

# IN4080 – 2022 FALL

## NATURAL LANGUAGE PROCESSING

Jan Tore Lønning

# Vectors, Distributions, Embeddings

Lecture 7, Oct 6

# Today

3

- Lexical semantics
- Vector models of documents
- tf-idf weighting
- Word-context matrices
- Word embeddings with dense vectors
- Word2vec

# The meaning of words

4

- Words (lecture 2)
  - Type – token
  - Word – lexeme – lemma
  - Meaning?

# Look into the dictionary

5

lemma sense

definition

pepper, *n.*

Pronunciation: Brit. /ˈpeɪpə/, U.S. /ˈpeɪpər/

Forms: OE *peopor* (*rare*), OE *pipcor* (*transmission error*), OE *pipor*, OE *pipur* (*rare*)  
Frequency (in current use):

Etymology: A borrowing from Latin. Etymon: Latin *pipēr*.  
< classical Latin *pipēr*, a loanword < Indo-Aryan (as is ancient Greek *πῖπερι*); compare Sai

1. The spice or the plant.

1. a. A hot pungent spice derived from the prepared fruits (peppercorns) of the pepper plant, *Piper nigrum* (see sense 2a), used from early times to season food, either whole or ground to powder (often in association with salt). Also (locally, chiefly with distinguishing word): a similar spice derived from the fruits of certain other species of the genus *Piper*; the fruits themselves.

The ground spice from *Piper nigrum* comes in two forms, the more pungent *black pepper*, produced from black peppercorns, and the milder *white pepper*, produced from white peppercorns: see BLACK *adj.* and *n.* Special uses 5a, PEPPERCORN *n.* 1a, and WHITE *adj.* and *n.* Special uses 7b(a).

2. a. The plant *Piper nigrum* (family Piperaceae), a climbing shrub indigenous to South Asia and also cultivated elsewhere in the tropics, which has alternate stalked entire leaves, with pendulous spikes of small green flowers opposite the leaves, succeeded by small berries turning red when ripe. Also more widely: any plant of the genus *Piper* or the family Piperaceae.

b. Usu. with distinguishing word: any of numerous plants of other families having hot pungent fruits or leaves which resemble pepper (1a) in taste and in some cases are used as a substitute for it.

c. U.S. The California pepper tree, *Schinus molle*. Cf. PEPPER TREE *n.* 3.

3. Any of various forms of capsicum, esp. *Capsicum annuum* var. *annuum*. Originally (chiefly with distinguishing word): any variety of the *C. annuum* Longum group, with elongated fruits having a hot, pungent taste, the source of cayenne, chilli powder, paprika, etc., or of the

- A word with several senses is called **polysemous**
- If two different words look and sound the same, they are called **homonyms**

- How to tell: one word or several?
  - Common origin
  - But not waterproof/easy to see











# Relations between senses

10

Term	Definition	Examples		
Synonymy	Have the same meaning in all(?) / some(?) contexts	<i>sofa-couch, bus-coach</i> <i>big-large</i>		
Antonym	Opposites with respect to a feature of meaning	<i>true-false, strong-weak, up-down</i>		
Hyponym-hypernym	The <hyponym> is a type-of the <hypernym>	<i>rose → flower, cow → animal,</i> <i>car → vehicle</i>		
Similarity		<i>cow-horse</i> <i>boy-girl</i>		

# Relations between senses

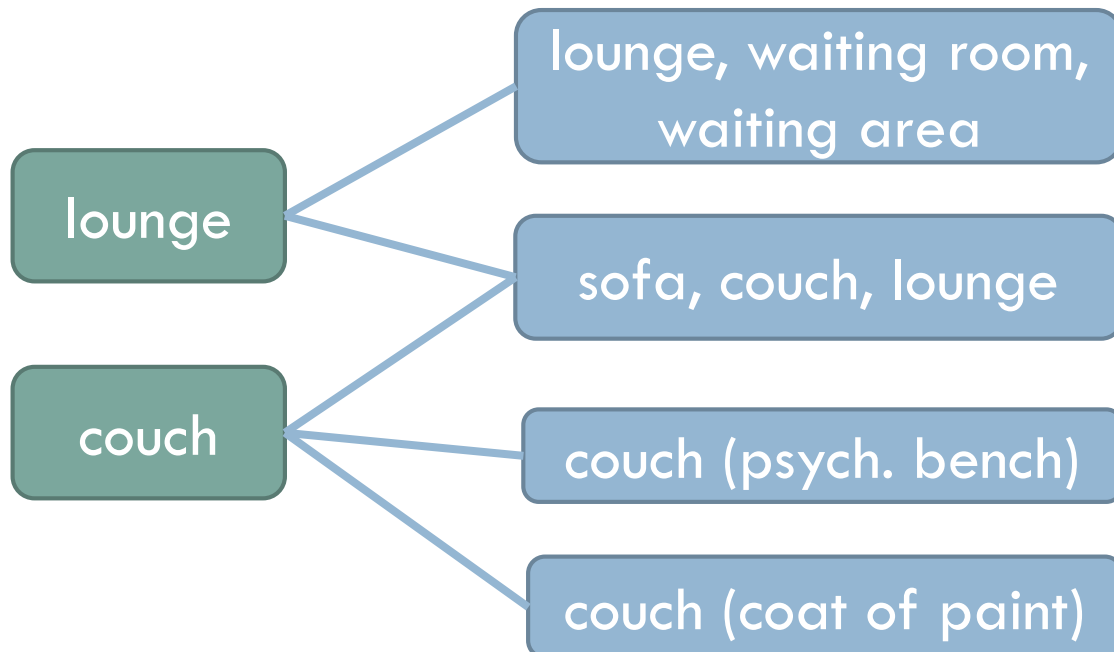
11

Term	Definition	Examples		
Synonymy	Have the same meaning in all(?) / some(?) contexts	<i>sofa-couch, bus-coach</i> <i>big-large</i>		
Antonym	Opposites with respect to a feature of meaning	<i>true-false, strong-weak, up-down</i>		
Hyponym-hypernym	The <hyponym> is a type-of the <hypernym>	<i>rose → flower, cow → animal,</i> <i>car → vehicle</i>		
Similarity		<i>cow-horse</i> <i>boy-girl</i>		
Related		<i>money-bank</i> <i>fish-water</i>		

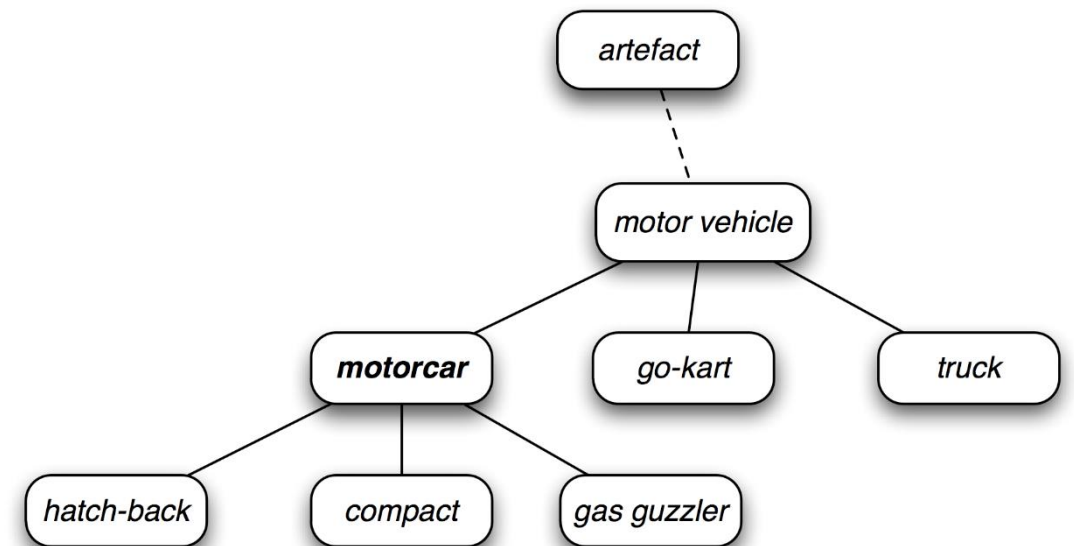
# Resources for lexical semantics: WordNet

12

- <https://wordnet.princeton.edu>
- To each word:
  - ▣ One or more synsets (synonymy)



- Hyponymy relations between the synsets



# What does ongchoi mean?

13

- Suppose you see these sentences:
  - ▣ *Ong choi is delicious sautéed with garlic.*
  - ▣ *Ong choi is superb over rice*
  - ▣ *Ong choi leaves with salty sauces*
- And you've also seen these:
  - ▣ *...spinach sautéed with garlic over rice*
  - ▣ *Chard stems and leaves are delicious*
  - ▣ *Collard greens and other salty leafy greens*
- Conclusion: Ongchoi is a leafy green like spinach, chard, or collard greens



# Similar

14

Related

(first-order association,  
syntagmatic)

ong choy

*delicious*

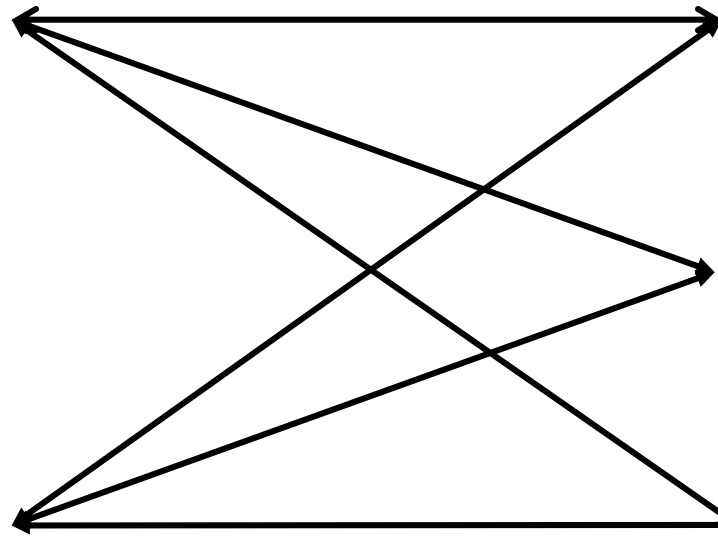
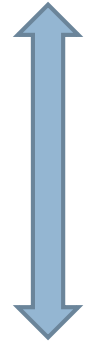
*sautéed with garlic*

spinach

*over rice*

Similar

(second-order  
association,  
paradigmatic)



