

# Word Representation



CS 288 Spring 2026  
UC Berkeley  
[cal-cs288.github.io/sp26](https://cal-cs288.github.io/sp26)

Berkeley **BAIR**  
EECS

# Today's Question: **How to Represent a Word?**

Computers don't understand words – they see them as symbols.

But can we make a computer understand *that ‘king’ and ‘queen’ are related*, or that *‘dog’ is similar to ‘cat’*?

# A naive option: One-hot encoding

Each word in the vocabulary is assigned a unique index and is represented as a binary vector with a single 1 in the position corresponding to the word's index and 0s elsewhere.

$$v_{\text{cat}} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \dots \end{bmatrix} \quad v_{\text{dog}} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \dots \end{bmatrix} \quad v_{\text{the}} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \dots \end{bmatrix} \quad v_{\text{language}} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \dots \end{bmatrix}$$

Lookup table:

Index	0	1	2	3	4
Word	cat	dog	the	Language	...

High-dimensional & Sparse

Does not capture semantic meaning of the words (all words are equally distant)

# Need for word meaning

- With words, a feature is a word identity (= string)
- Requires **exact same word** to be in the training and testing set  
  
“terrible”  $\neq$  “horrible”
- If we can represent **word meaning** in vectors:
  - The previous word was vector [35, 22, 17, ...]
  - Now in the test set we might see a similar vector [34, 21, 14, ...]
  - We can generalize to **similar but unseen words!!!**

# What do words mean?

- **Synonyms:** couch/sofa, car/automobile, filbert/hazelnut
- **Antonyms:** dark/light, rise/fall, up/down
- Some words are not synonyms but they share some element of meaning
  - cat/dog, car/bicycle, cow/horse
- Some words are not similar but they are **related**
  - coffee/cup, house/door, chef/menu
- **Affective meanings or connotations:**

**valence:** the pleasantness of the stimulus

**arousal:** the intensity of emotion provoked by the stimulus

**dominance:** the degree of control exerted by the stimulus

vanish	disappear	9.8
belief	impression	5.95
muscle	bone	3.65
modest	flexible	0.98
hole	agreement	0.3

SimLex-999

	Valence	Arousal	Dominance
courageous	8.05	5.5	7.38
music	7.67	5.57	6.5
heartbreak	2.45	5.65	3.58
cub	6.71	3.95	4.24

# Lexical resources

**WordNet Search - 3.1**  
- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations  
Display options for sense: (gloss) "an example sentence"

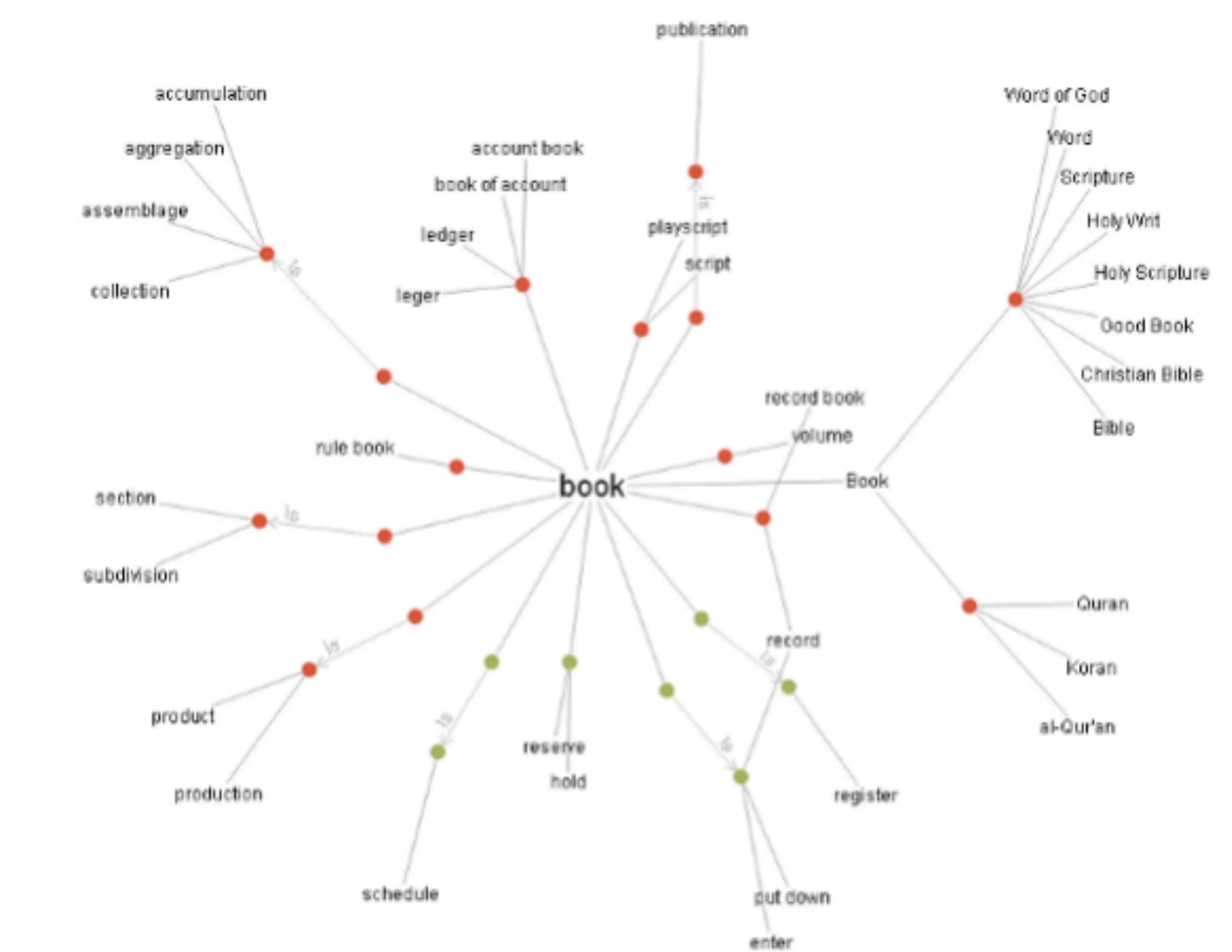
**Noun**

- S: (n) **mouse** (any of numerous small rodents typically resembling diminutive rats having pointed snouts and small ears on elongated bodies with slender usually hairless tails)
- S: (n) **shiner, black eye, mouse** (a swollen bruise caused by a blow to the eye)
- S: (n) **mouse** (person who is quiet or timid)
- S: (n) **mouse, computer mouse** (a hand-operated electronic device that controls the coordinates of a cursor on your computer screen as you move it around on a pad; on the bottom of the device is a ball that rolls on the surface of the pad) "*a mouse takes much more room than a trackball*"

**Verb**

- S: (v) **sneak, mouse, creep, pussyfoot** (to go stealthily or furtively) "*..instead of sneaking around spying on the neighbor's house*"
- S: (v) **mouse** (manipulate the mouse of a computer)

<https://wordnet.princeton.edu/>



(-) Huge amounts of human labor to create and maintain

# Distributional hypothesis

# A simple linguistic intuition

“Words that occur in similar contexts tend to have similar meanings.”

Example:

- “The **cat** sat on the \_\_\_\_.”
- “The **dog** lay on the \_\_\_\_.”

# Distributional semantics

Words that occur in similar **contexts** tend to have similar meanings



J.R.Firth 1957

- “You shall know a word by the company it keeps”
- One of the most successful ideas of modern statistical NLP!

When a word  $w$  appears in a text, its context is the set of words that appear nearby (within a fixed-size window).

**This idea turns semantics into a statistical problem over text.**

# Distributional semantics

*...government debt problems turning into banking crises as happened in 2009...*

*...saying that Europe needs unified banking regulation to replace the hodgepodge...*

*...India has just given its banking system a shot in the arm...*



These **context words** will represent **banking**

# Quick quiz: Word guessing!

(Example from Eisenstein's book)

Everybody likes tezgūino.

We make tezgūino out of corn.

A bottle of tezgūino is on the table.

Don't have tezgūino before you drive.

# Distributional semantics

How we can do the same thing computationally?

- Count the words in the context of a target word
- See what other words occur in those contexts

We can represent a word's context using vectors!

# Word-word co-occurrence matrix

Solution: Let's use **word-word co-occurrence counts** to represent the meaning of words!  
Each word is represented by the corresponding **row vector**

**context words:**

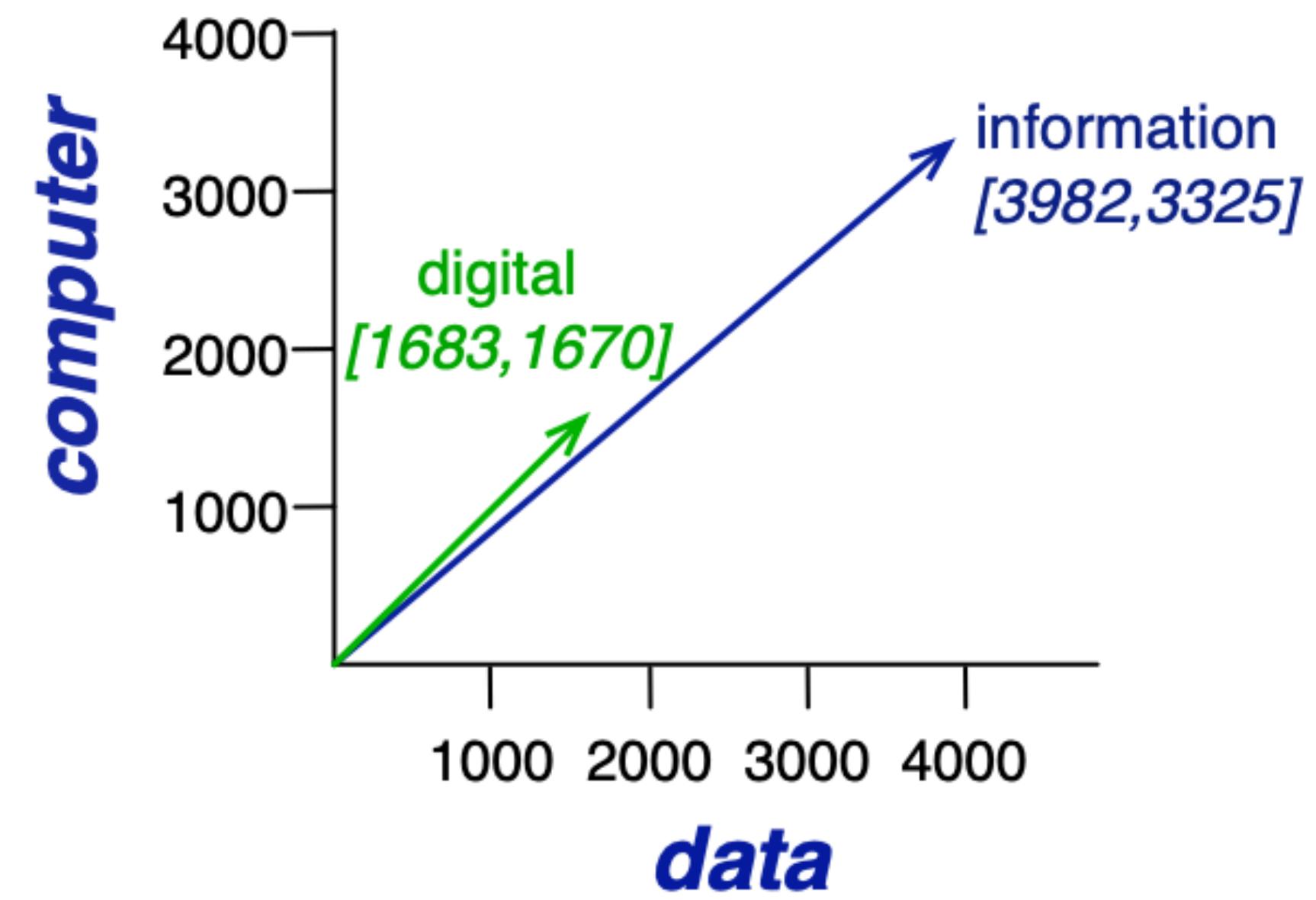
4 words to the left +  
4 words to the right

is traditionally followed by **cherry** pie, a traditional dessert  
often mixed, such as **strawberry** rhubarb pie. Apple pie  
computer peripherals and personal **digital** assistants. These devices usually  
a computer. This includes **information** available on the internet

	aardvark	...	computer	data	result	pie	sugar	...
cherry	0	...	2	8	9	442	25	...
strawberry	0	...	0	0	1	60	19	...
digital	0	...	1670	1683	85	5	4	...
information	0	...	3325	3982	378	5	13	...

Very high-dimensional; Most entries are 0s  $\Rightarrow$  sparse vectors

# Measuring similarity



A common similarity metric: **cosine** of the angle between the two vectors (the larger, the more similar the two vectors are)

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^{|V|} u_i v_i}{\sqrt{\sum_{i=1}^{|V|} u_i^2} \sqrt{\sum_{i=1}^{|V|} v_i^2}}$$

Q: Why cosine similarity instead of dot product  $\mathbf{u} \cdot \mathbf{v}$ ?

# Measuring similarity

What is the range of  $\cos(u, v)$  if  $u, v$  are **count vectors**?

- (A)  $[-1, 1]$
- (B)  $[0, 1]$
- (C)  $(0, 1)$
- (D)  $(-1, 1)$
- (E)  $(-\infty, +\infty)$

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^{|V|} u_i v_i}{\sqrt{\sum_{i=1}^{|V|} u_i^2} \sqrt{\sum_{i=1}^{|V|} v_i^2}}$$

The answer is (b). Cosine similarity ranges between -1 and 1 in general. In this model, all the values of  $u_i, v_i$  are non-negative.

