18, JM3

# Training the skip-gram model

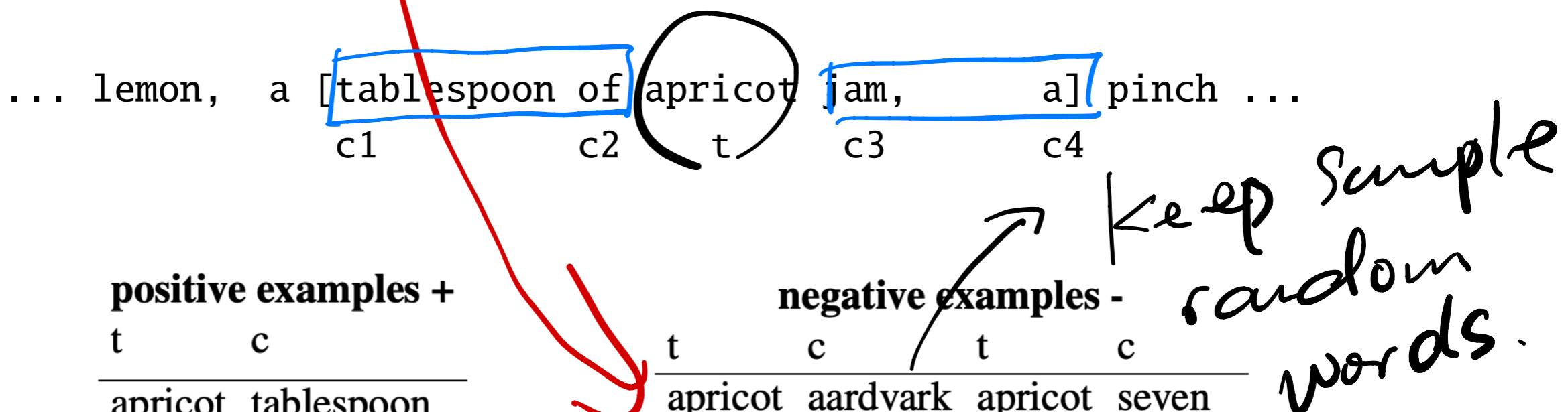
- Train to maximise likelihood of **raw text**
- Too slow in practice, due to normalisation over  $|V|$

$$P(w_k | w_j) = \frac{\exp(c_{w_k} \cdot v_{w_j})}{\sum_{w' \in V} \exp(c_{w'} \cdot v_{w_j})}$$

*softmax.*      *slow.*

- Reduce problem to binary classification, distinguish real context words from non-context words aka “**negative samples**”
  - ▶ words drawn randomly from  $V$       all words.

# Negative Sampling



positive examples +  
t      c  
-----  
apricot tablespoon  
apricot of  
apricot jam  
apricot a

negative examples -  
t      c  
-----  
apricot aardvark  
apricot my  
apricot forever  
apricot dear  
apricot if