# The distributional hypothesis

15

- Words that occur in similar contexts have similar meanings

# Today

16

- Lexical semantics
- **Vector models of documents**
- tf-idf weighting
- Word-context matrices
- Word embeddings with dense vectors
- Word2vec



# Shakespeare (from J & M)

17

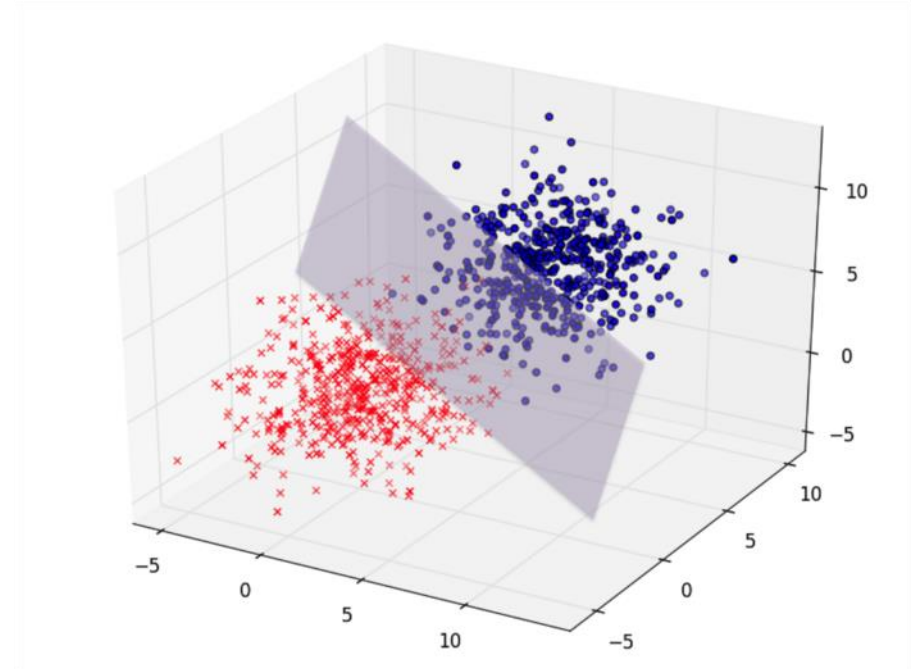
	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

- Vectors are similar for the two comedies
- Different than the historical dramas
- Comedies have more **fools** and **wit** and fewer **battles**.
- Notice similarity to text classification
- Mandatory 1B, multinomial
  - ▣ The document represented by a vector with the occurrences of 35,000 terms

# Document classification

18

- The word vectors were used as basis for classification
- If two documents had the same vectors they were put in the same class
- Documents are similar = on the same side of the separating hyperplane

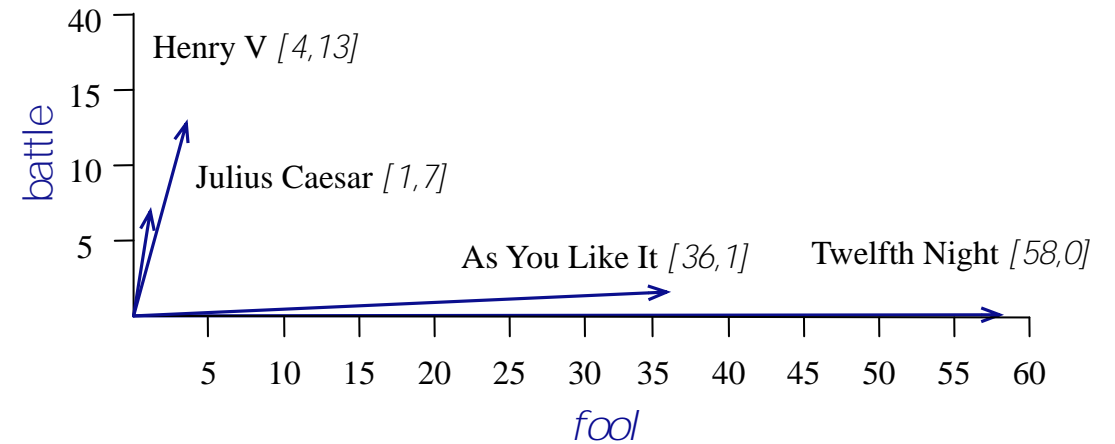


A problem to draw 35,000 dimensions

# Information retrieval (IR)

19

- Documents placed in the same  $n$ -dimensional space as in classification
- Retrieve documents similar to a given document

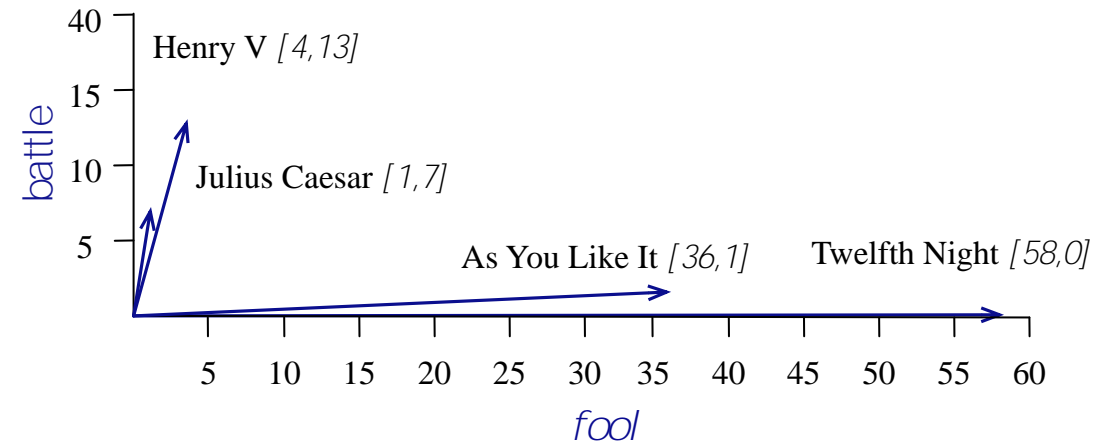


# Cosine similarity

20

- Several possible ways to define similarity, e.g.,
  - ▣ Euclidean
  - ▣ Manhattan
- Most common: cosine
  - ▣ Do the arrows point in the same direction?

$$\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\vec{v}}{|\vec{v}|} \cdot \frac{\vec{w}}{|\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$



# Let us try: $\cos(v_1, v_2)$

21

## Full vectors

	AYLI	TwNi	JuCa	HenV
AYLI	1.000	0.950	0.945	0.949
TwNi	0.950	1.000	0.809	0.822
JuCa	0.945	0.809	1.000	0.999
HenV	0.949	0.822	0.999	1.000

## *battles & fools*

	AYLI	TwNi	JuCa	HenV
AYLI	1.000	1.000	0.169	0.321
TwNi	1.000	1.000	0.141	0.294
JuCa	0.169	0.141	1.000	0.988
HenV	0.321	0.294	0.988	1.000

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

# Today

22

- Lexical semantics
- Vector models of documents
- **tf-idf weighting**
- Word-context matrices
- Word embeddings with dense vectors
- Word2vec

# Inverse document frequency

23

- Are all words equally important?
- Intuition: A word occurring in a large proportion of documents is not a good discriminator.
- $idf_t = \log \frac{N}{df_t}$ 
  - $df_t$  the number of documents containing  $t$ .

	Example: 1,000,000 documents, log_10							
$df_t$	1,000,000	100,000	10,000	1,000	100	10	1	0
$idf_t$	0	1	2	3	4	5	6	?

# *tf-idf*

24

- Tf-idf weighting
- $tf_{t,d} \times idf_t$ 
  - ▣  $tf_{t,d}$  is the frequency of the term  $t$  in the document  $d$



# Variants of $tf-idf$

25

- $idf_t = \log \frac{N}{df_t}$
- scikit-learn:
  - $\log \frac{N}{df_t} + 1$ 
    - `TfidfTransformer(smooth_idf=False)`
  - $idf_t = \log \frac{N+1}{df_t+1} + 1$ 
    - `TfidfTransformer(smooth_idf=True)`
    - Here  $\log_{10}$  vs  $\log_e$  matters
- and others in the literature
- Avoids  $\log(0)$

## $tf_{t,d}$ - Alternatives

- Linear:
  - raw counts
  - relative frequencies (L1)
  - length normalization (L2)
  - all yield same cosine-similarity
- Sublinear
  - (J&M):  $tf_{t,d} = \log(count(t, d) + 1)$
  - (sklearn)  $tf_{t,d} = \log(count(t, d)) + 1$ ,  
0 when  $count(t, d) = 0$ 
    - `TfidfTransformer(sub_linear=True)`
  - Similar effect as log in  $idf_t$

# The effect of tf-idf

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022

**Figure 6.8** A tf-idf weighted term-document matrix for four words in four Shakespeare plays, using the counts in Fig. 6.2. For example the 0.049 value for *wit* in *As You Like It* is the product of  $tf = \log_{10}(20 + 1) = 1.322$  and  $idf = .037$ . Note that the idf weighting has eliminated the importance of the ubiquitous word *good* and vastly reduced the impact of the almost-ubiquitous word *fool*.

# The effect of tf-idf

27

	As You Like It	Twelfth Night	Julius Caesar	Henry V
<b>battle</b>	0.074	0	0.22	0.28
<b>good</b>	0	0	0	0
<b>fool</b>	0.019	0.021	0.0036	0.0083
<b>wit</b>	0.049	0.044	0.018	0.022

	AYLI	TwNi	JuCa	HenV
AYLI	1.000	0.577	0.860	0.861
TwNi	0.577	1.000	0.081	0.083
JuCa	0.860	0.081	1.000	1.000
HenV	0.861	0.083	1.000	1.000

- Maybe not the expected effect
- Probably because of few terms.