# Sparse vs. dense vectors

# Sparse vs. dense vectors

- The vectors in the word-word occurrence matrix are
  - Long: vocabulary size
  - Sparse: most are 0's
- Alternative: we want to represent words as short (50-300 dimensional) & dense (real-valued) vectors
- This is the basis for modern NLP systems!

$$v_{\text{cat}} = \begin{pmatrix} -0.224 \\ 0.130 \\ -0.290 \\ 0.276 \end{pmatrix}$$

$$v_{\text{dog}} = \begin{pmatrix} -0.124 \\ 0.430 \\ -0.200 \\ 0.329 \end{pmatrix}$$

$$v_{\text{the}} = \begin{pmatrix} 0.234 \\ 0.266 \\ 0.239 \\ -0.199 \end{pmatrix}$$

$$v_{\text{language}} = \begin{pmatrix} 0.290 \\ -0.441 \\ 0.762 \\ 0.982 \end{pmatrix}$$

# Why dense vectors?

- Short vectors are easier to use as features in ML systems
- Dense vectors may generalize better than explicit counts
- They can't capture high-order co-occurrence
  - $w_1$  co-occurs with "car",  $w_2$  co-occurs with "automobile"
  - They should be similar but they aren't because "car" and "automobile" are distinct dimensions
- In practice, they work better!

# How to get short dense vectors?

- **Count-based methods:** Singular value decomposition (SVD) of count matrix

$$\begin{bmatrix} X \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} W \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_V \end{bmatrix} \begin{bmatrix} C \\ |V| \times |V| \end{bmatrix}$$

embedding  
for  
word i

$$W \quad |V| \times k$$

$$\begin{bmatrix} X \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} W \\ |V| \times k \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_k \end{bmatrix} \begin{bmatrix} C \\ k \times |V| \end{bmatrix}$$

We can approximate the full matrix by only keeping the top  $k$  (e.g., 100) singular values!

# How to get short dense vectors?

- **Count-based methods:** Singular value decomposition (SVD) of count matrix
- **Prediction-based methods:**
  - Vectors are created by training a classifier to predict whether a word c (“pie”) is likely to appear in the context of a word w (“cherry”)
  - Examples: **word2vec** (Mikolov et al., 2013), **Glove** (Pennington et al., 2014), **FastText** (Bojanowski et al., 2017)

Also called **word embeddings!**

*Don't count, predict!* A systematic comparison of  
context-counting vs. context-predicting semantic vectors

Marco Baroni and Georgiana Dinu and Germán Kruszewski  
Center for Mind/Brain Sciences (University of Trento, Italy)

(Baroni et al., 2014)

# Word embeddings

# Word embeddings

= **Learned** representations from text for representing words

- Input: a large text corpora,  $V, d$ 
  - $V$ : a pre-defined vocabulary
  - $d$ : dimension of word vectors (e.g. 300)
  - Text corpora:
    - Wikipedia + Gigaword 5: 6B tokens
    - Twitter: 27B tokens
    - Common Crawl: 840B tokens
- Output:  $f : V \rightarrow \mathbb{R}^d$

$$v_{\text{cat}} = \begin{pmatrix} -0.224 \\ 0.130 \\ -0.290 \\ 0.276 \end{pmatrix} \quad v_{\text{dog}} = \begin{pmatrix} -0.124 \\ 0.430 \\ -0.200 \\ 0.329 \end{pmatrix}$$
$$v_{\text{the}} = \begin{pmatrix} 0.234 \\ 0.266 \\ 0.239 \\ -0.199 \end{pmatrix} \quad v_{\text{language}} = \begin{pmatrix} 0.290 \\ -0.441 \\ 0.762 \\ 0.982 \end{pmatrix}$$

Each word is represented by a low-dimensional (e.g.,  $d = 300$ ), real-valued vector

Each coordinate/dimension of the vector doesn't have a particular interpretation

# Trained word embeddings available

- word2vec: <https://code.google.com/archive/p/word2vec/>
- GloVe: <https://nlp.stanford.edu/projects/glove/>
- FastText: <https://fasttext.cc/>

## Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License](#) v1.0 whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
  - [Wikipedia 2014 + Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
  - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
  - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
  - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data

Differ in algorithms, text corpora, dimensions, cased/uncased...  
Applied to many other languages

# Word embeddings

- Basic property: similar words have similar vectors

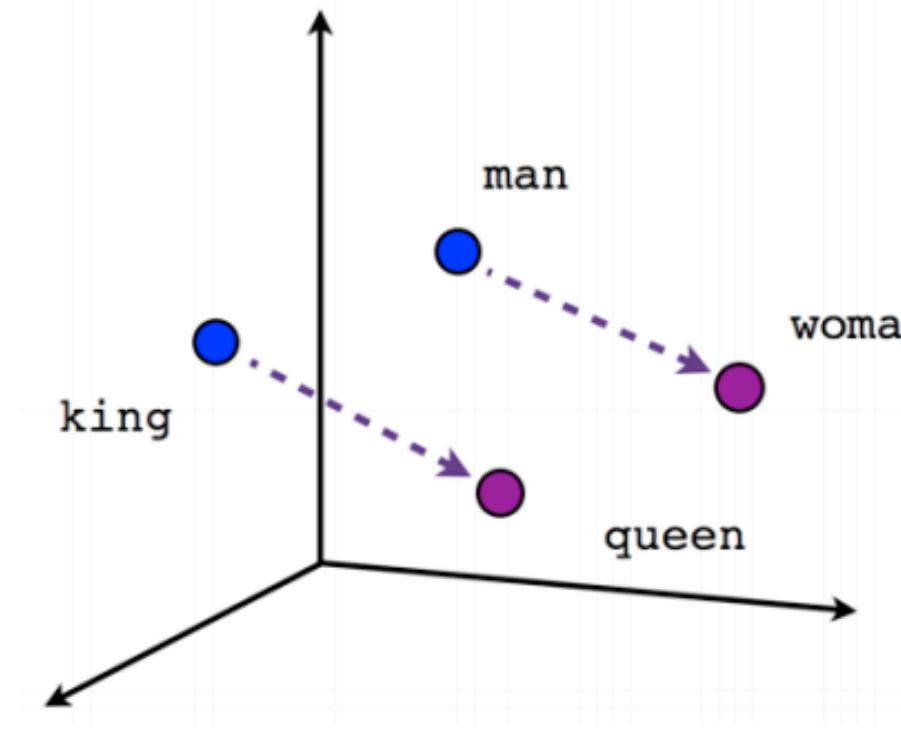
word  $w^* = \text{"sweden"}$   
 $\arg \max_{w \in V} \cos(e(w), e(w^*))$

Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408

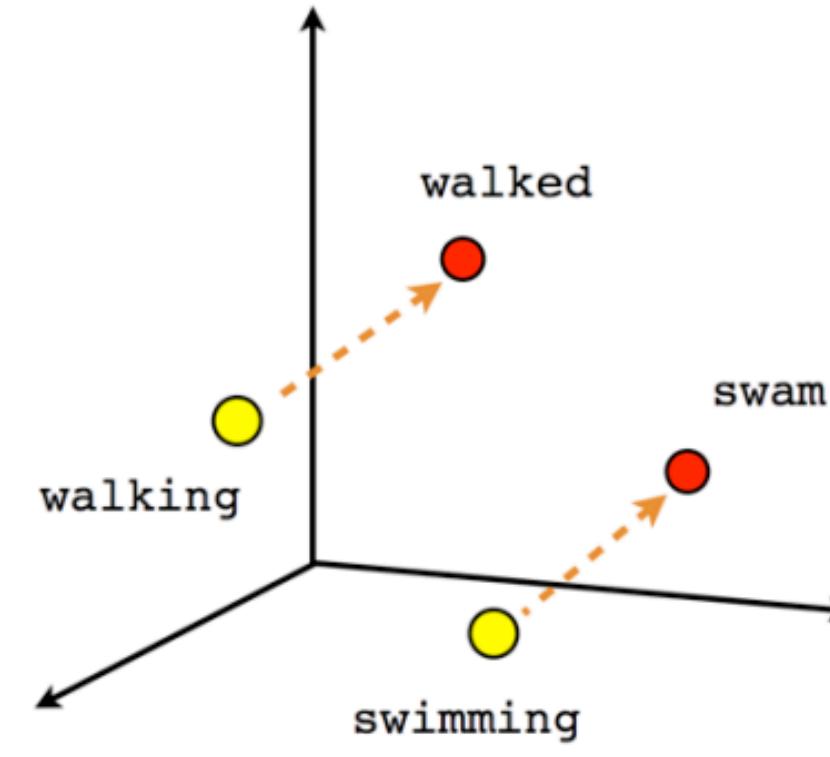
$\cos(u, v)$  ranges between -1 and 1

# Word embeddings

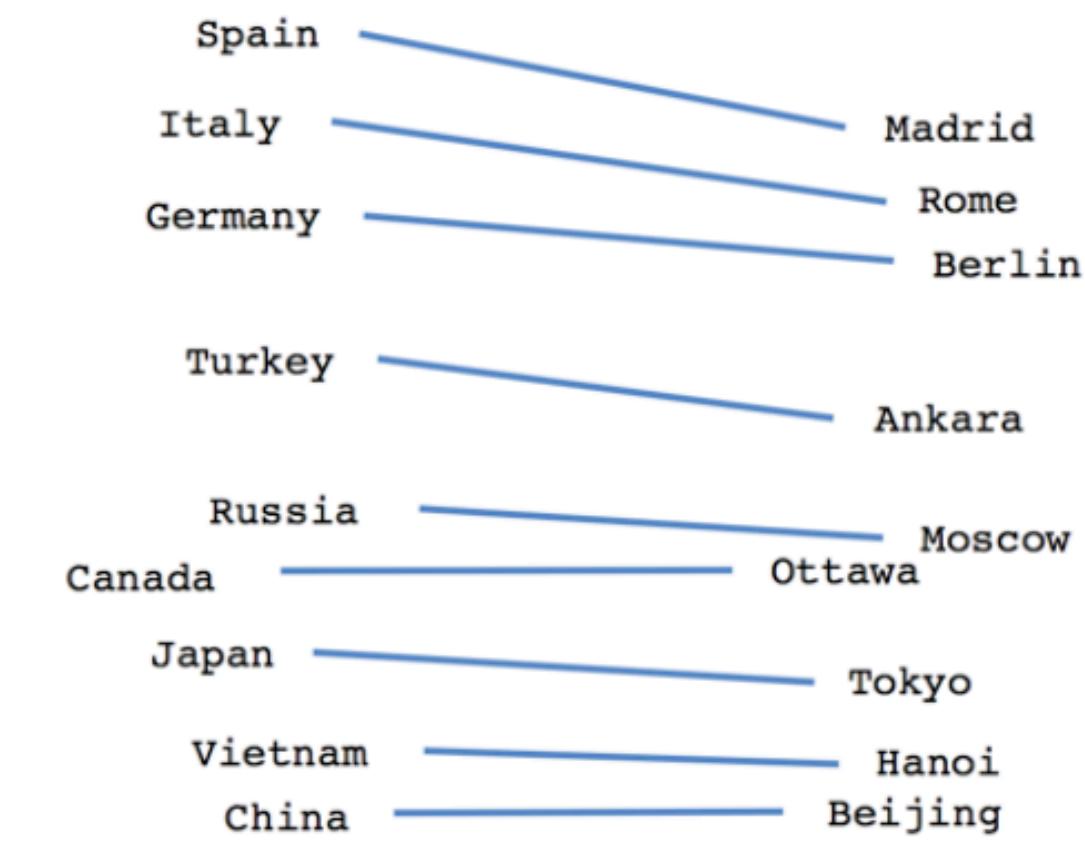
- They have some other nice properties too!