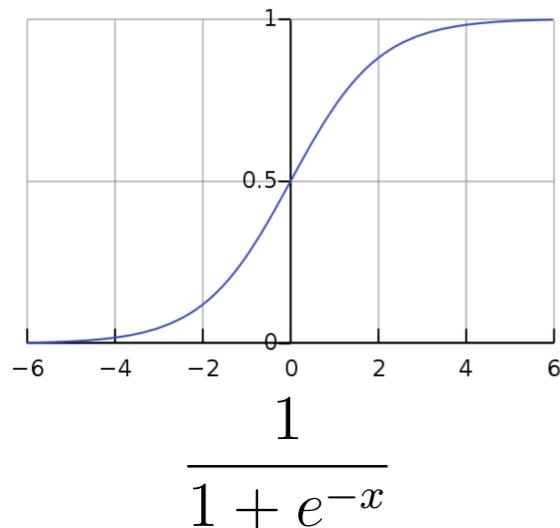
*target word E contexts W. Emb.*

$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c / \tau}}$$

maximise similarity between target word and real context words

$$P(-|t, c) = 1 - \frac{1}{1 + e^{-t \cdot c}}$$

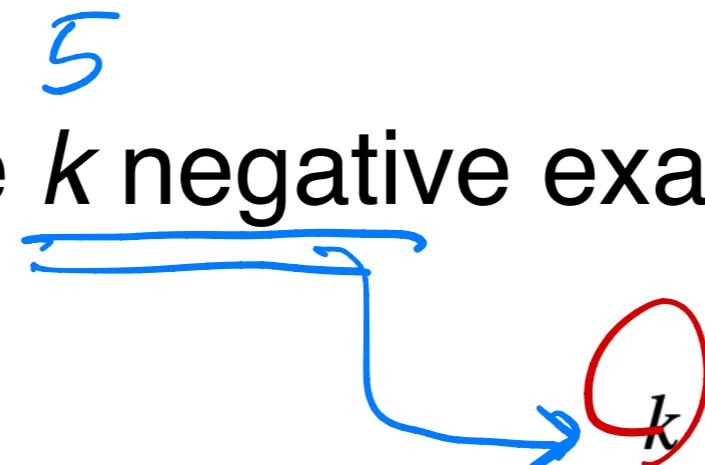
minimise similarity between target word and non-context words



# Skip-gram Loss

$$L(\theta) = \sum_{(t,c) \in +} \log P(+|t, c) + \sum_{(t,c) \in -} \log P(-|t, c)$$

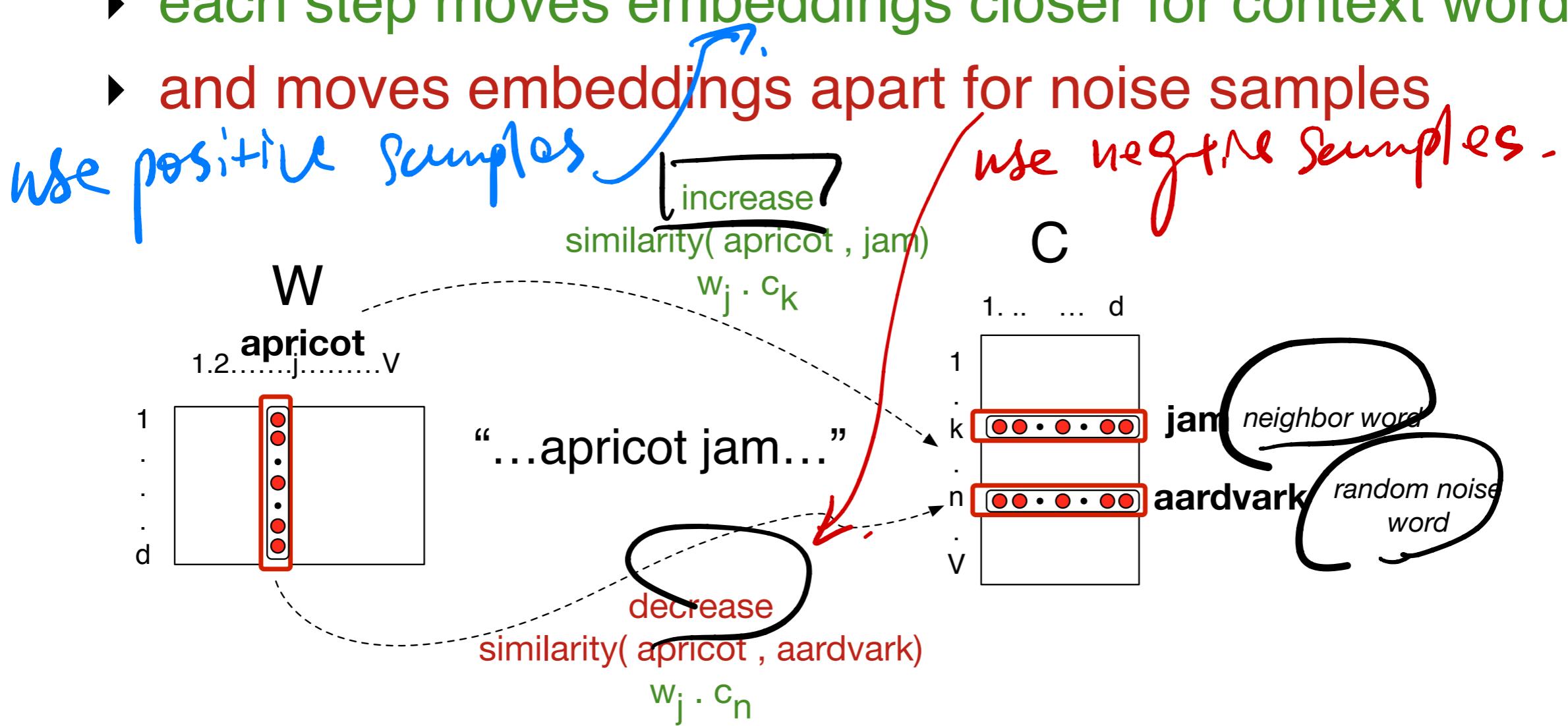
- In practice, use  $k$  negative examples