# Today

28

- Lexical semantics
- Vector models of documents
- tf-idf weighting
- **Word-context matrices**
- Word embeddings with dense vectors
- Word2vec

# Vector repr. of words 1: A vector of documents

	As You Like It	Twelfth Night	Julius Caesar	Henry V
<b>battle</b>	1	0	7	13
<b>good</b>	114	80	62	89
<b>fool</b>	36	58	1	4
<b>wit</b>	20	15	2	3

*battle* is "the kind of word that occurs in Julius Caesar and Henry V"

*fool* is "the kind of word that occurs in comedies, especially Twelfth Night"

# Vector repr. of words 2: Word-context matrix

30

- Two **words** are similar in meaning if their context vectors are similar

sugar, a sliced lemon, a tablespoonful of their enjoyment. Cautiously she sampled her first well suited to programming on the digital for the purpose of gathering data and **apricot pineapple computer. information** jam, a pinch each of, and another fruit whose taste she likened In finding the optimal R-stage policy from necessary for the study authorized in the

	aardvark	computer	data	pinch	result	sugar	...
apricot	0	0	0	1	0	1	
pineapple	0	0	0	1	0	1	
digital	0	2	1	0	1	0	
information	0	1	6	0	4	0	

# Word-context matrix

31

## Document-term matrix

- Objects: a set of documents,  $D$
- Features: a set of terms,
  - ▣  $T = \{t_1, t_2, \dots, t_n\}$
- Each document  $d$  is identified with a vector
  - ▣  $(v_1, v_2, \dots, v_n)$
  - ▣ where  $v_i$  is calculated from the frequency of  $t_i$  in  $d$ .

## Word-context matrix

- Objects: a vocabulary of words,  $V$
- Features: a set of words,
  - ▣  $C = \{c_1, c_2, \dots, c_n\}$
- A set of texts,  $T$
- A definition of the context of an occurrence of  $w$  in  $T$
- Each word  $w$  in  $V$  is identified with a vector
  - ▣  $(v_1, v_2, \dots, v_n)$
  - ▣ where  $v_i$  is calculated from the frequency of  $c_i$  in all the contexts of  $w$  in  $T$

# Word-context matrix

32

## Comments

- $C=V$ , or  $C$  is a smaller set of the most frequent terms
  - ▣ To avoid too large repr.
- Context, alternatives:
  - ▣ A sentence
  - ▣ A window of  $k$  tokens on each side
  - ▣ A document
  - ▣ Defined by grammatical relations (after parsing)

## Word-context matrix

- Objects: a vocabulary of words,  $V$
- Features: a set of words,
  - ▣  $C = \{c_1, c_2, \dots, c_n\}$
- A set of texts,  $T$
- A definition of the context of an occurrence of  $w$  in  $T$
- Each word  $w$  in  $V$  is identified with a vector
  - ▣  $(v_1, v_2, \dots, v_n)$
  - ▣ where  $v_i$  is calculated from the frequency of  $c_i$  in all the contexts of  $w$  in  $T$



# So-far

33

- A word  $w$  can be represented by a context vector  $v_w$  where position  $j$  in the vector reflects the frequency of occurrences of  $w_j$  with  $w$ .
- Can be used for
  - ▣ studying similarities between words.
  - ▣ document similarities
- But the vectors are *sparse*
  - ▣ Long: 20-50,000
  - ▣ Many entries are 0
- Even though *car* and *automobile* get similar vectors, because both co-occur with e.g., *drive*, in the vector for *drive* there is no connection between the *car* element and the *automobile* element.

# Today

34

- Lexical semantics
- Vector models of documents
- tf-idf weighting
- Word-context matrices
- **Word embeddings with dense vectors**
- Word2vec

# Dense vectors

35

## How?

- Shorter vectors.
  - ▣ (length 50-1000)
  - ▣ “low-dimensional” space
- Dense (most elements are not 0)
- Intuitions:
  - ▣ Similar words should have similar vectors.
  - ▣ Words that occur in similar contexts should be similar.

## Properties

- Generalize better than sparse vectors.
- Input for deep learning
  - ▣ Fewer weights (or other weights)
- Capture semantic similarities better.
- Better for sequence modelling:
  - ▣ Language models, etc.

# Word embeddings

36

- Each word is represented as a vector of reals
- Words are more or less similar
- A word can be similar to one word in some dimensions and other words in other dimensions



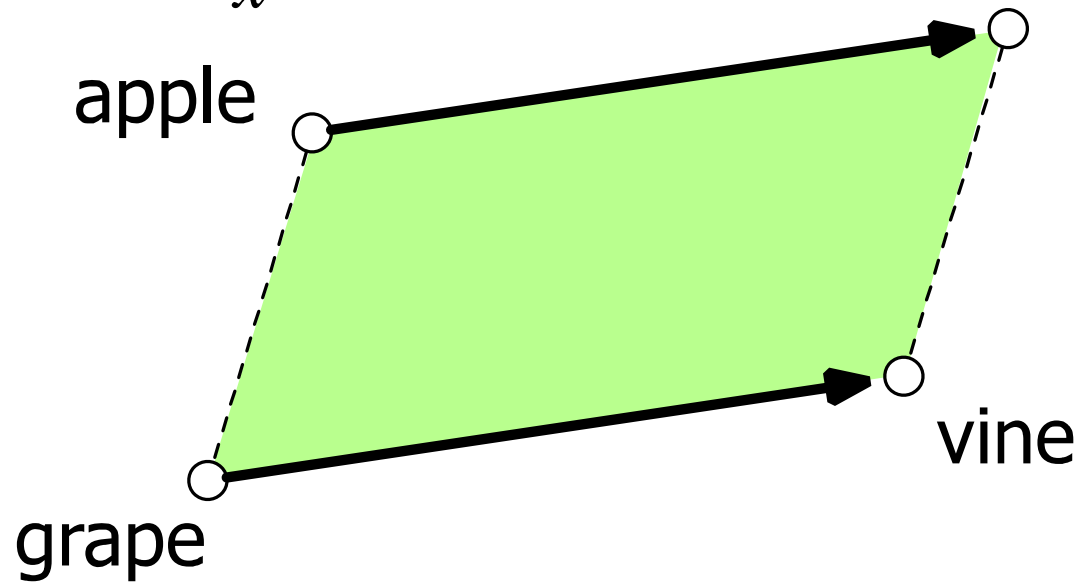
Figure from

<https://medium.com/@jayeshbahire>

# Analogical relations

- The classic parallelogram model of analogical reasoning (Rumelhart and Abrahamson 1973)
- For a problem  $a:a^*::b:b^*$ , the parallelogram method is:

$$\hat{b}^* = \operatorname{argmax}_x \operatorname{distance}(x, a^* - a + b)$$

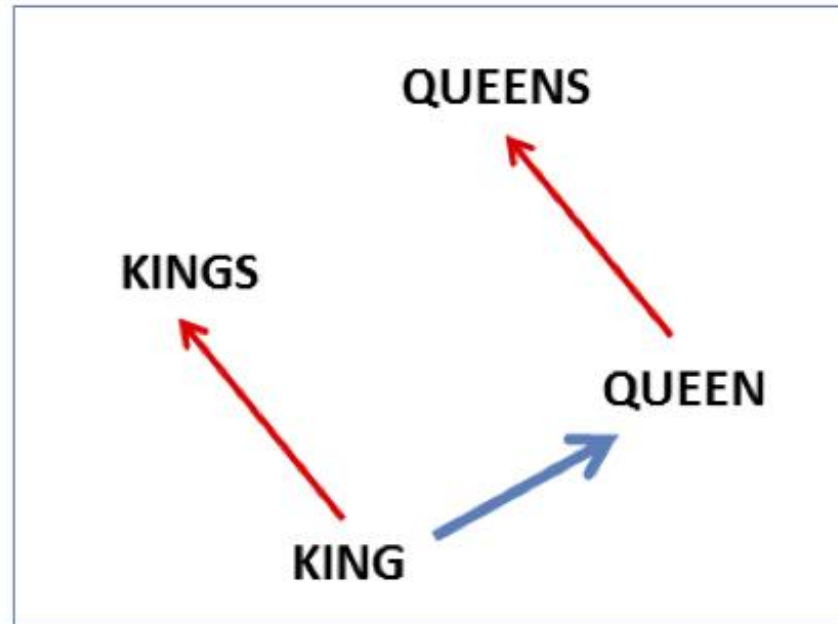
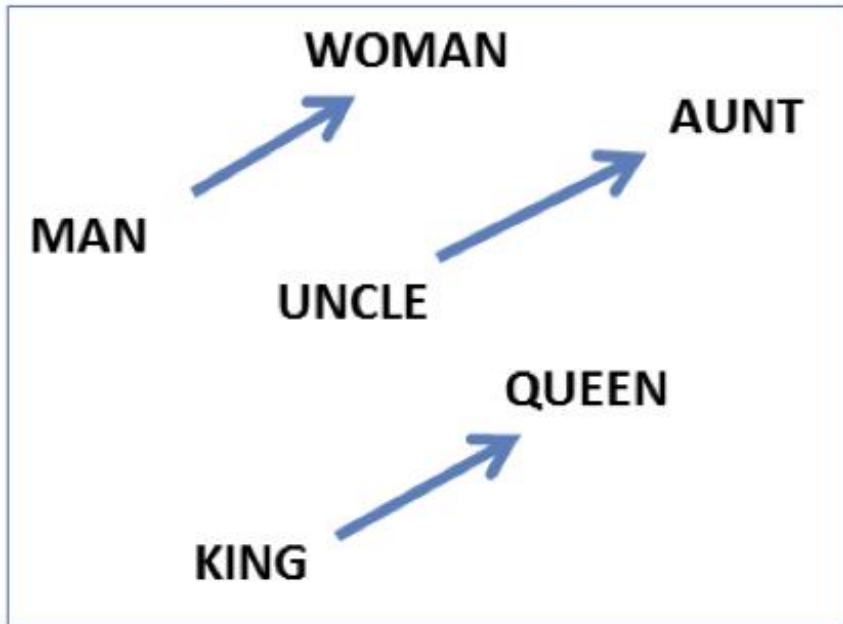