Male-Female



Verb tense



Country-Capital

$$v_{\text{man}} - v_{\text{woman}} \approx v_{\text{king}} - v_{\text{queen}}$$

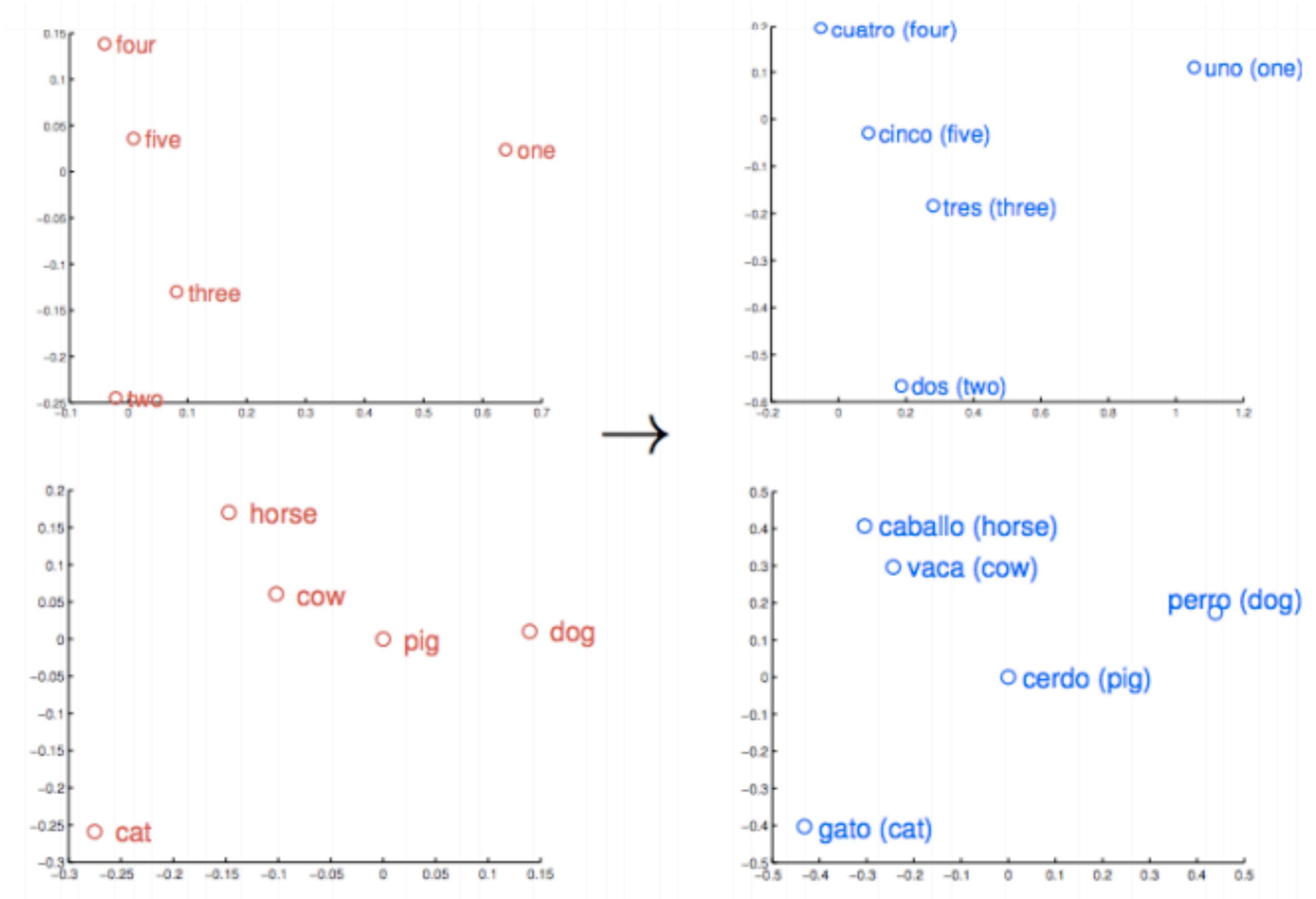
$$v_{\text{Paris}} - v_{\text{France}} \approx v_{\text{Rome}} - v_{\text{Italy}}$$

Word analogy test:  $a : a^* :: b : b^*$

$$b^* = \arg \max_{w \in V} \cos(e(w), e(a^*) - e(a) + e(b))$$

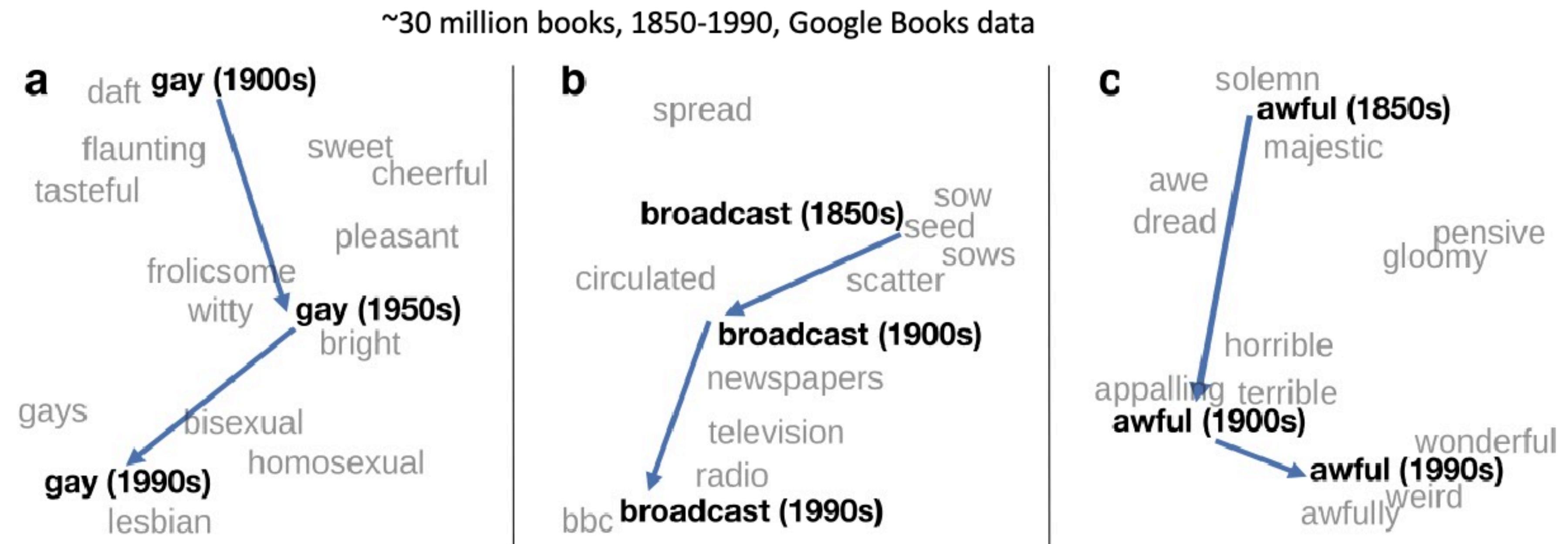
# Word embeddings

- They have some other nice properties too!



# Embeddings as a window onto historical semantics

Train embeddings on different decades of historical text to see meanings shift



William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. Proceedings of ACL.

# Embeddings reflect cultural bias!

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In *NeurIPS*, pp. 4349-4357. 2016.

Ask “Paris : France :: Tokyo : x”

- x = Japan

Ask “father : doctor :: mother : x”

- x = nurse

Ask “man : computer programmer :: woman : x”

- x = homemaker

Algorithms that use embeddings as part of e.g., hiring searches for programmers, might lead to bias in hiring

# Word embeddings: the learning problem

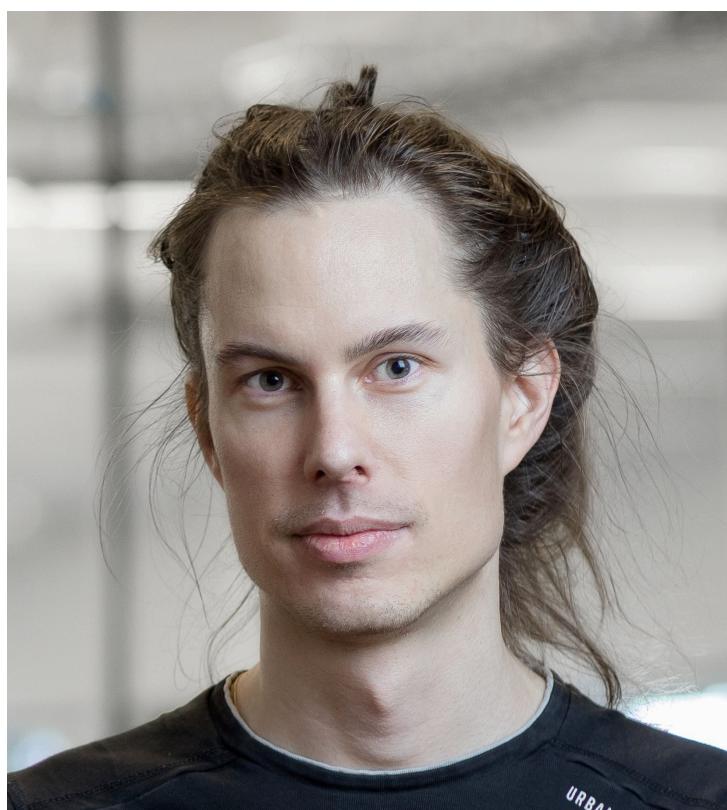
**Learning** vectors from text for representing words

- **Input:**
  - a large text corpus,
  - vocabulary  $V$
  - vector dimension  $d$  (e.g., 300)
- **Output:**  $f: V \rightarrow \mathbb{R}^d$

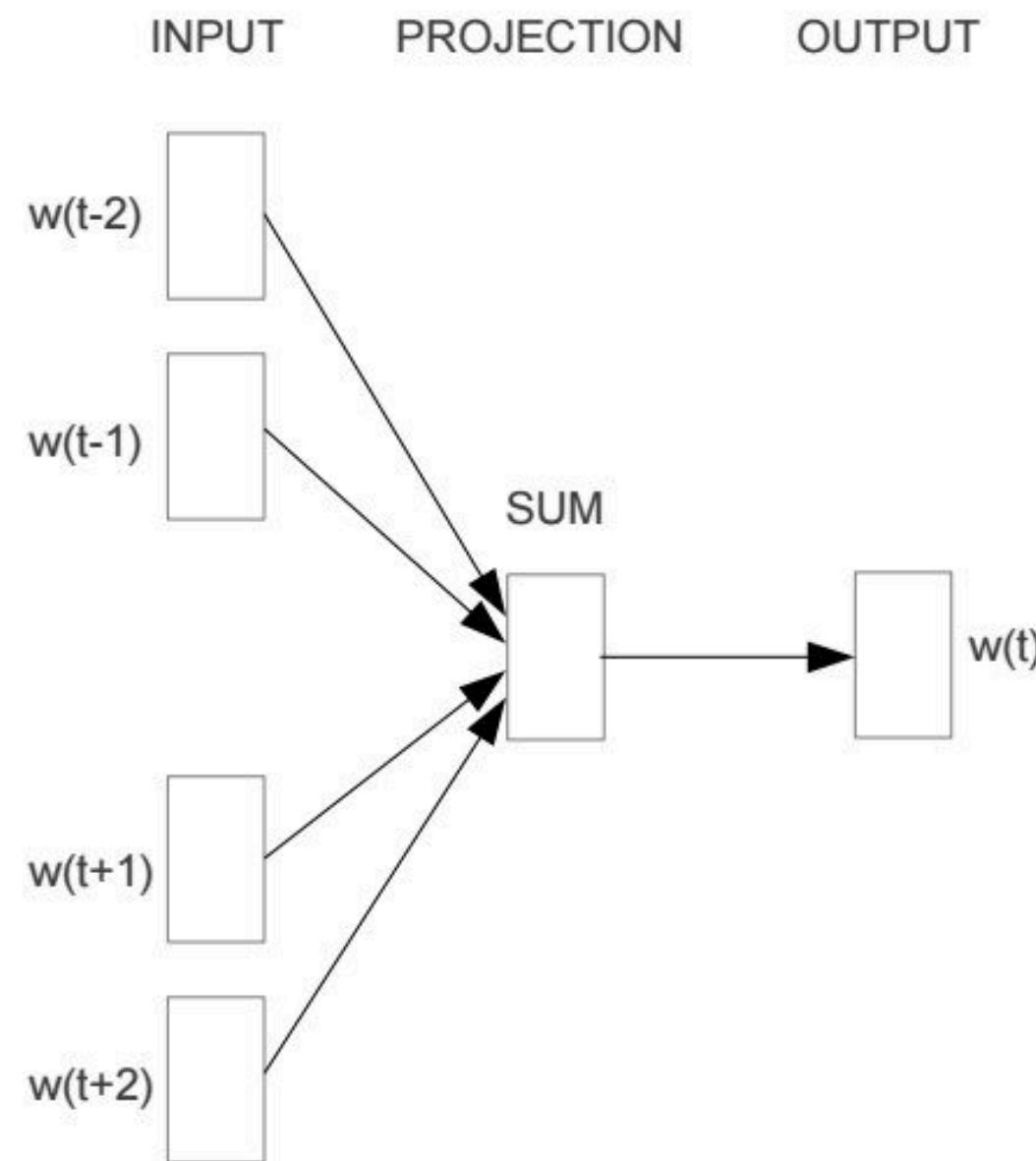
$$v_{\text{cat}} = \begin{pmatrix} -0.224 \\ 0.130 \\ -0.290 \\ 0.276 \end{pmatrix} \quad v_{\text{dog}} = \begin{pmatrix} -0.124 \\ 0.430 \\ -0.200 \\ 0.329 \end{pmatrix}$$
$$v_{\text{the}} = \begin{pmatrix} 0.234 \\ 0.266 \\ 0.239 \\ -0.199 \end{pmatrix} \quad v_{\text{language}} = \begin{pmatrix} 0.290 \\ -0.441 \\ 0.762 \\ 0.982 \end{pmatrix}$$