$$L(\theta) = \log P(+|t, c) + \sum_{i=1}^k \log P(-|t, n_i)$$

# Training Illustration

- Iterative process (stochastic gradient descent)
  - each step moves embeddings closer for context words
  - and moves embeddings apart for noise samples



# Desiderata

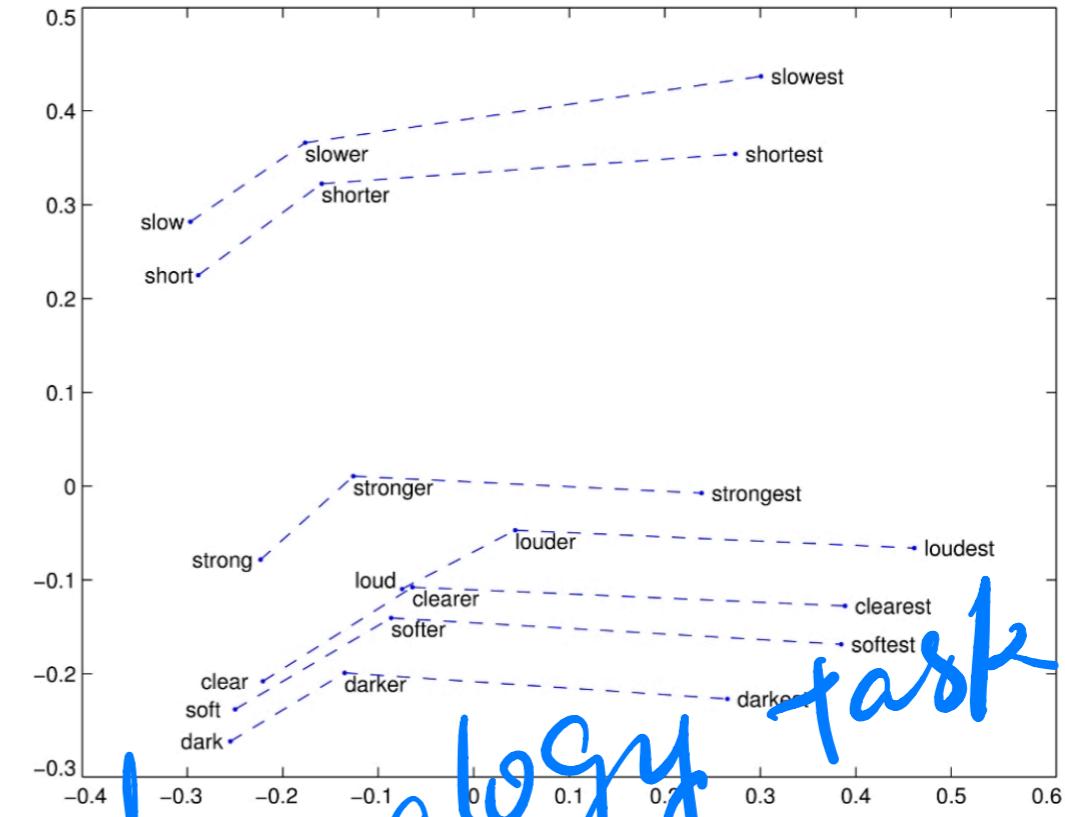
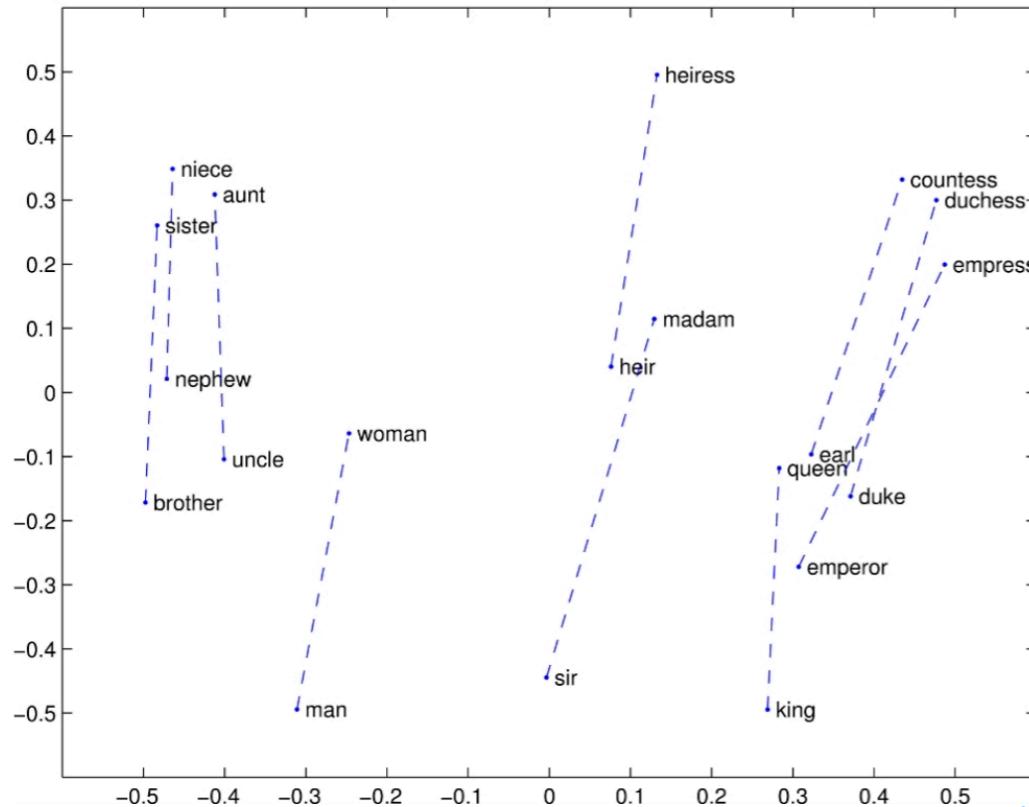
*Features*