# Analogy: Embeddings capture relational meaning!

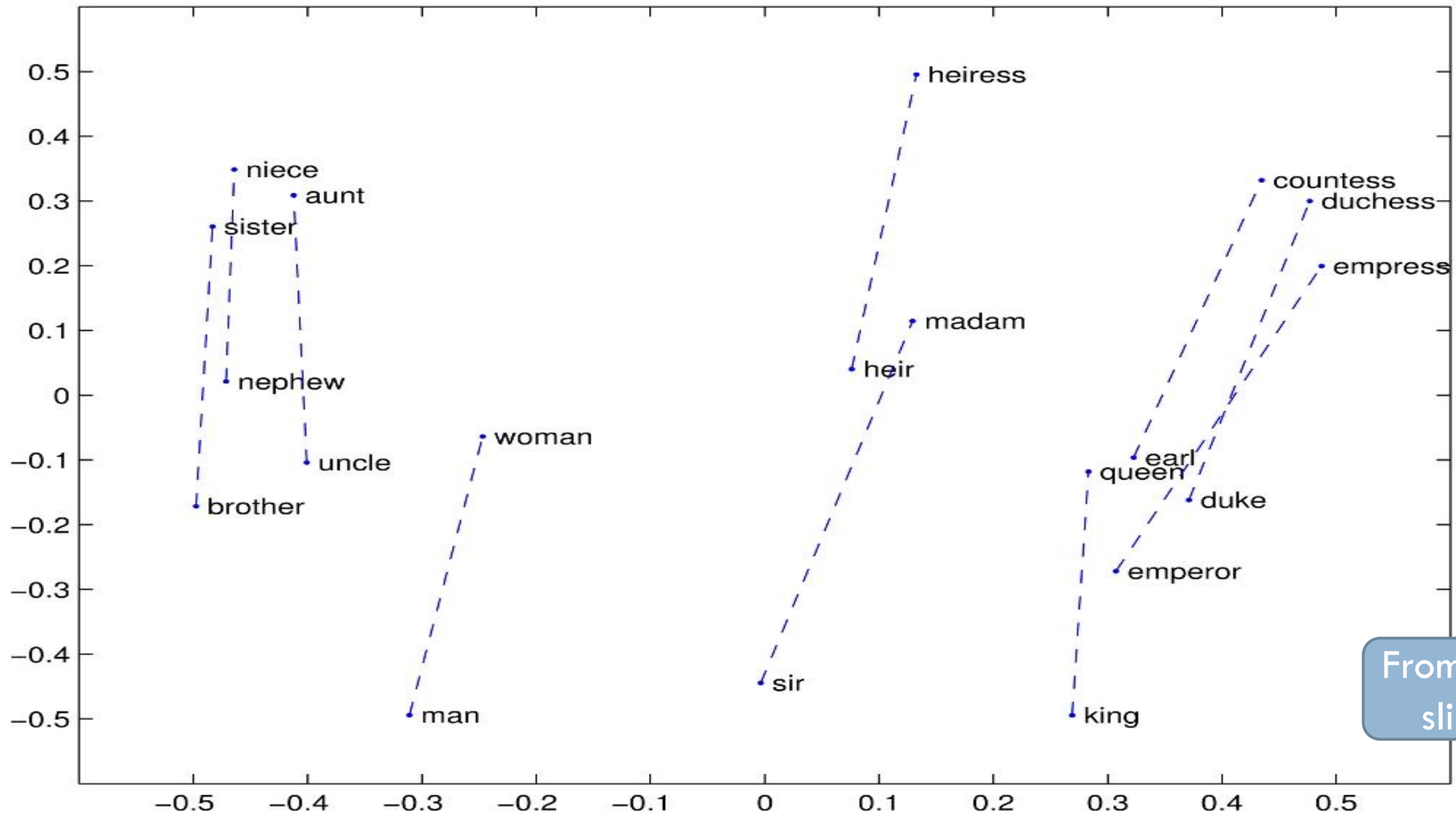
38

$\text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'}) \approx \text{vector}(\text{'queen'})$

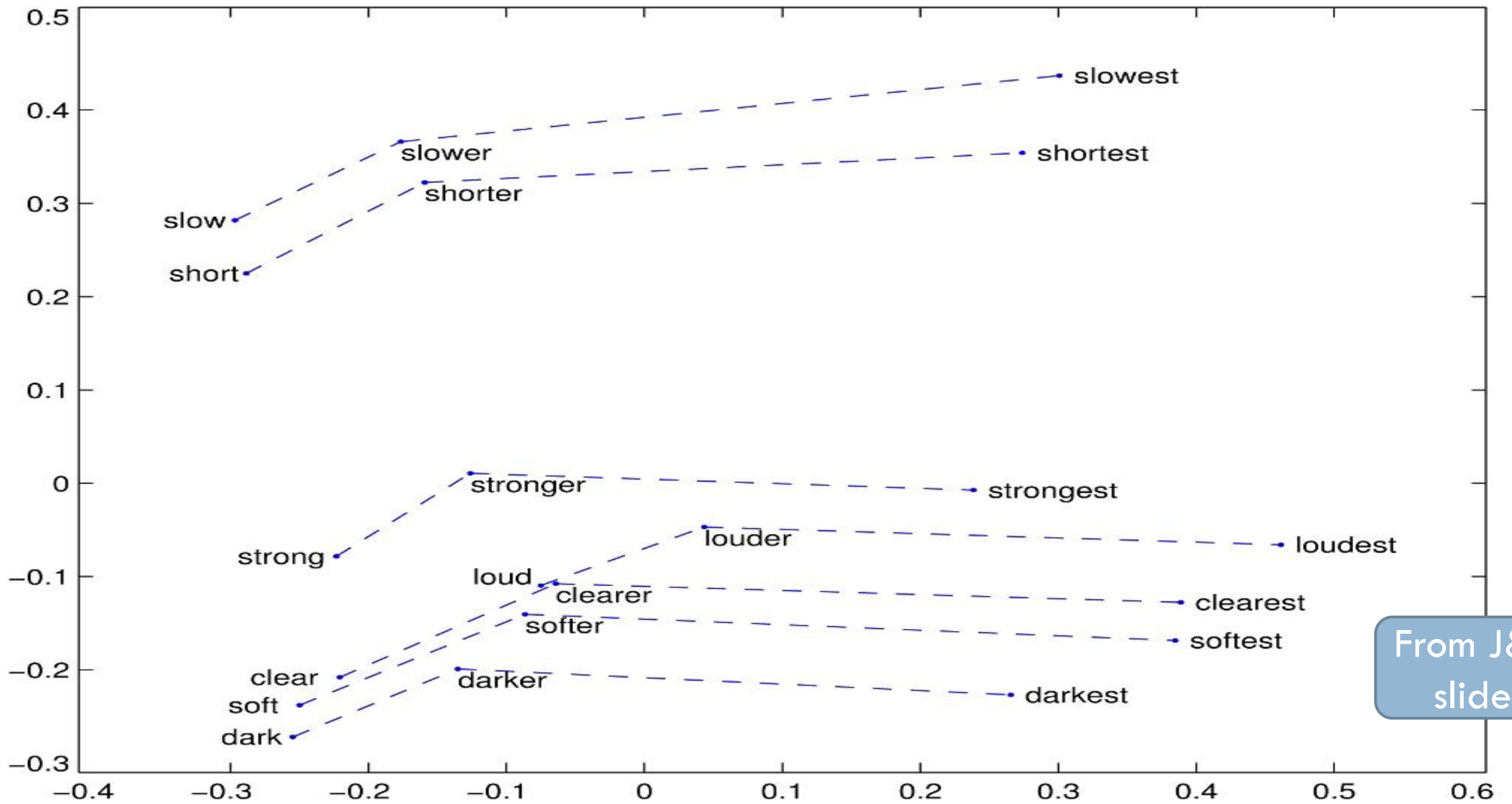
$\text{vector}(\text{'Paris'}) - \text{vector}(\text{'France'}) + \text{vector}(\text{'Italy'}) \approx \text{vector}(\text{'Rome'})$



From J&M  
slides



From J&M  
slides



From J&M  
slides



# Demo

41

- <http://vectors.nlpl.eu/explore/embeddings/en/>

# Track change of meaning of words

42



~30 million books, 1850-1990, Google Books data

From J&M  
slides

# Bias

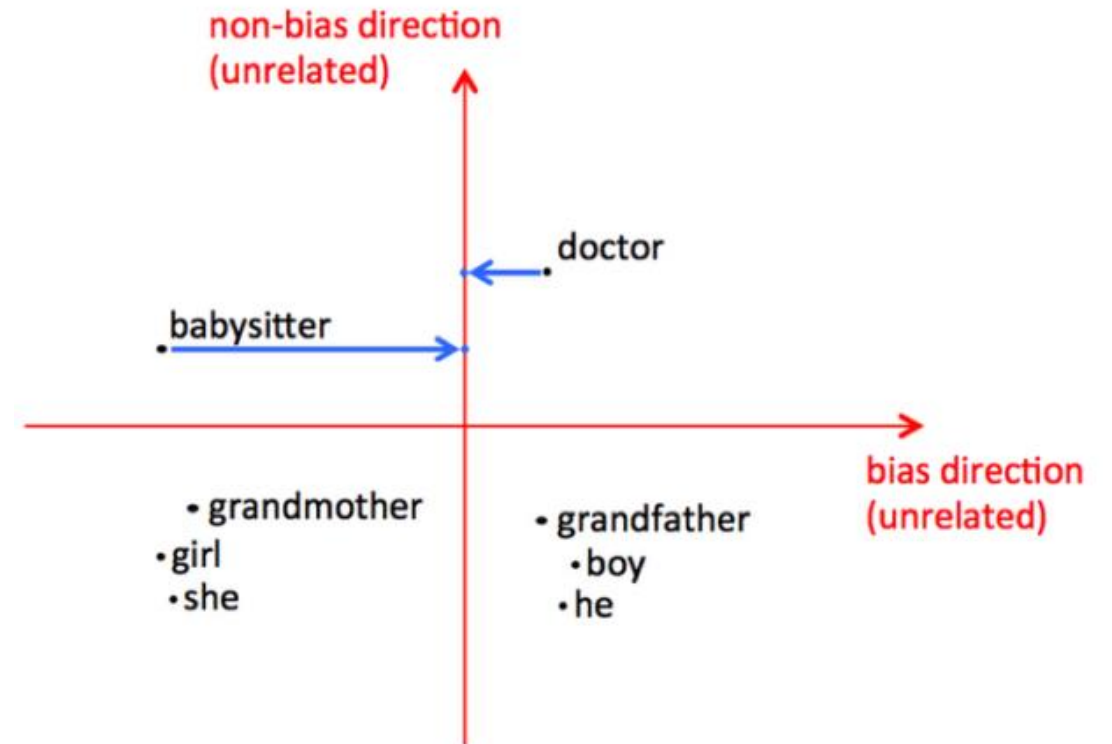
43

- *Man is to computer programmer as woman is to homemaker.*
- Different adjectives associated with:
  - ▣ male and female terms
  - ▣ typical black names and typical white names
- Embeddings may be used to study historical bias

# Debiasing (research topic)

44

- Goal: neutralize the biases
- Some positive results
- But also reports that is is not fully possible
  
- Is debiasing a goal?
- When should we (not) debias?