# Word2vec

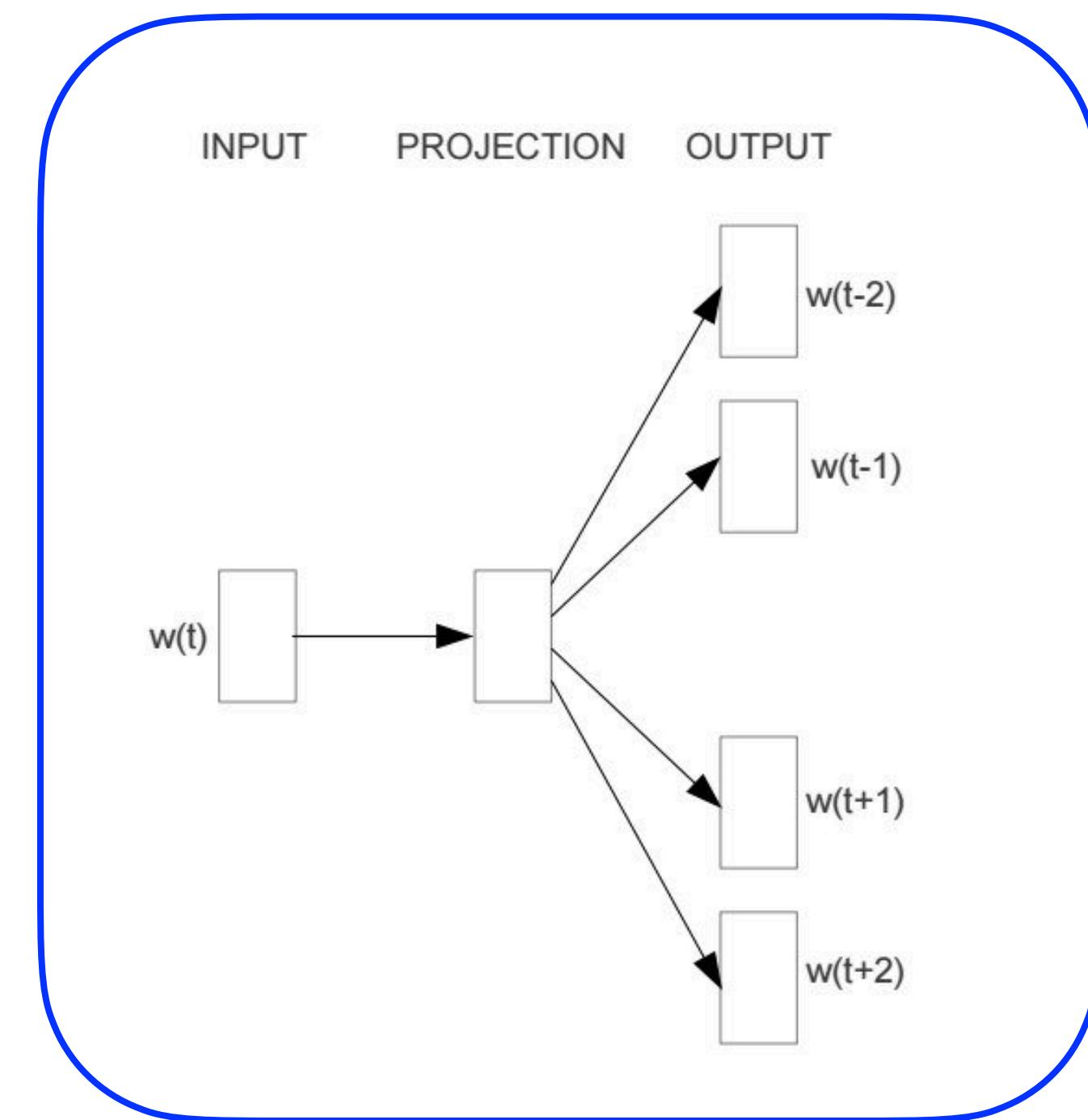
- (Mikolov et al 2013a): Efficient Estimation of Word Representations in Vector Space
- (Mikolov et al 2013b): Distributed Representations of Words and Phrases and their Compositionality



Tomáš Mikolov



Continuous Bag of Words (CBOW)



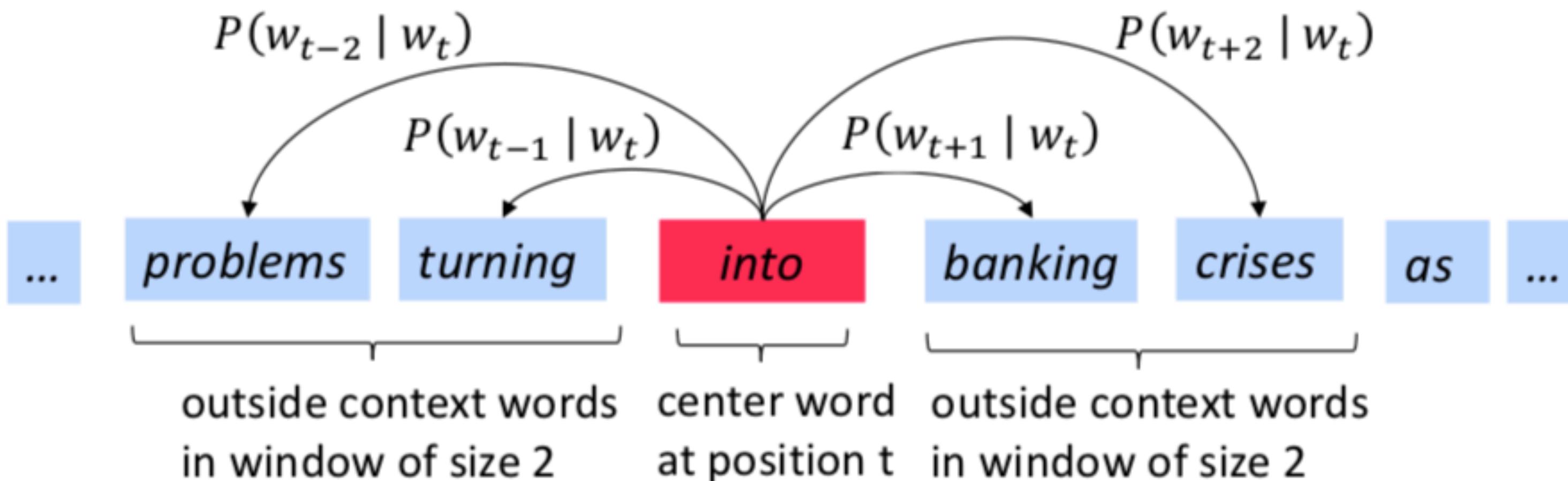
Skip-gram

# Skip-gram

- Key idea: Use each word to predict other words in its context
- $P(b | a)$  = given the center word is  $a$ , what is the probability that  $b$  is a context word?
- Context: a fixed window of size  $2m$  ( $m = 2$  in the example)

A classification problem!

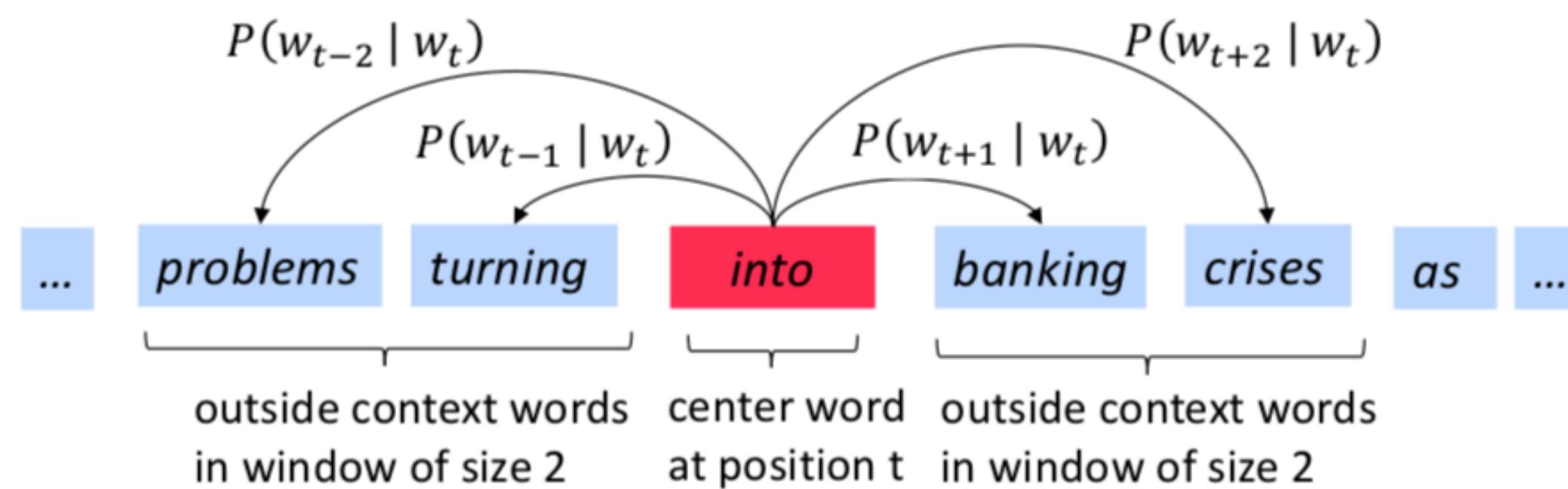
$P(\cdot | a)$  is a probability distribution defined over  $V$ :  
$$\sum_{w \in V} P(w | a) = 1$$



**Key trick: We turn unlabeled text into supervised learning data**

# Skip-gram: Intuition (1/2)

**Key trick: We turn unlabeled text into supervised learning data**

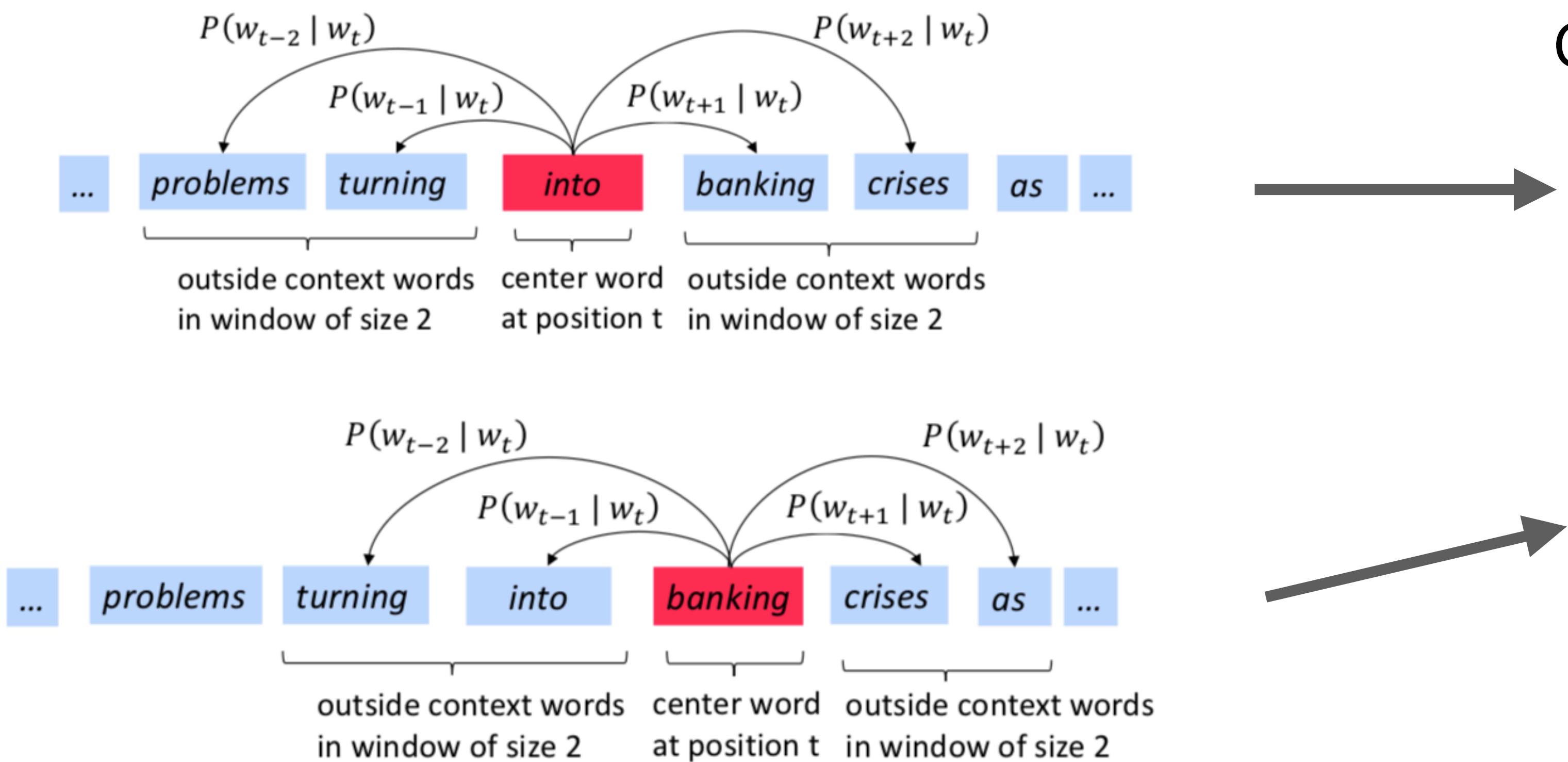


Convert into training data:

(into, problems)  
(into, turning)  
(into, banking)  
(into, crises)