- Unsupervised
  - ▶ Raw, unlabelled corpus
- Efficient
  - ▶ Negative sampling (~~avoid softmax over full vocabulary~~)
  - ▶ Scales to very very large corpus *pre-trained*
- Useful representation:
  - ▶ How do we evaluate word vectors?

# Evaluating Word Vectors

- Lexicon style tasks
    - ▶ *WordSim-353* are pairs of nouns with judged relatedness
    - ▶ *SimLex-999* also covers verbs and adjectives
    - ▶ *TOEFL* asks for closest synonym as multiple choice
    - ▶ ...
  - Test compatibility of word pairs using cosine similarity in vector space
- compares to human judgement.*

# Embeddings Exhibit Meaningful Geometry



Word analogy task

- Word analogy task

- ▶ *Man* is to *King* as *Woman* is to ???
- ▶  $v(Man) - v(King) = v(Woman) - v(\text{???})$
- ▶  $v(\text{??}) = \boxed{v(Woman) - v(Man) + v(King)}$



closest  
vector.  
build model.

# Evaluating Word Vectors

*test in the task*

- Best evaluation is in other **downstream tasks**
  - ▶ Use bag-of-word embeddings as a feature representation in a classifier
  - ▶ First layer of most deep learning models is to embed input text; **use pre-trained word vectors as embeddings** *not randomly initialize all parameters*
- Recently **contextual word vectors** shown to work even better
  - ▶ ELMO & **BERT** (next lecture!)



# Pointers to Software

- Word2Vec
  - ▶ C implementation of Skip-gram and CBOW  
<https://code.google.com/archive/p/word2vec/>
- GenSim
  - ▶ Python library with many methods include LSI, topic models and Skip-gram/CBOW  
<https://radimrehurek.com/gensim/index.html>
- GLOVE
  - ▶ <http://nlp.stanford.edu/projects/glove/>

# Further Reading

- ▶ JM3, Ch 6