# Evaluation of embeddings

- Extrinsic evaluation:
  - ▣ Evaluate contribution as part of an application
- Intrinsic evaluation:
  - ▣ Evaluate against a resource
- Some datasets
  - ▣ WordSim-353:
    - Broader "semantic relatedness"
  - ▣ SimLex-999:
    - Narrower: similarity
    - Manually annotated for similarity

Word1	Word2	POS	Sim-score
old	new	A	1.58
smart	intelligent	A	9.2
plane	jet	N	8.1
woman	man	N	3.33
word	dictionary	N	3.68
create	build	V	8.48
get	put	V	1.98
keep	protect	V	5.4

Part of SimLex-999

# Today

46

- Lexical semantics
- Vector models of documents
- tf-idf weighting
- Word-context matrices
- Word embeddings with dense vectors
- **Word2vec**

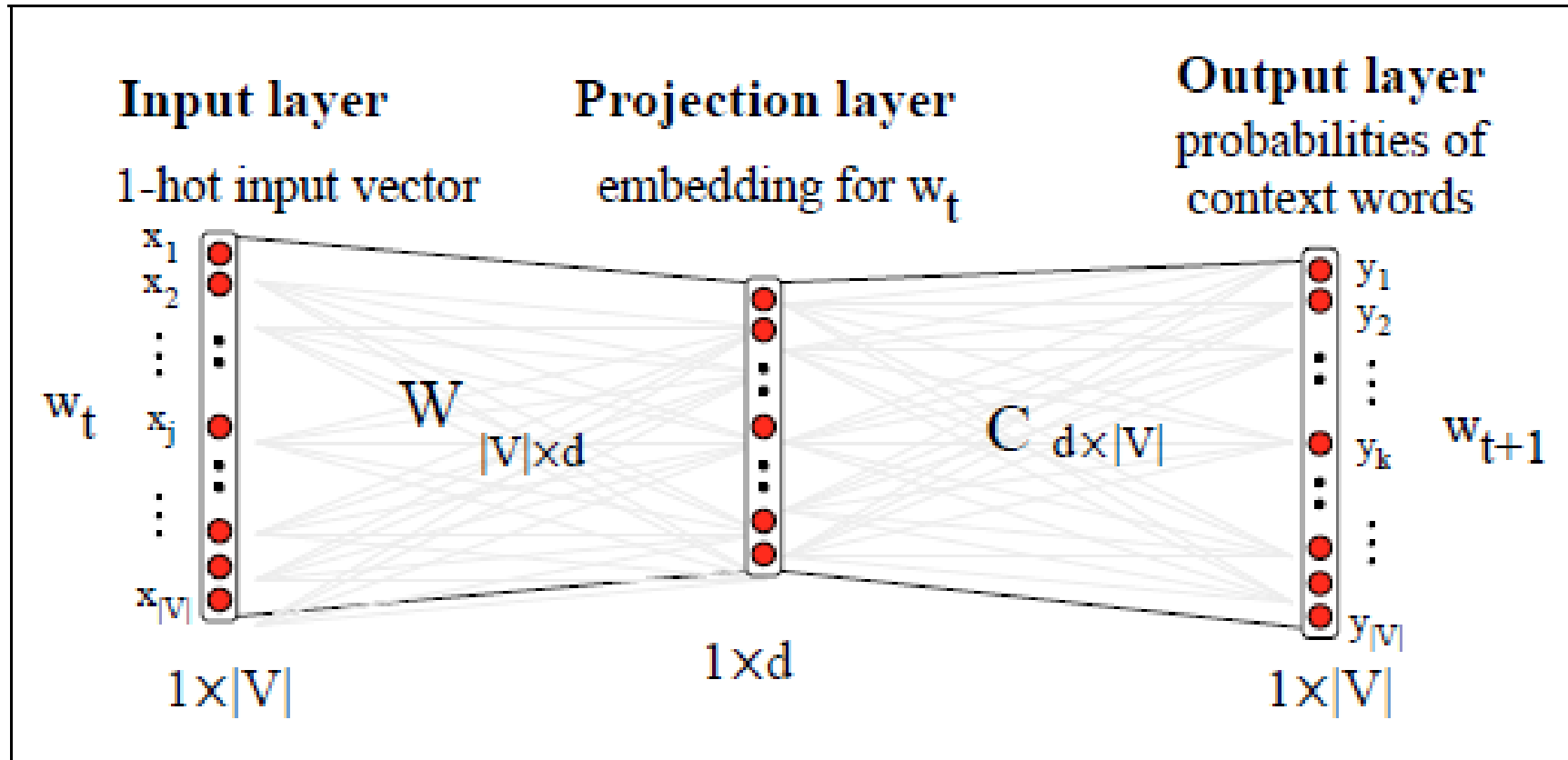
# Idea

47

- Instead of counting, use a neural network to learn a LM
- Simplest form: a bigram model:
  - ▣ For a given word  $w_{i-1}$ , try to predict the next word  $w_i$
  - ▣ i.e. try to estimate  $P(w_i | w_{i-1})$

# Model

48



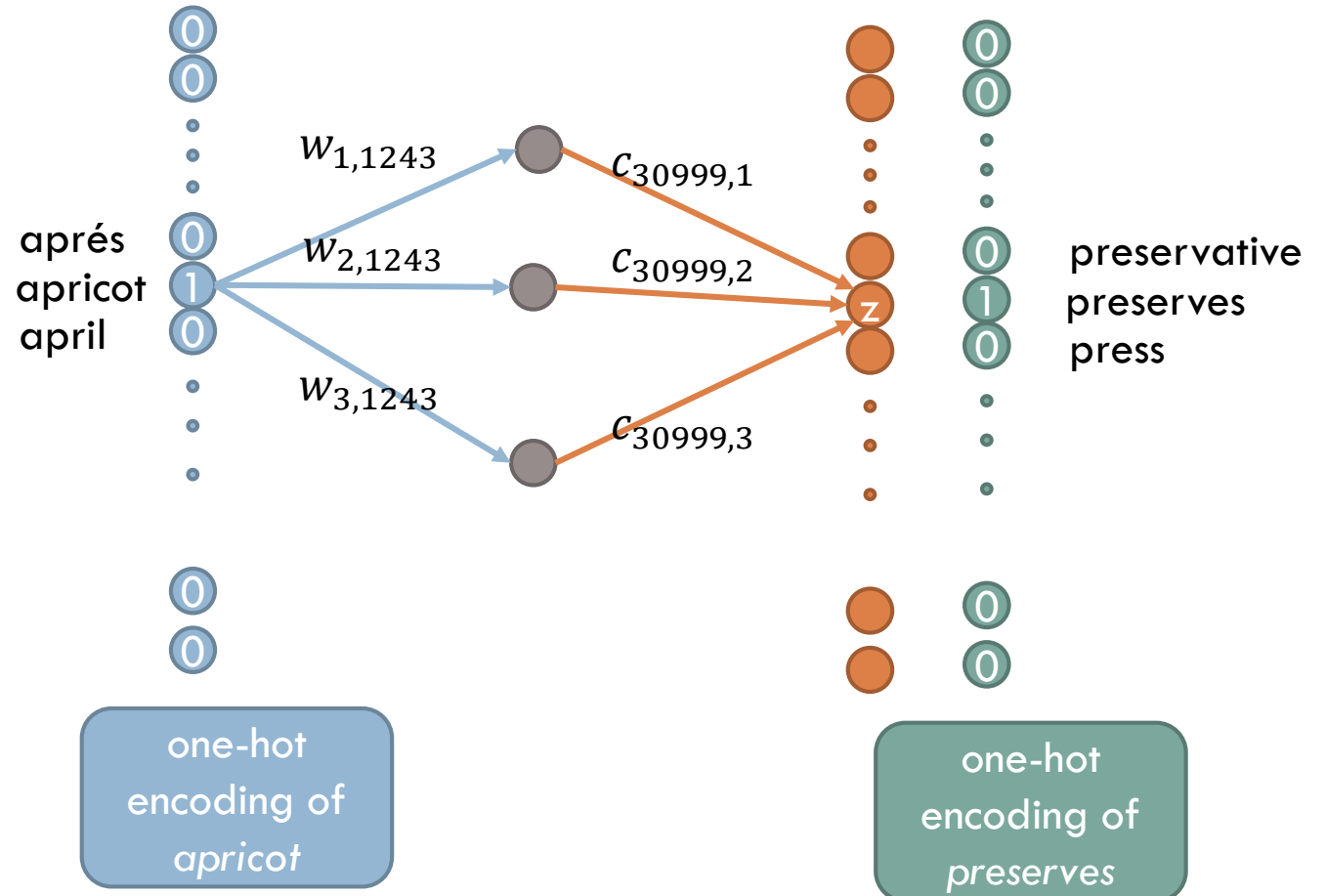
**Figure 16.5** The skip-gram model viewed as a network (Mikolov et al. 2013, Mikolov et al. 2013a).