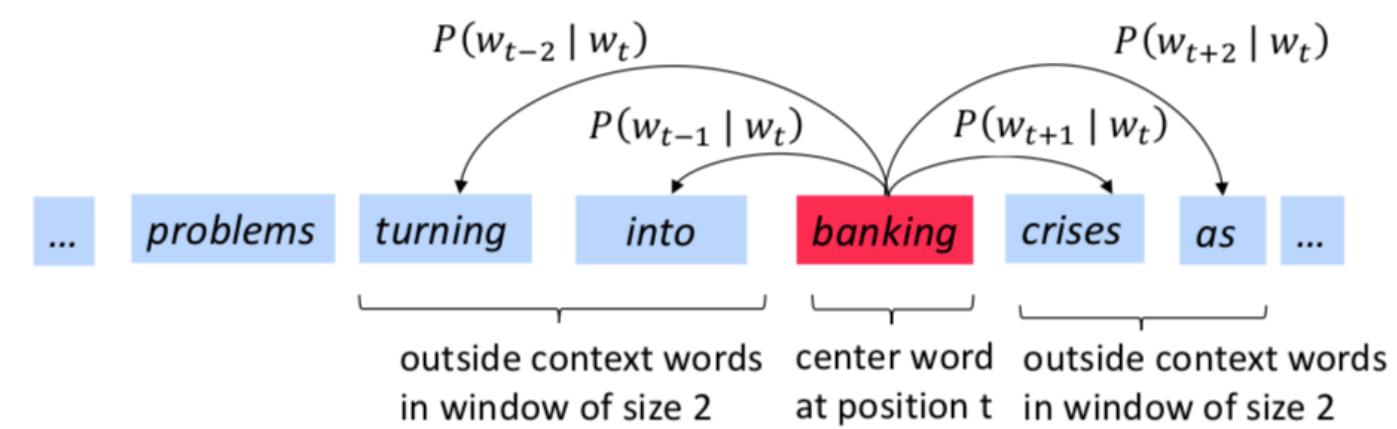
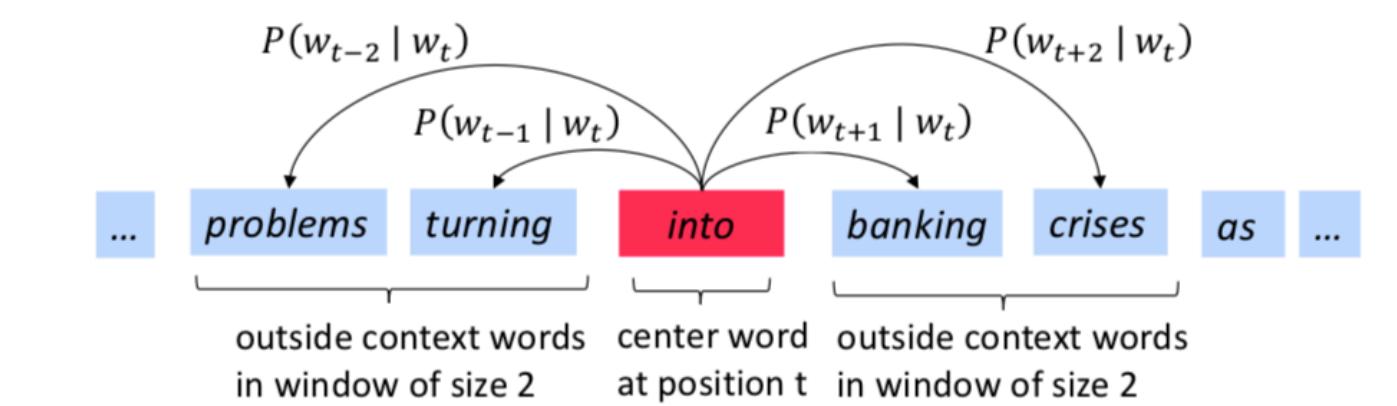
# Skip-gram: Intuition (1/2)

**Key trick: We turn unlabeled text into supervised learning data**



# Skip-gram: Intuition (2/2)

**Key trick: We turn unlabeled text into supervised learning data**



Convert into training data:

(into, problems)

(into, turning)

(into, banking)

(into, crises)

(banking, turning)

(banking, into)

(banking, crises)

(banking, as)

...

Our goal is to find parameters that can maximize

$$\frac{P(\text{problems} \mid \text{into}) \times P(\text{turning} \mid \text{into}) \times P(\text{banking} \mid \text{into}) \times P(\text{crises} \mid \text{into}) \times}{P(\text{turning} \mid \text{banking}) \times P(\text{into} \mid \text{banking}) \times P(\text{crises} \mid \text{banking}) \times P(\text{as} \mid \text{banking}) \dots}$$

# Skim-gram: Objective function

- For each position  $t = 1, 2, \dots, T$ , predict context words within context size  $m$ , given center word  $w_t$ :

$$\prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} \mid w_t; \theta)$$

# Skim-gram: Objective function

- For each position  $t = 1, 2, \dots, T$ , predict context words within context size  $m$ , given center word  $w_t$ :

$$\mathcal{L}(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} \mid w_t; \theta)$$

all the parameters to  
be optimized

# Skim-gram: Objective function

- For each position  $t = 1, 2, \dots, T$ , predict context words within context size  $m$ , given center word

$w_t$ :

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

all the parameters to  
be optimized

- It is equivalent to minimizing the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log \mathcal{L}(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} | w_t; \theta)$$

# How to define $P(w_{t+j} \mid w_t; \theta)$ ?

- Use two sets of vectors for each word in the vocabulary

$\mathbf{u}_a \in \mathbb{R}^d$  : vector for center word  $a, \forall a \in V$

$\mathbf{v}_b \in \mathbb{R}^d$  : vector for context word  $b, \forall b \in V$

- Use inner product  $\mathbf{u}_a \cdot \mathbf{v}_b$  to measure how likely word  $a$  appears with context word  $b$

$$P(w_{t+j} \mid w_t) = \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

Does this term  
seem familiar?

# ... vs multinomial logistic regression

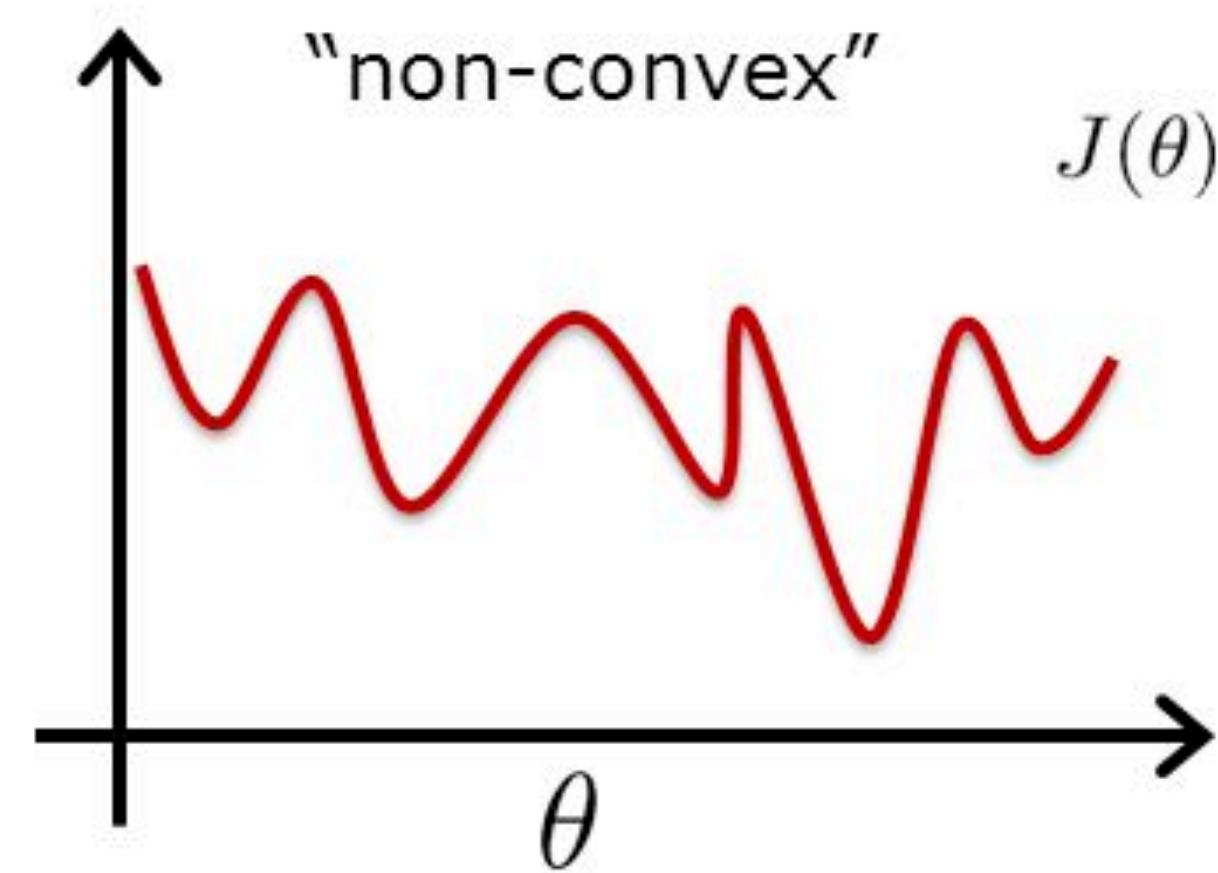
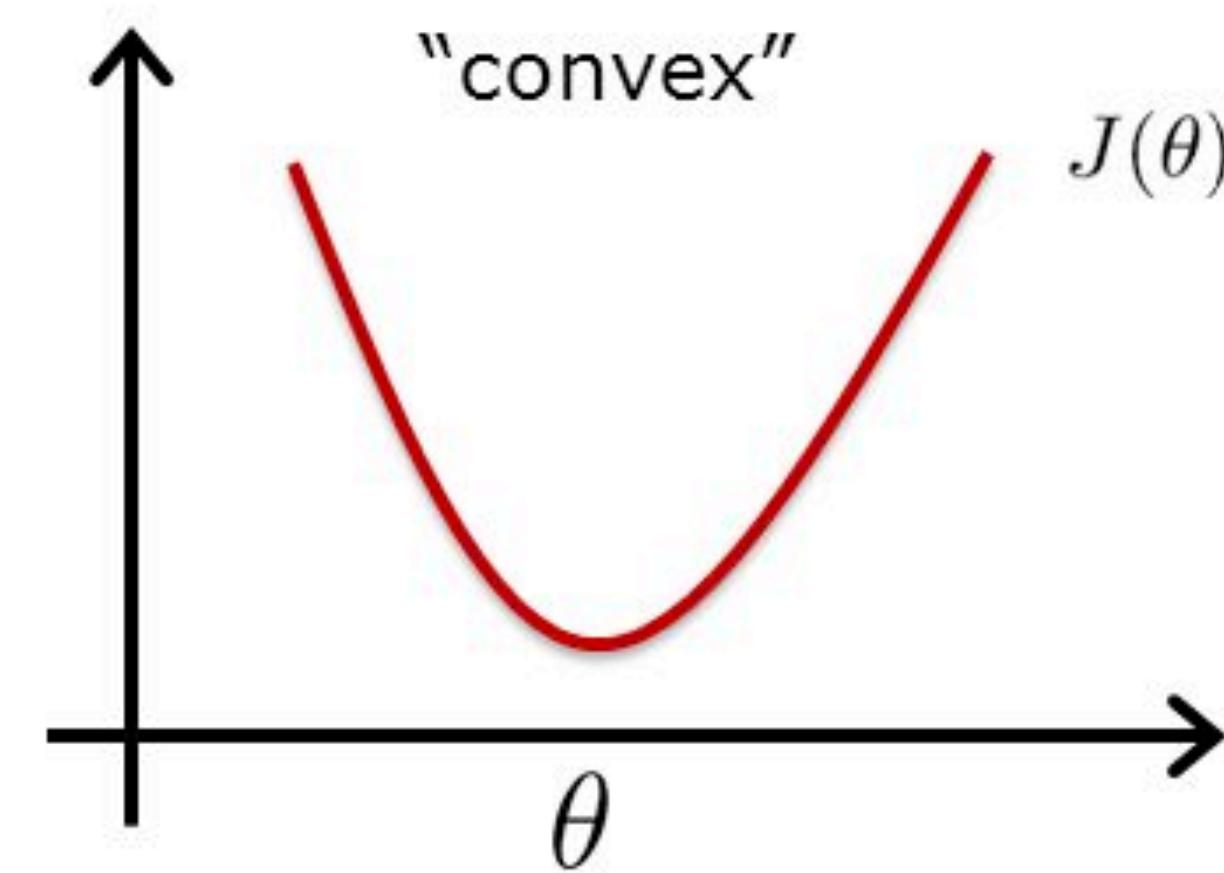
- Essentially a  $|V|$ -way classification problem
- Recall: multinomial logistic regression:

$$P(y = c \mid x) = \frac{\exp(\mathbf{w}_c \cdot \mathbf{x} + b_c)}{\sum_{j=1}^m \exp(\mathbf{w}_j \cdot \mathbf{x} + b_j)}$$

$$P(w_{t+j} \mid w_t) = \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

- If we fix  $\mathbf{u}_{w_t}$ , it is reduced to a multinomial logistic regression problem.
- However, since we have to learn both  $\mathbf{u}$  and  $\mathbf{v}$  together, the training objective is **non-convex**.

# ... vs multinomial logistic regression



- It is hard to find a global minimum
- But can still use stochastic gradient descent to optimize  $\theta$ :

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} J(\theta)$$

# Important note

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

- In this formulation, we don't care about the classification task itself like we do for the logistic regression model we saw previously.
- The key point is that the parameters used to optimize this training objective—when the training corpus is large enough—can give us very good representations of words (following the principle of **distributional hypothesis**)!

# Quick quiz

How many parameters does this model have (i.e. what is size of  $\theta$ )?

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

- (a)  $d|V|$
- (b)  $2d|V|$
- (c)  $2m|V|$
- (d)  $2md|V|$

$V :=$  Vocabulary  
 $d :=$  dimension of embedding  
 $m :=$  size of context window

The answer is (b).

Each word has two d-dimensional vectors, so it is  $2 \times |V| \times d$ .

# Word2vec formulation

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

Q: Why do we need two vectors for each word?

- Because one word is not likely to appear in its own context window, e.g.,  $P(\text{dog} \mid \text{dog})$  should be low. If we use one set of vectors only, it essentially needs to minimize  $\mathbf{u}_{\text{dog}} \cdot \mathbf{u}_{\text{dog}}$

Q: Which set of vectors are used as word embeddings?

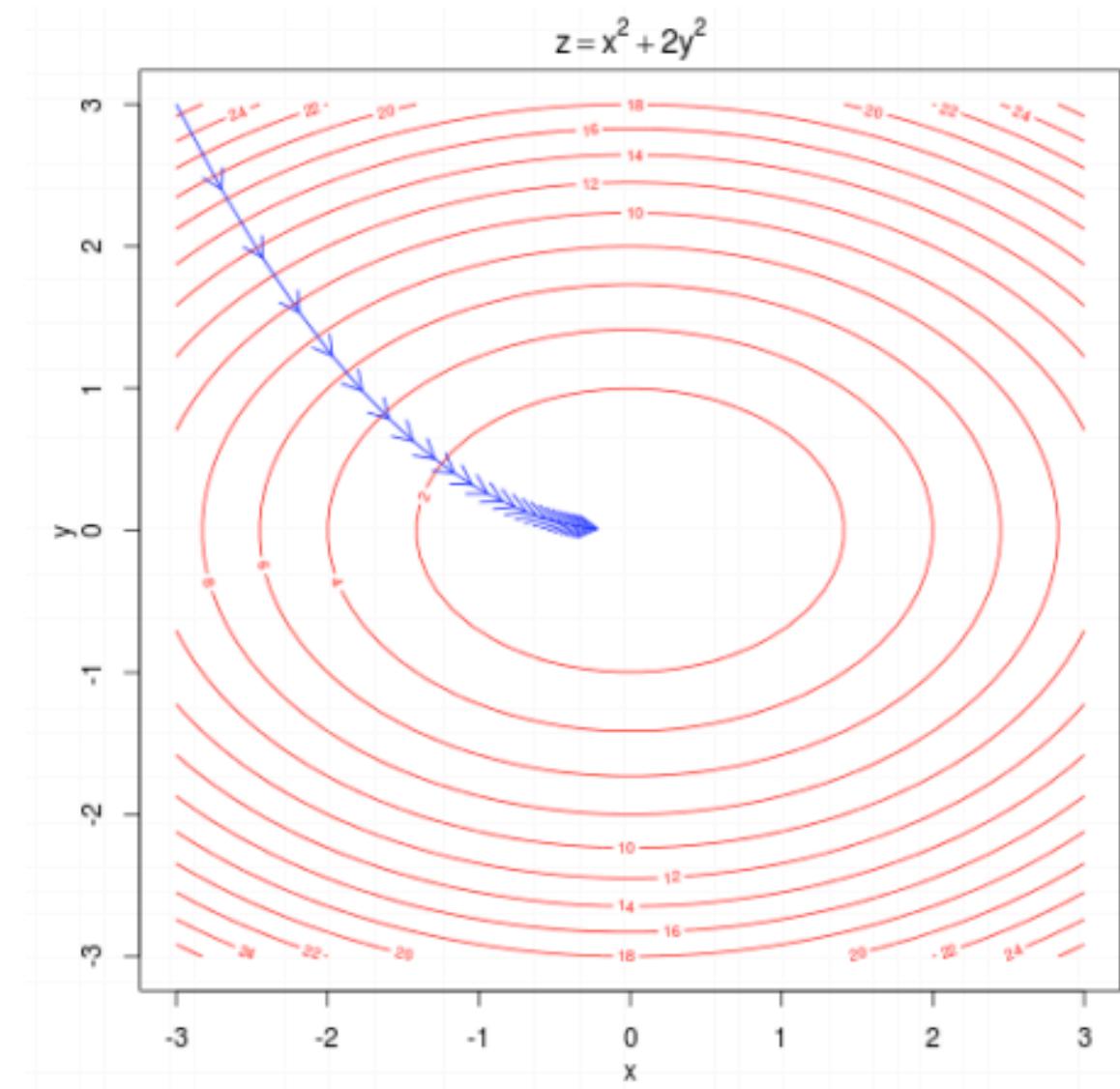
- This is an empirical question. Typically just  $\mathbf{u}_w$ , but you can also concatenate the two vectors..

# How to train this model?

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

- To train such a model, we need to compute the vector gradient  $\nabla_{\theta} J(\theta) = ?$
- Remember that  $\theta$  represents all  $2d|V|$  model parameters, in one vector.

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix}$$



# Warm up: Vectorized gradients

$$f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{a}$$
$$\mathbf{x}, \mathbf{a} \in \mathbb{R}^n$$

$$\frac{\partial f}{\partial \mathbf{x}} = \mathbf{a}$$

$$f = x_1 a_1 + x_2 a_2 + \dots + x_n a_n$$

$$\frac{\partial f}{\partial \mathbf{x}} = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

# Warm up: Vectorized gradients exercises

Let  $f = \exp(\mathbf{w} \cdot \mathbf{x})$ , what is the value of  $\frac{\partial f}{\partial \mathbf{x}}$ ? (Assume  $\mathbf{w}, \mathbf{x} \in \mathbb{R}^n$ )

- (a)  $\mathbf{w}$
- (b)  $\exp(\mathbf{w} \cdot \mathbf{x})$
- (c)  $\exp(\mathbf{w} \cdot \mathbf{x})\mathbf{w}$
- (d)  $\mathbf{x}$

Hints:

- The derivative of  $\exp(\mathbf{z})$  is  $\exp(\mathbf{z})$
- The derivative of  $\mathbf{w} \cdot \mathbf{x}$  with respect to  $\mathbf{x}$  is  $\mathbf{w}$

The answer is (c).

$$\frac{\partial}{\partial x_i} = \frac{\exp(\sum_{k=1}^n w_i x_i)}{\partial x_i} = \exp(\sum_{k=1}^n w_i x_i) w_i$$

# Let's compute gradients for word2vec

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

Consider one pair of center/context words  $(t, c)$ :

$$y = -\log \left( \frac{\exp(\mathbf{u}_t \cdot \mathbf{v}_c)}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)} \right)$$

We need to compute the gradient of  $y$  with respect to

$$\mathbf{u}_t \text{ and } \mathbf{v}_k, \forall k \in V$$

# Overall algorithm

- Input: text corpus, embedding size  $d$ , vocabulary  $V$ , context size  $m$
- Initialize  $\mathbf{u}_i, \mathbf{v}_i$  randomly  $\forall i \in V$
- Run through the training corpus and for each training instance  $(t, c)$ :
  - Update  $\mathbf{u}_t \leftarrow \mathbf{u}_t - \eta \frac{\partial y}{\partial \mathbf{u}_t}; \quad \frac{\partial y}{\partial \mathbf{u}_t} = -\mathbf{v}_c + \sum_{k \in V} P(k | t) \mathbf{v}_k$
  - Update  $\mathbf{v}_k \leftarrow \mathbf{v}_k - \eta \frac{\partial y}{\partial \mathbf{v}_k}, \forall k \in V; \quad \frac{\partial y}{\partial \mathbf{v}_k} = \begin{cases} (P(k | t) - 1) \mathbf{u}_t & k = c \\ P(k | t) \mathbf{u}_t & k \neq c \end{cases}$

Q: Can you think of any issues with this algorithm?

# Overall algorithm: Problem

**Problem:** every time you get one pair of  $(t, c)$ , you need to update  $v_k$  with all the words in the vocabulary! This is very expensive computationally.

$$\frac{\partial y}{\partial \mathbf{u}_t} = -\mathbf{v}_c + \sum_{k \in V} P(k|t) \mathbf{v}_k \quad ; \quad \frac{\partial y}{\partial \mathbf{v}_k} = \begin{cases} (P(k|t) - 1) \mathbf{u}_t & k = c \\ P(k|t) \mathbf{u}_t & k \neq c \end{cases}$$

$$P(k|t) = \frac{\exp(\mathbf{u}_t \cdot \mathbf{v}_k)}{\sum_{j \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_j)}$$

**Key question:** Do we really need to compare the center word against every word in the vocabulary every time?

# Skip-gram with negative sampling (SGNS)

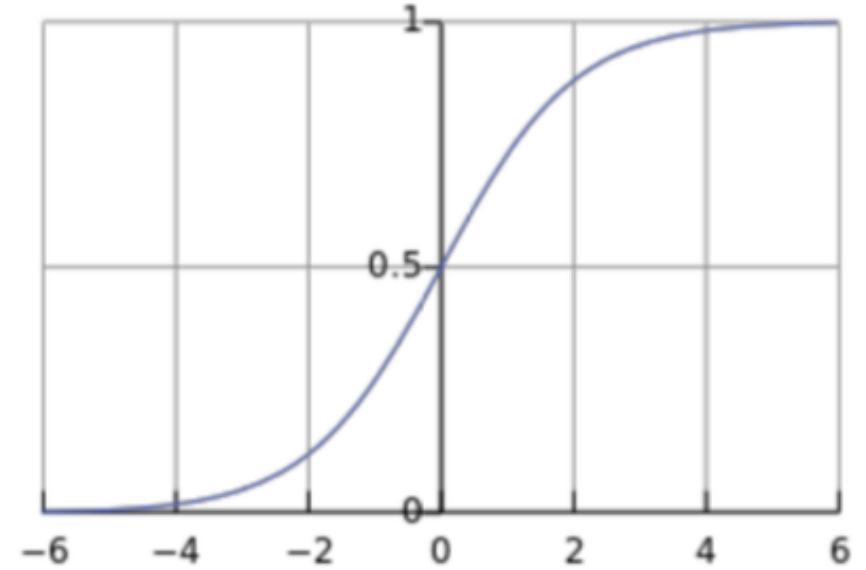
**Negative sampling:** instead of considering all the words in  $V$ , let's randomly sample  $K$  (5-20) negative examples.

Softmax:

$$y = -\log \left( \frac{\exp(\mathbf{u}_t \cdot \mathbf{v}_c)}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)} \right)$$

Negative sampling:  $y = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{j \sim P(w)} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_j))$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



# Skip-gram with negative sampling (SGNS)

**Key idea: Convert the  $|V|$ -way classification into a set of binary classification tasks.**

Every time we get a pair of words  $(t, c)$ , we don't predict  $c$  among all the words in the vocabulary. Instead, we predict  $(t, c)$  is a positive pair, and  $(t, c')$  is a negative pair for a small number of sampled  $c'$ .

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

negative examples -			
t	c	t	c
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

$$y = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{j \sim P(w)} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_j))$$

$P(w)$ : sampling according to  
the frequency of words

Similar to **binary logistic regression**, but we need to optimize  $\mathbf{u}$  and  $\mathbf{v}$  together.

$$P(y = 1 \mid t, c) = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c) \quad p(y = 0 \mid t, c') = 1 - \sigma(\mathbf{u}_t \cdot \mathbf{v}_{c'}) = \sigma(-\mathbf{u}_t \cdot \mathbf{v}_{c'})$$

# Understanding SGNS

In skip-gram with negative sampling (SGNS), how many parameters need to be updated in  $\theta$  for every  $(t, c)$  pair?

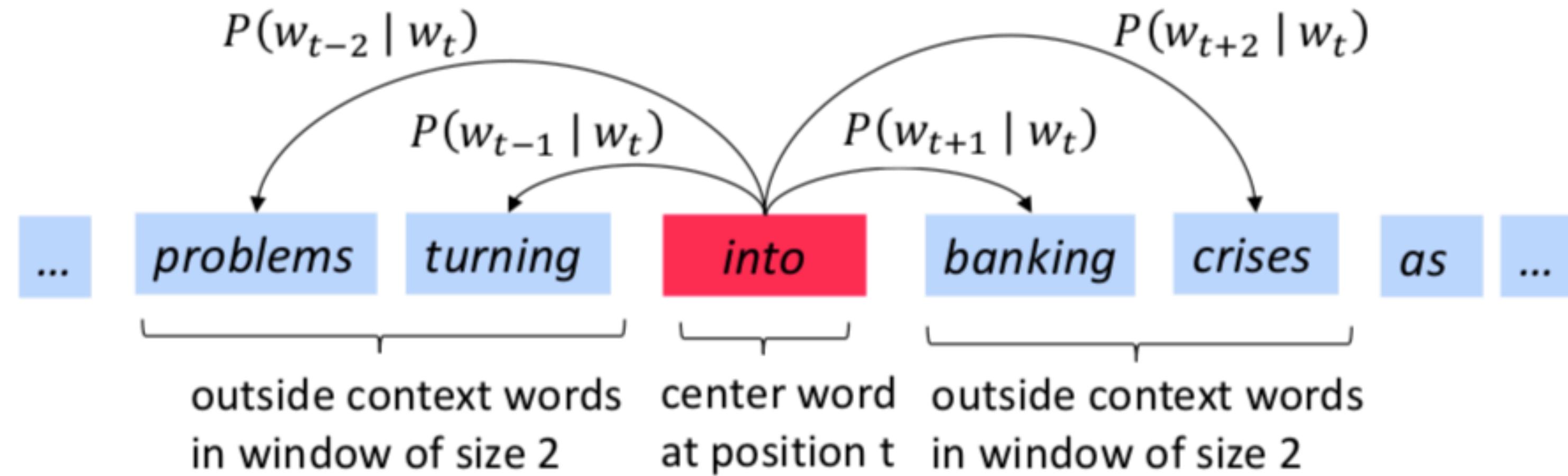
$$y = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^K \mathbb{E}_{j \sim P(w)} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_j))$$

- (a)  $Kd$
- (b)  $2Kd$
- (c)  $(K + 1)d$
- (d)  $(K + 2)d$

The answer is (d).