# Model: zoom in

49

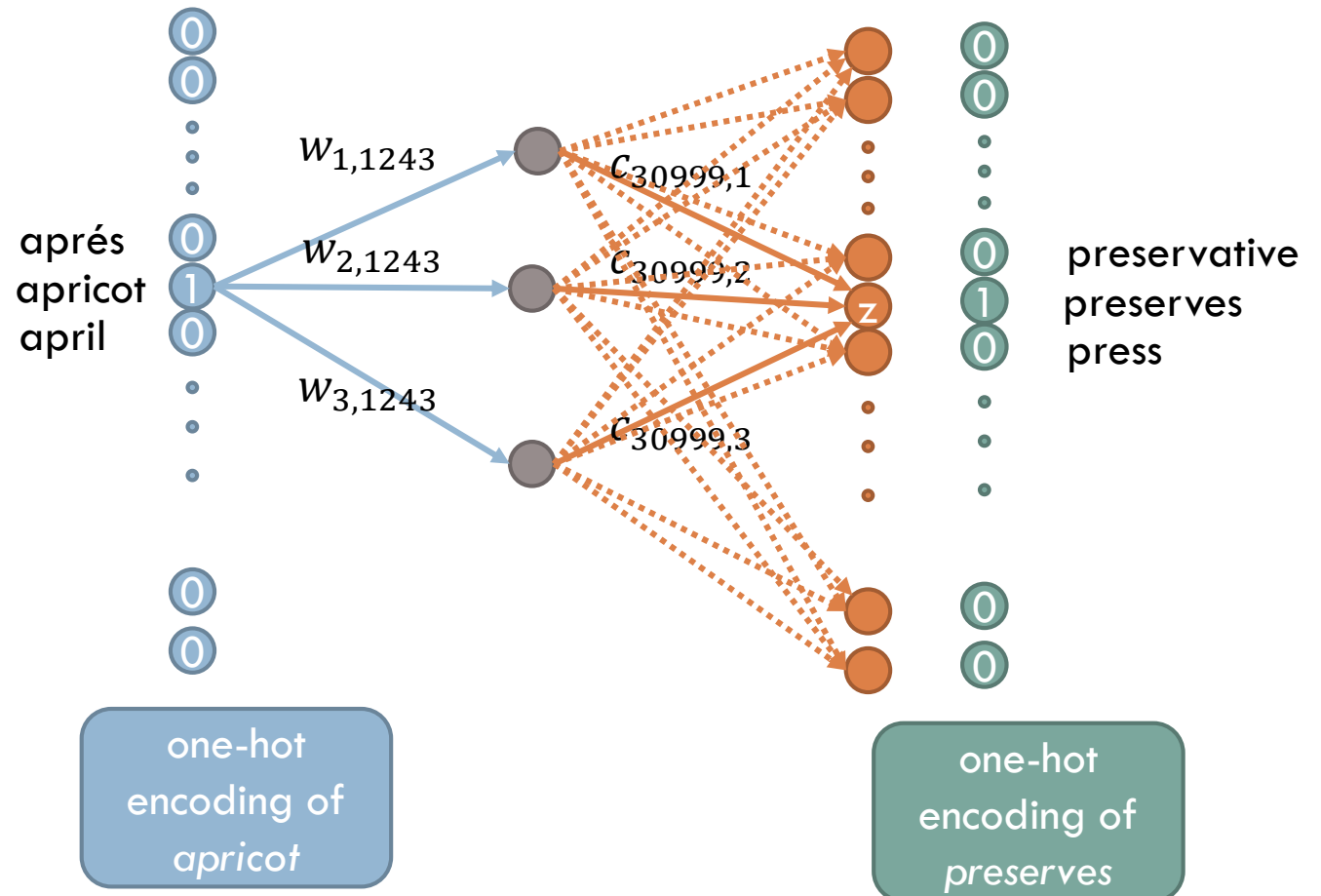
- *apricot* is word 1243
  - word-embedding:
    - $\mathbf{w} = (w_{1,1243}, \dots, w_{d,1243})$
- *preserves* is word 30999
  - context-embedding:
    - $\mathbf{c} = (c_{30999,1}, \dots, c_{30999,d})$
- $\mathbf{z} = \mathbf{w} \cdot \mathbf{c} = \sum_{i=1}^d w_{i,1243} c_{i,30999}$



# Model

50

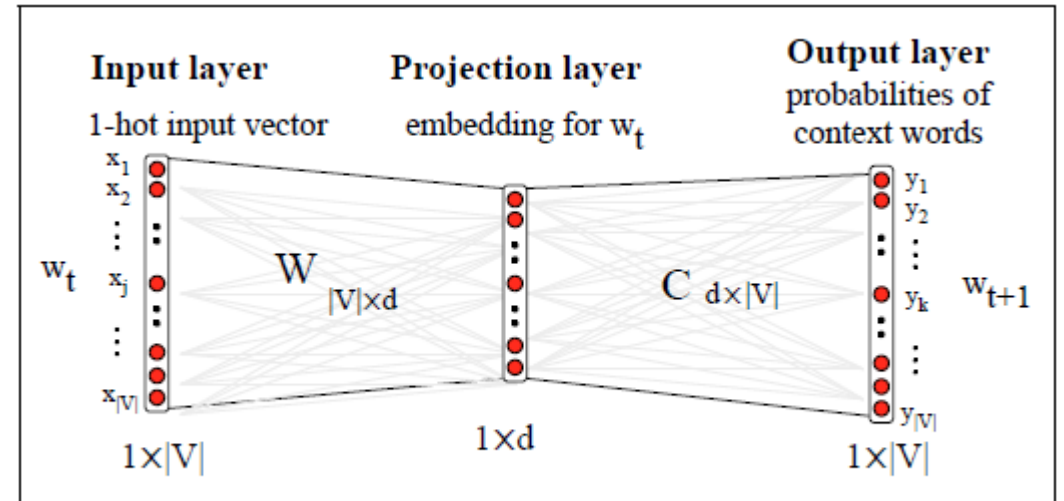
- Input and output word are represented by sparse one-hot vectors
- Dim  $d$  typically 50-300
- Idea for training:
  - ▣ Consider all possible next words for  $w'$  for this word
  - ▣ Use softmax to get a probability distribution of all next words



# Embeddings from this

51

- Idea: Use the weight matrix  $W_{|V| \times d}$  as embeddings, i.e.:
- Represent word  $j$  by  $(w_{j,1}, w_{j,2}, \dots, w_{j,d}) =$  the weights that sends this word to the hidden layer
- Why? since similar words will predict more or less the same words, they will get similar embeddings

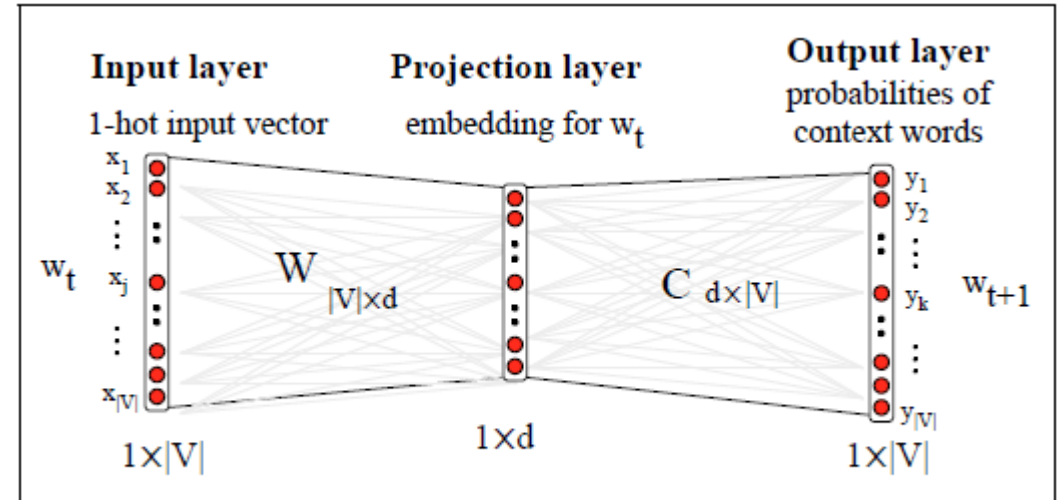


**Figure 16.5** The skip-gram model viewed as a network (Mikolov et al. 2013, Mikolov et al. 2013a).

# Observations

52

- Since two words that are similar are predicted by the same words, there will also be similarities between similar words in  $C_{d \times |V|}$
- This will help the training of  $W_{|V| \times d}$
- We could alternatively use  $C_{d \times |V|}$  as the embeddings