We need to calculate gradients with respect to  $\mathbf{u}_t$  and  $(K + 1) \mathbf{v}_i$  (one positive and K negatives).

# Skip-gram: Summary (1/2)



**Key trick: We turn unlabeled text into supervised learning data**

$$P(w_{t+j} | w_t) = \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

# Skip-gram: Summary (2/2)

- The original skip-gram objective:

$$P(w_{t+j} | w_t) = \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

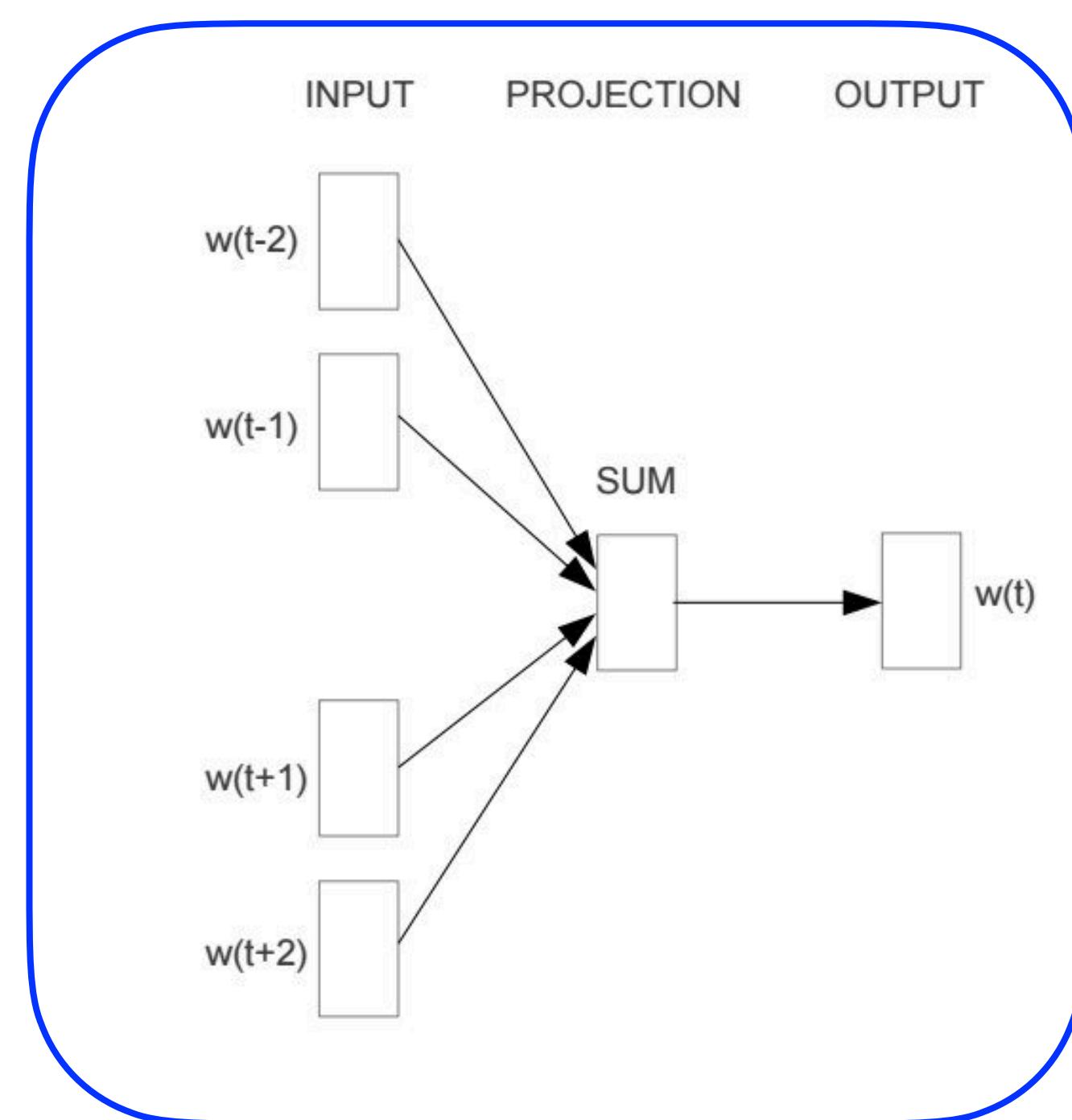
- Maximize:  $\mathcal{L}(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t; \theta)$
- Equivalent to minimize  $-\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} | w_t; \theta)$
- Elegant but computationally infeasible
- The problem comes entirely from the softmax normalization  $\rightarrow$  Negative sampling! Empirically, this approximation works extremely well

# Word2vec

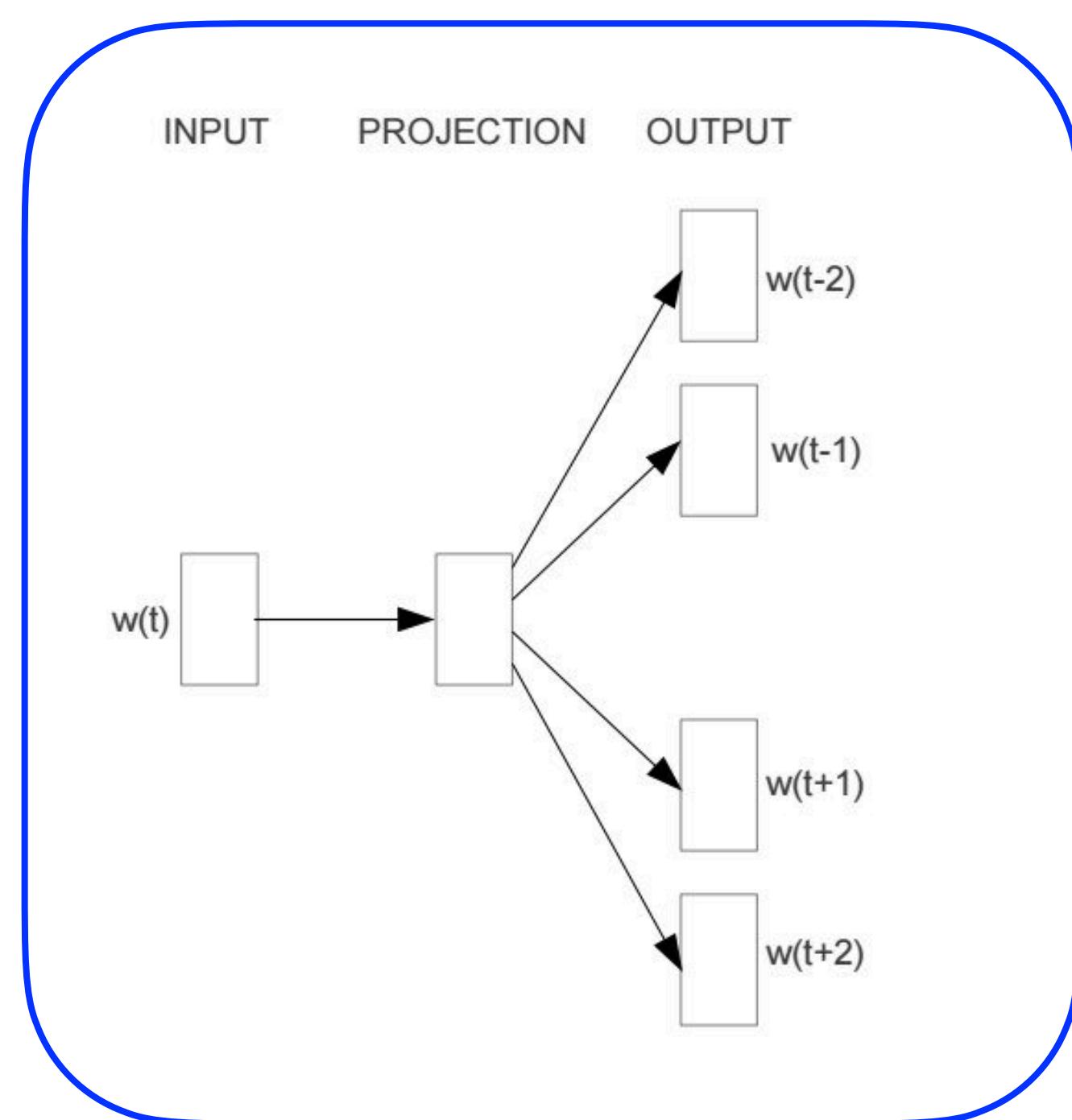
- (Mikolov et al 2013a): Efficient Estimation of Word Representations in Vector Space
- (Mikolov et al 2013b): Distributed Representations of Words and Phrases and their Compositionality



Tomáš Mikolov



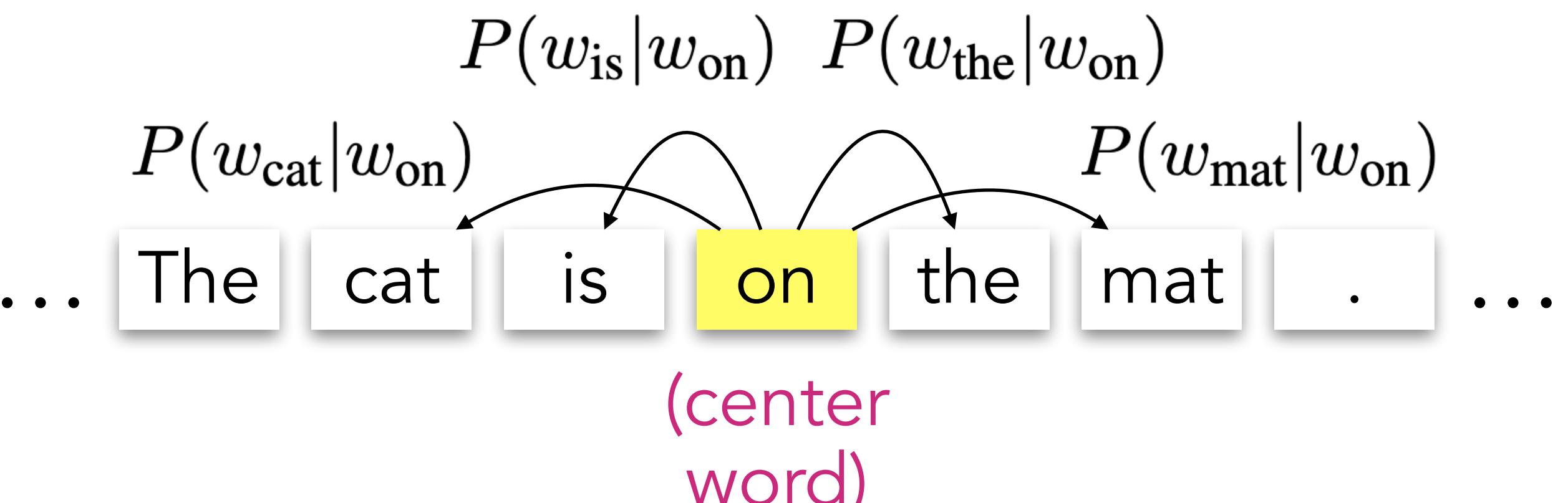
Continuous Bag of Words (CBOW)