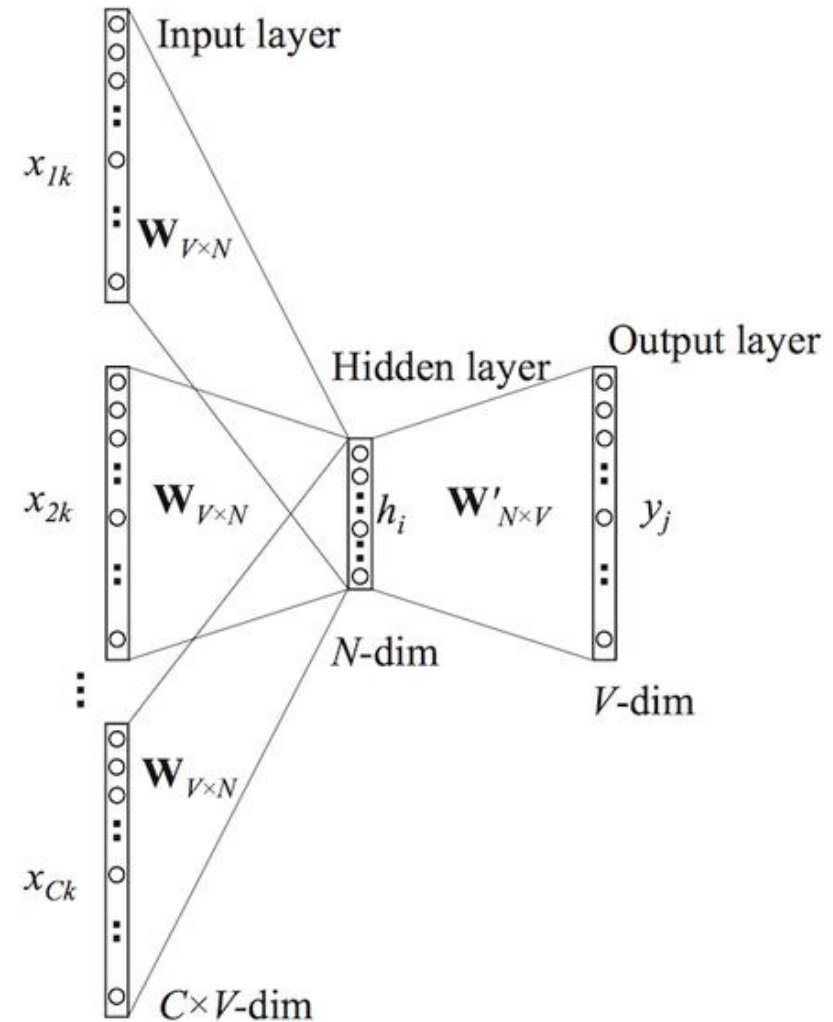


**Figure 16.5** The skip-gram model viewed as a network (Mikolov et al. 2013, Mikolov et al. 2013a).

# CBOW

53

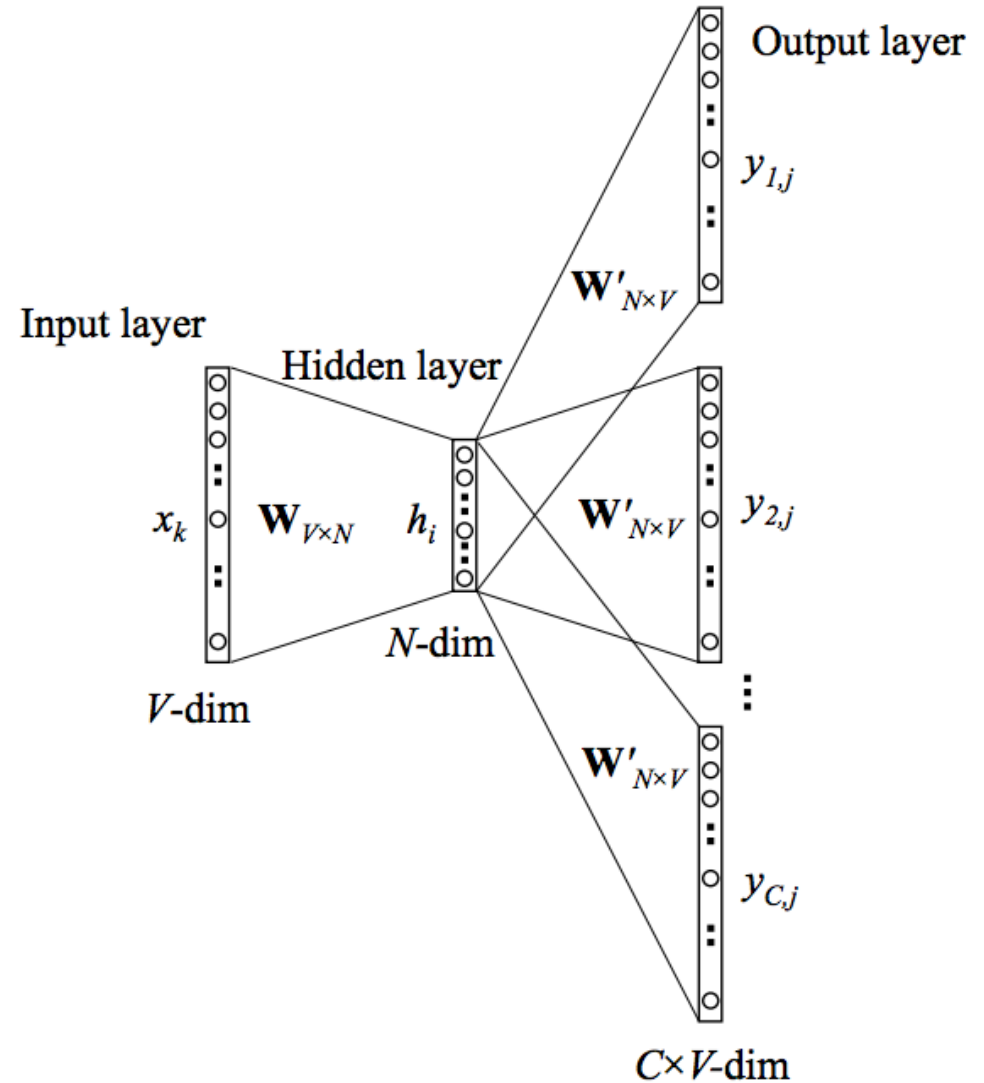
- We could generalize to predicting from a number of preceding words, e.g. 3, as indicated in the figure.
- Observe this is order-independent
- Continuous bag of words model (CBOW):
  - ▣ Predict  $w_t$  from a window  $(w_{t-k}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+k})$



# Skip-gram

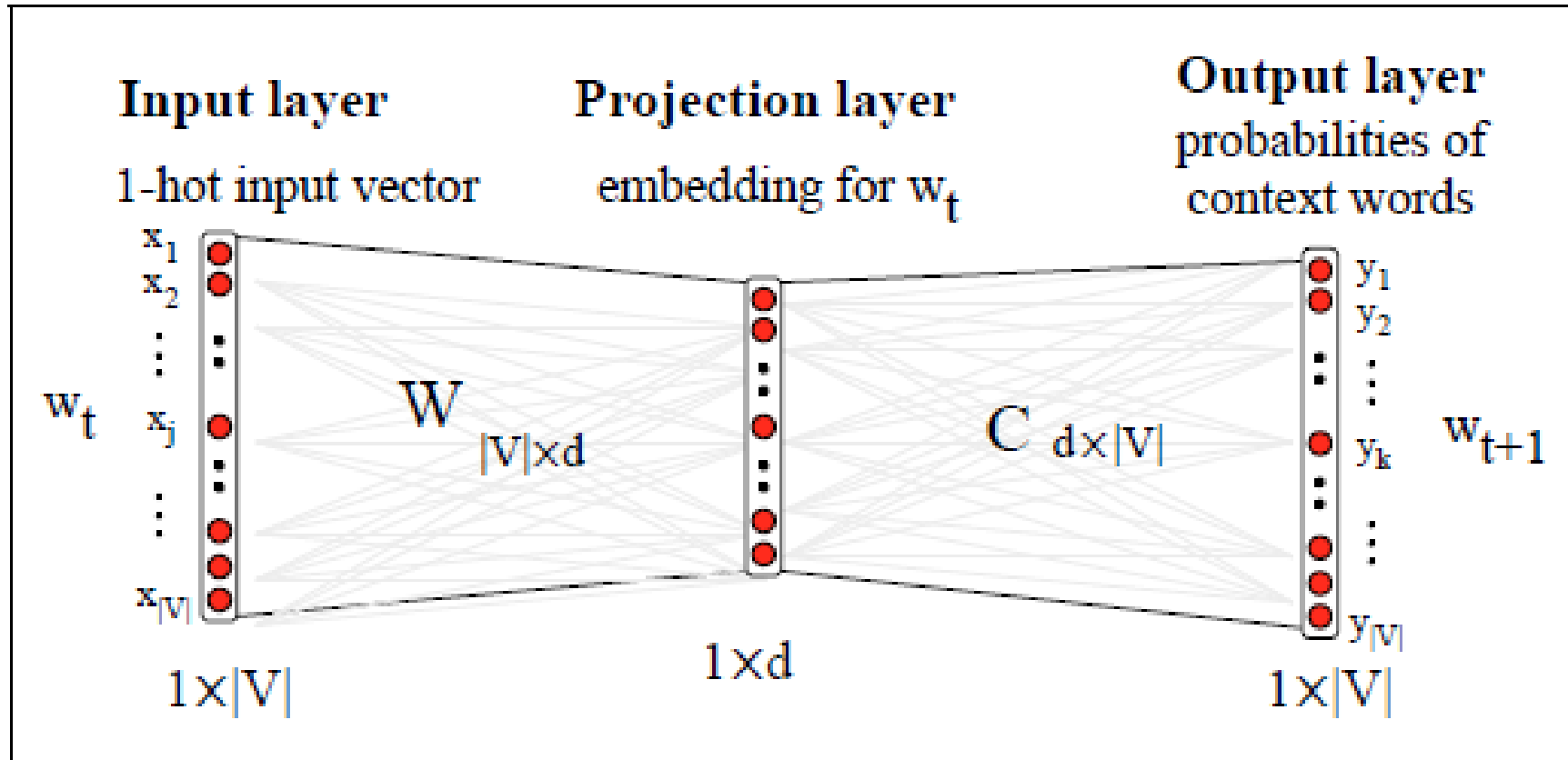
54

- From  $w_t$  predict all the words in a window  
( $w_{t-k}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+k}$ )
- Assume independence of the context words, i.e. from  $w_t$  predict each of the words  $w$  in  $\{w_{t-k}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+k}\}$
- The size of the window will influence which embeddings you get



# Skip-gram model

55

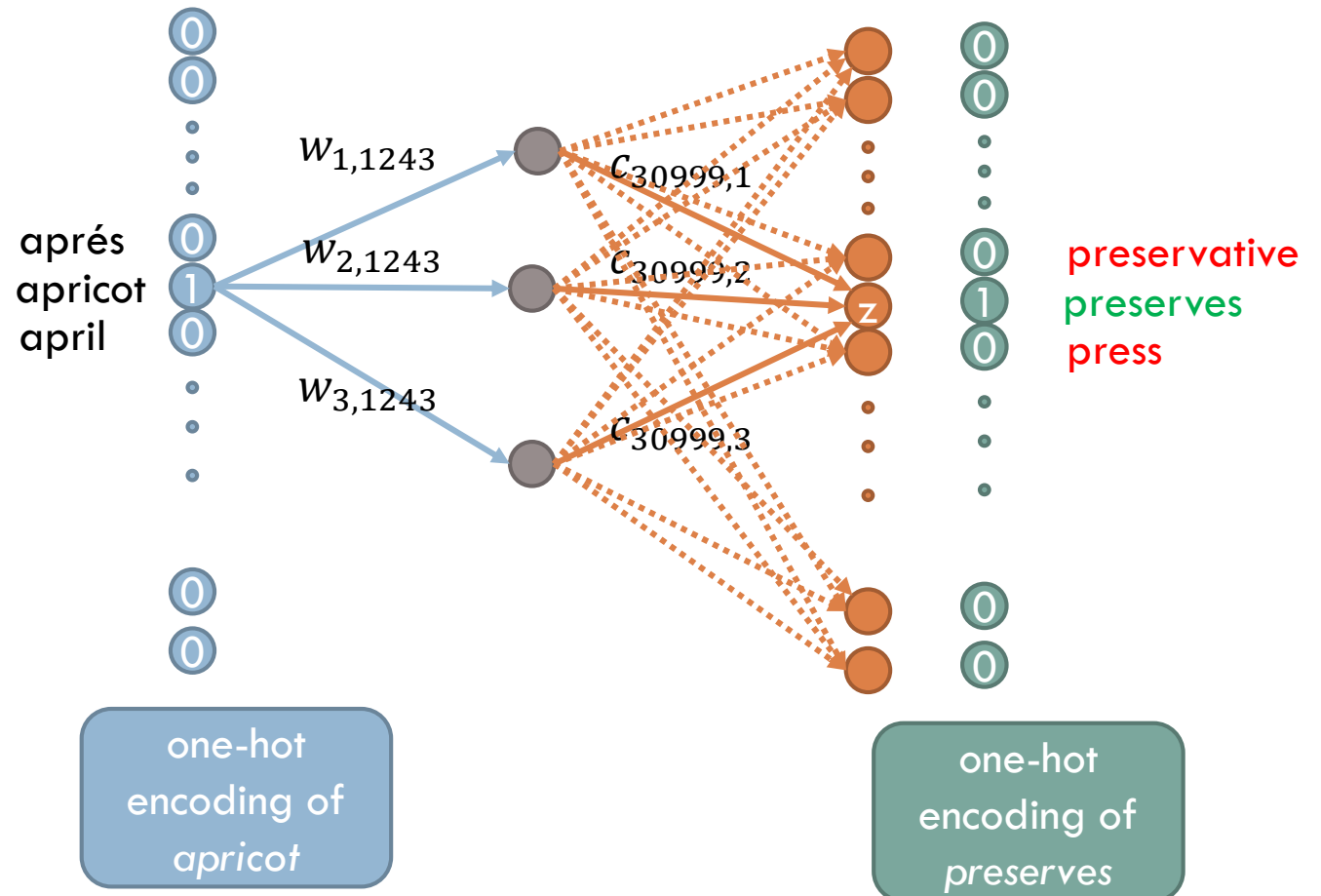


**Figure 16.5** The skip-gram model viewed as a network (Mikolov et al. 2013, Mikolov et al. 2013a).

# Softmax is expensive

56

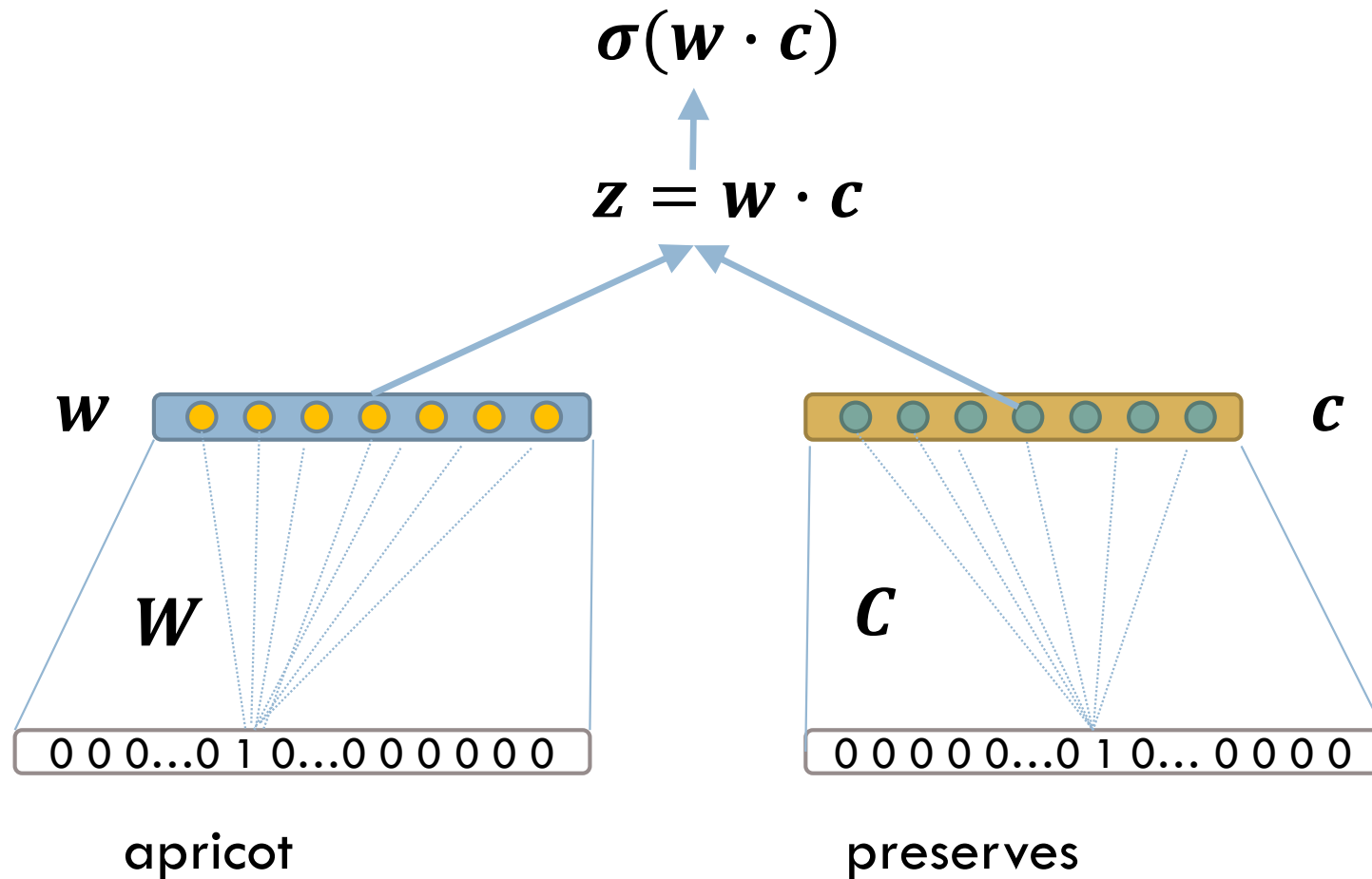
- The use of softmax is expensive
- For one observation, *apricot preserves*, one must change all the  $C_{i,j}$ -s to
  - ▣ increase the probability for *preserves*
  - ▣ decrease all the other probabilities
- $d \times |C|$ , say  $300 \times 50,000$





# Prediction as classification

57



□ To predict preserves from apricot, corresponds to a classification task where

- `class(apricot, preserves)=+`
- `class(apricot, w)=-` for all other `w`

# Skip-gram with negative sampling

58

1. Treat the target word and a neighboring context word as a positive example.
2. Randomly sample other words in the lexicon to get negative samples
  - ▣ sample accordance to frequency
  - ▣ adjusted for high-frequent and low-frequent words: 
$$P_{\alpha}(w) = \frac{\text{count}(w)^{\alpha}}{\sum_{w'} \text{count}(w')^{\alpha}}$$
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the weights as the embeddings

# Skip-Gram Training Data

- Training sentence:

- ... lemon, a tablespoon of **apricot** preserves or a ...

- c1                   c2                   t                   c3                   c4

- Training data: input/output pairs centering on *apricot*

- Assume a +/- 2 word window

# Skip-Gram Training Data

60

- ... lemon, a **tablespoon** of **apricot** preserves or a ...
- **c1** **c2** **t** **c3** **c4**
- For each positive example, we'll create  $k$  negative examples.
  - ▣ Using *noise* words: Any random word that isn't  $t$

positive examples +	
t	c
apricot	tablespoon
apricot	of
apricot	preserves
apricot	or

negative examples -			
t	c	t	c
apricot	aardvark	apricot	twelve
apricot	puddle	apricot	hello
apricot	where	apricot	dear
apricot	coaxial	apricot	forever

# Learning

61

- Like Logistic Regression
- Start with randomly initialized weights for  $\mathbf{W}$  and  $\mathbf{C}$
- For the training items  $(\mathbf{w}, \mathbf{c})$ , calculate  $\hat{y} = \sigma(\mathbf{c} \cdot \mathbf{w}) = \frac{1}{1+e^{-\mathbf{c} \cdot \mathbf{w}}}$
- Compare to the gold labels using cross-entropy loss
  - ▣ The gold label is 1 if  $\mathbf{c}$  is a context word and 0 if  $\mathbf{c}$  is a negative example
  - ▣ This is like Logistic regression
- Use the derivative of the loss with respect to  $\mathbf{c}$ :  $\frac{\partial}{\partial \mathbf{c}} L_{ce}$  to update  $\mathbf{c}$
- and the derivative of the loss with respect to  $\mathbf{w}$  to update  $\mathbf{w}$

# Update equations in SGD

62

- We skip the derivation, but these are the resulting update equations

$$\mathbf{c}_{pos}^{t+1} = \mathbf{c}_{pos}^t - \eta [\sigma(\mathbf{c}_{pos}^t \cdot \mathbf{w}^t) - 1] \mathbf{w}^t$$

$$\mathbf{c}_{neg}^{t+1} = \mathbf{c}_{neg}^t - \eta [\sigma(\mathbf{c}_{neg}^t \cdot \mathbf{w}^t)] \mathbf{w}^t$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \left[ [\sigma(\mathbf{c}_{pos} \cdot \mathbf{w}^t) - 1] \mathbf{c}_{pos} + \sum_{i=1}^k [\sigma(\mathbf{c}_{neg_i} \cdot \mathbf{w}^t)] \mathbf{c}_{neg_i} \right]$$

- $\hat{y} = \sigma(\mathbf{c} \cdot \mathbf{w})$
- Similar to the logistic regression, where we update weights
- Her we update both the  $w$ -s and the  $c$ -s.

# Result

63

- We learn two separate embedding matrices  $W$  and  $C$
- We can use  $W$  as representations for the words
  - ▣ (or combine with  $C$  in some ways)
  
- What have we learned:
  - ▣ If two words  $w_1$  and  $w_2$  occur in similar contexts
    - = with the same (or similar) context words, e.g.  $c$ ,
  - ▣ then both  $w_1$  and  $w_2$  should have a large cosine with  $c$ ,
    - hence get similar vectors.

# Use of embeddings

64

- Embeddings are used as representations for words as input in all kinds of NLP tasks using deep learning:
  - ▣ Text classification
  - ▣ Language models
  - ▣ Named-entity recognition
  - ▣ Machine translation
  - ▣ etc.
- These embeddings are nowadays called **static**
- Since 2018, **Transformers**:
  - ▣ The embedding of each word depends on the context
  - ▣ Superior results in all tasks
- IN5550, Spring



# Resources

65

- gensim
  - ▣ Easy-to-use tool for training own models
- Word2vec
  - ▣ <https://code.google.com/archive/p/word2vec/>
- <https://fasttext.cc/>
- <https://nlp.stanford.edu/projects/glove/>
- <http://vectors.nlpl.eu/repository/>
  - ▣ Pretrained embeddings, also for Norwegian