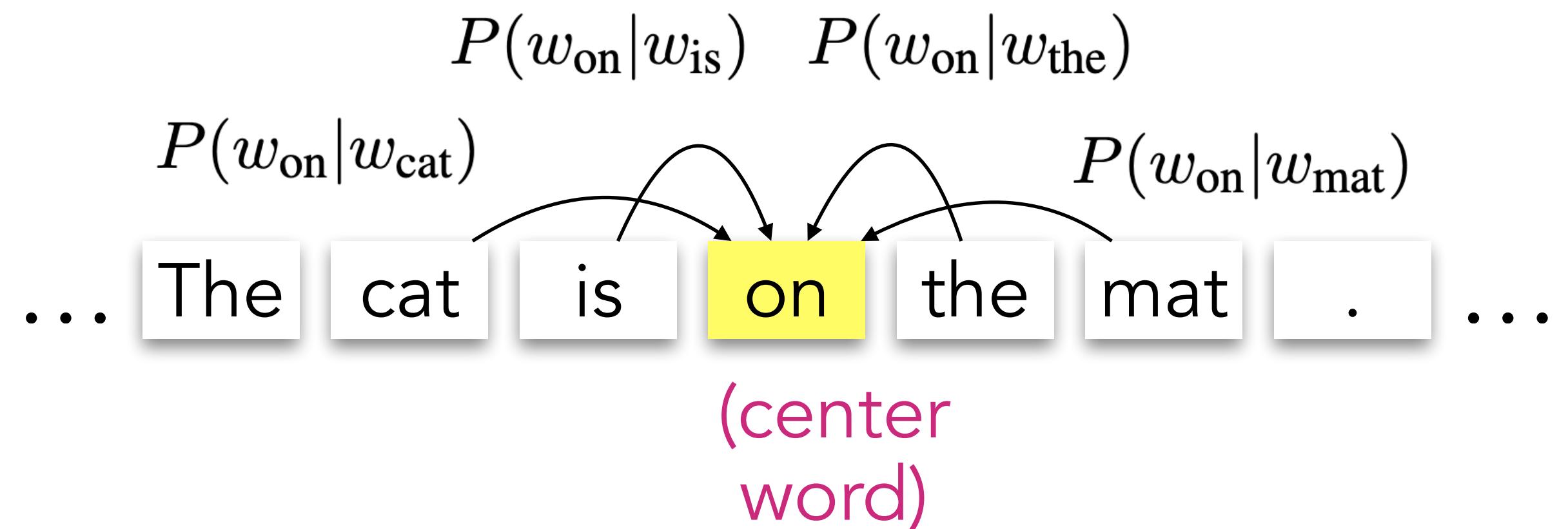
Skip-gram

# Skip-gram vs. CBOW

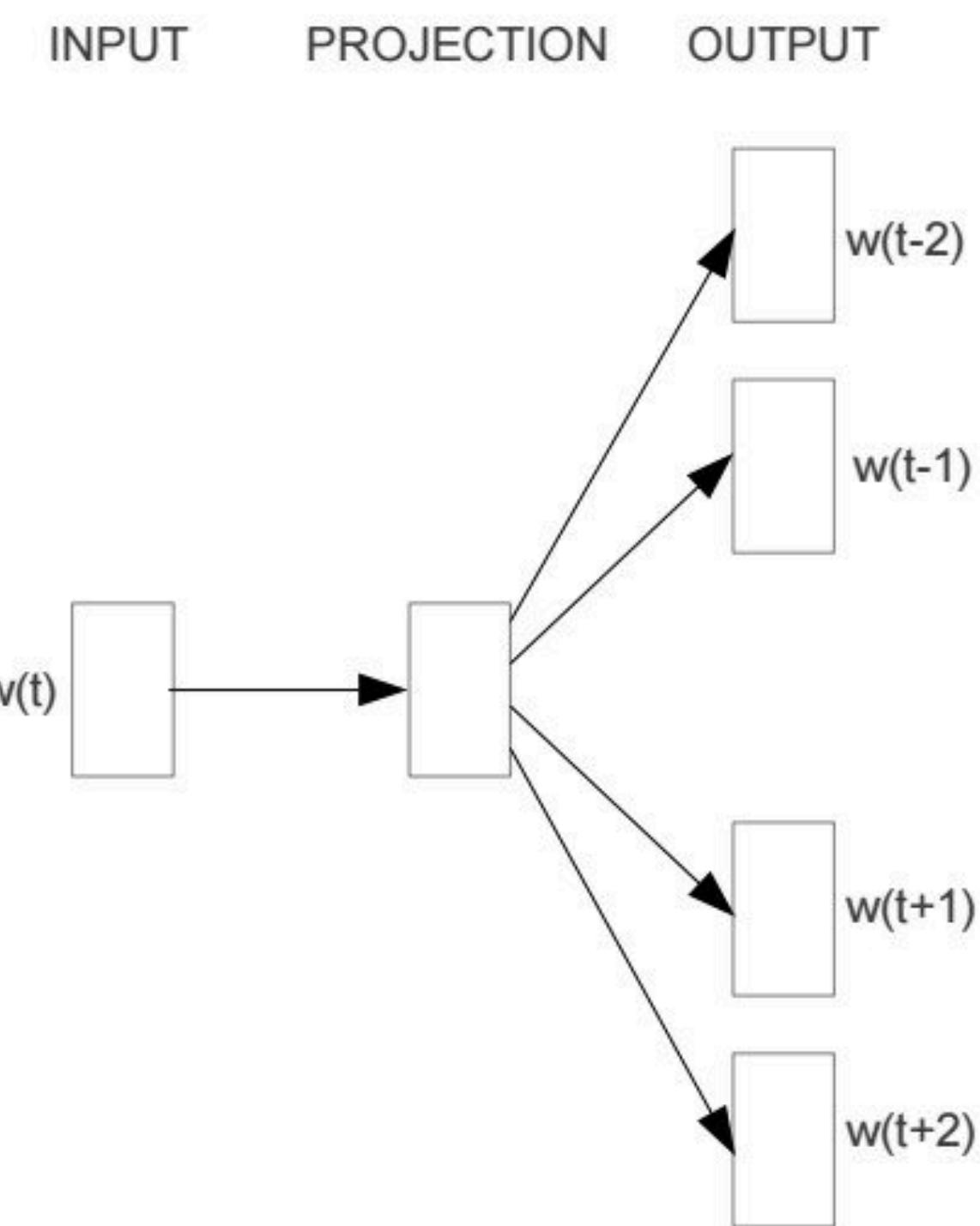
Skip-gram: Trained to predict the context words given a target word.



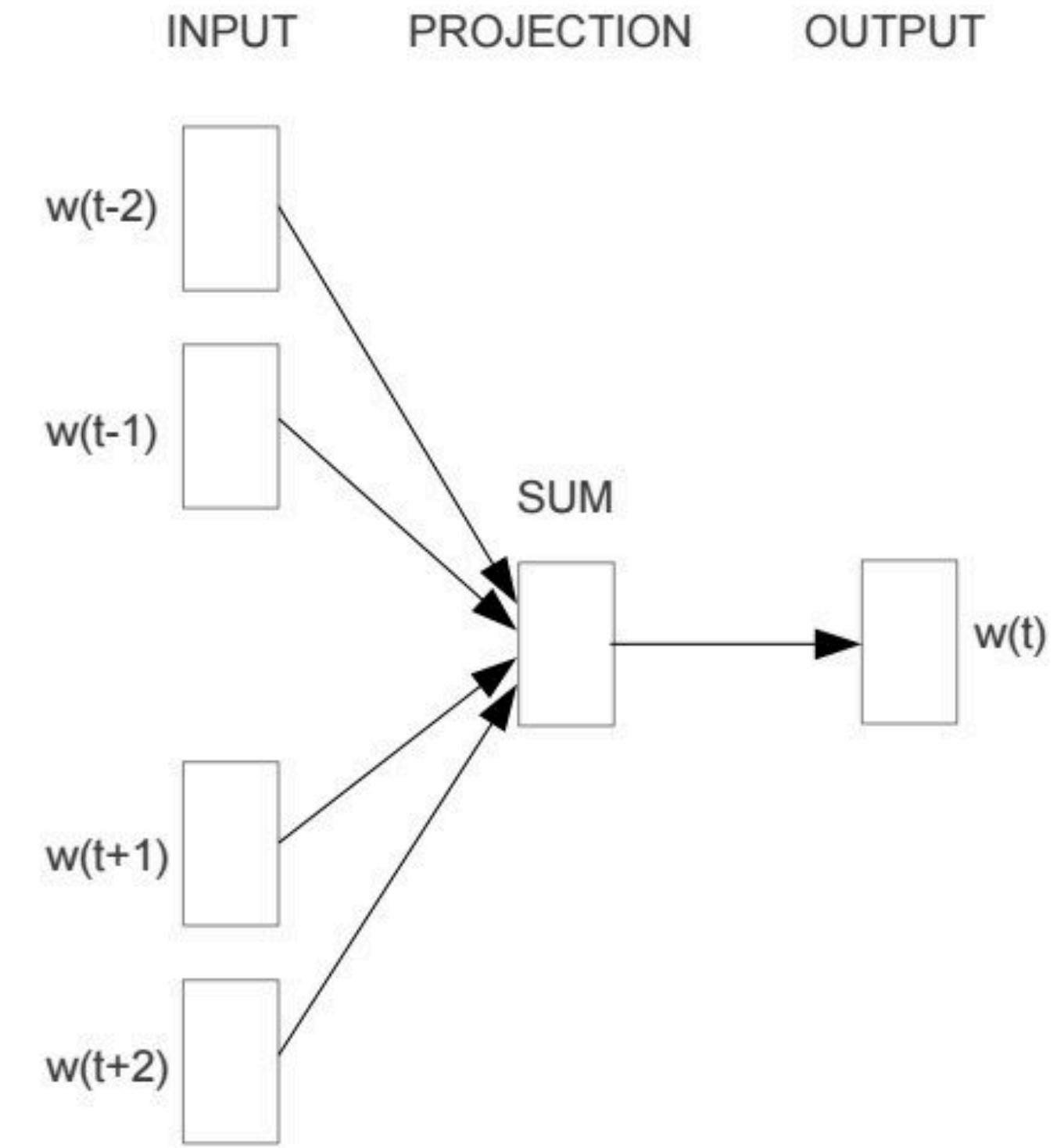
CBOW: Trained to predict the target word given its context words.



# Skip-gram vs. CBOW



Skip-gram



Continuous Bag of Words (CBOW)

# CBOW Objective

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

=

$$P(w_t | w_{t-m}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+m})$$

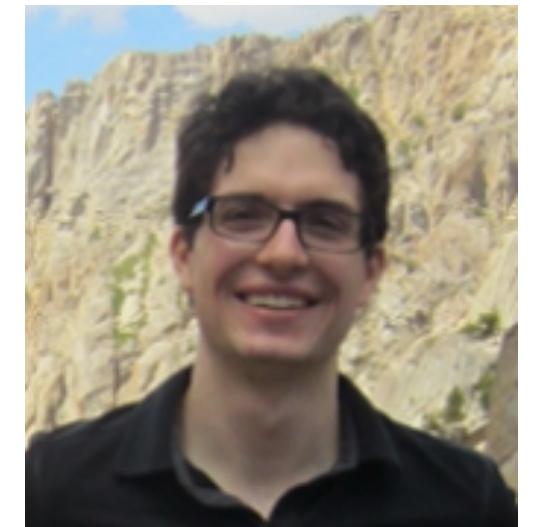
Skim-gram used:  $P(w_{t+j}|w_t) = \frac{\exp(\mathbf{v}_{w_{t+j}}^\top \mathbf{v}_{w_t})}{\sum_{w' \in \mathcal{V}} \exp(\mathbf{v}_{w'^\top} \mathbf{v}_{w_t})}$

where

$$\bar{\mathbf{v}}_c = \frac{1}{2m} \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{w_{t+j}}$$

# Other word embeddings?

# GloVe: Global Vectors



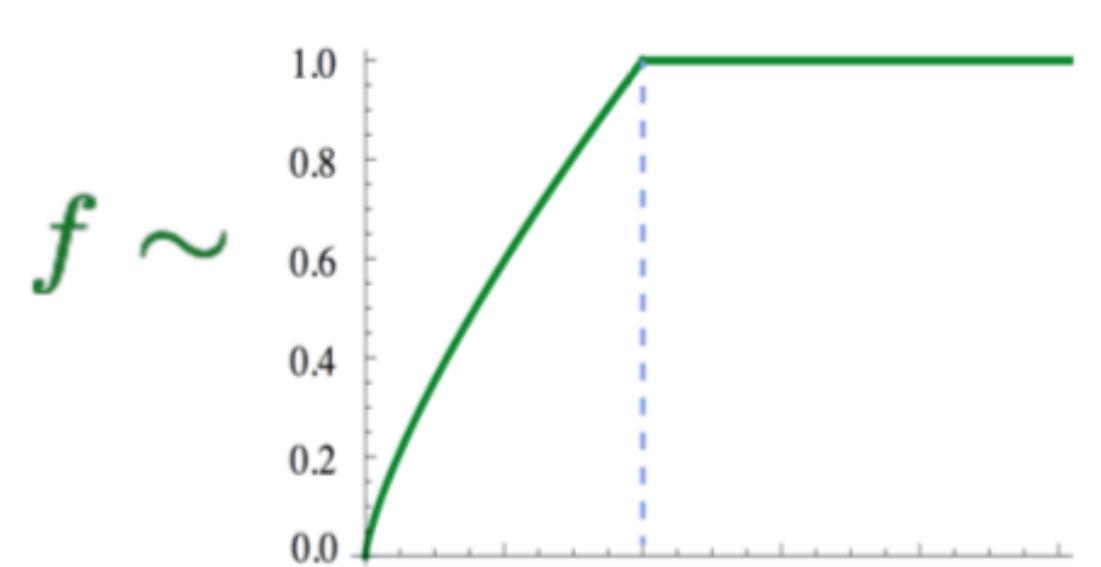
- Take the global co-occurrence statistics:  $X_{i,j}$
- Key idea: let's approximate  $\mathbf{u}_i \cdot \mathbf{v}_j$  using their co-occurrence counts directly

$$J(\theta) = \sum_{i,j \in V} f(X_{i,j}) \left( \mathbf{u}_i \cdot \mathbf{v}_j + b_i + \tilde{b}_j - \log X_{i,j} \right)$$

- $f$ : Weighting function, want to give more importance to more common pairs, but capped at a certain point

Advantages:

- Training faster
- Scalable to very large corpora



# FastText: Subword Embeddings

- Similar to Skip-gram, but break words into n-grams with  $n = 3$  to  $6$

where: 3-grams: <wh, whe, her, ere, re>

4-grams: <whe, wher, here, ere>

5-grams: <wher, where, here>

6-grams: <where, where>

- Replace  $\mathbf{u}_i \cdot \mathbf{v}_j$  by  $\sum_{g \in n\text{-grams}(w_i)} \mathbf{u}_g \cdot \mathbf{v}_j$

# Trained word embeddings available

- word2vec: <https://code.google.com/archive/p/word2vec/>
- GloVe: <https://nlp.stanford.edu/projects/glove/>
- FastText: <https://fasttext.cc/>

## Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License v1.0](#) whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
  - [Wikipedia 2014](#) + [Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
  - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
  - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
  - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data

Differ in algorithms, text corpora, dimensions, cased/uncased...  
Applied to many other languages

# Easy to use!

```
from gensim.models import KeyedVectors
# Load vectors directly from the file
model = KeyedVectors.load_word2vec_format('data/GoogleGoogleNews-vectors-negative300.bin', binary=True)
# Access vectors for specific words with a keyed lookup:
vector = model['easy']
```

---

```
In [17]: model.similarity('straightforward', 'easy')
Out[17]: 0.5717043285477517

In [18]: model.similarity('simple', 'impossible')
Out[18]: 0.29156160264633707

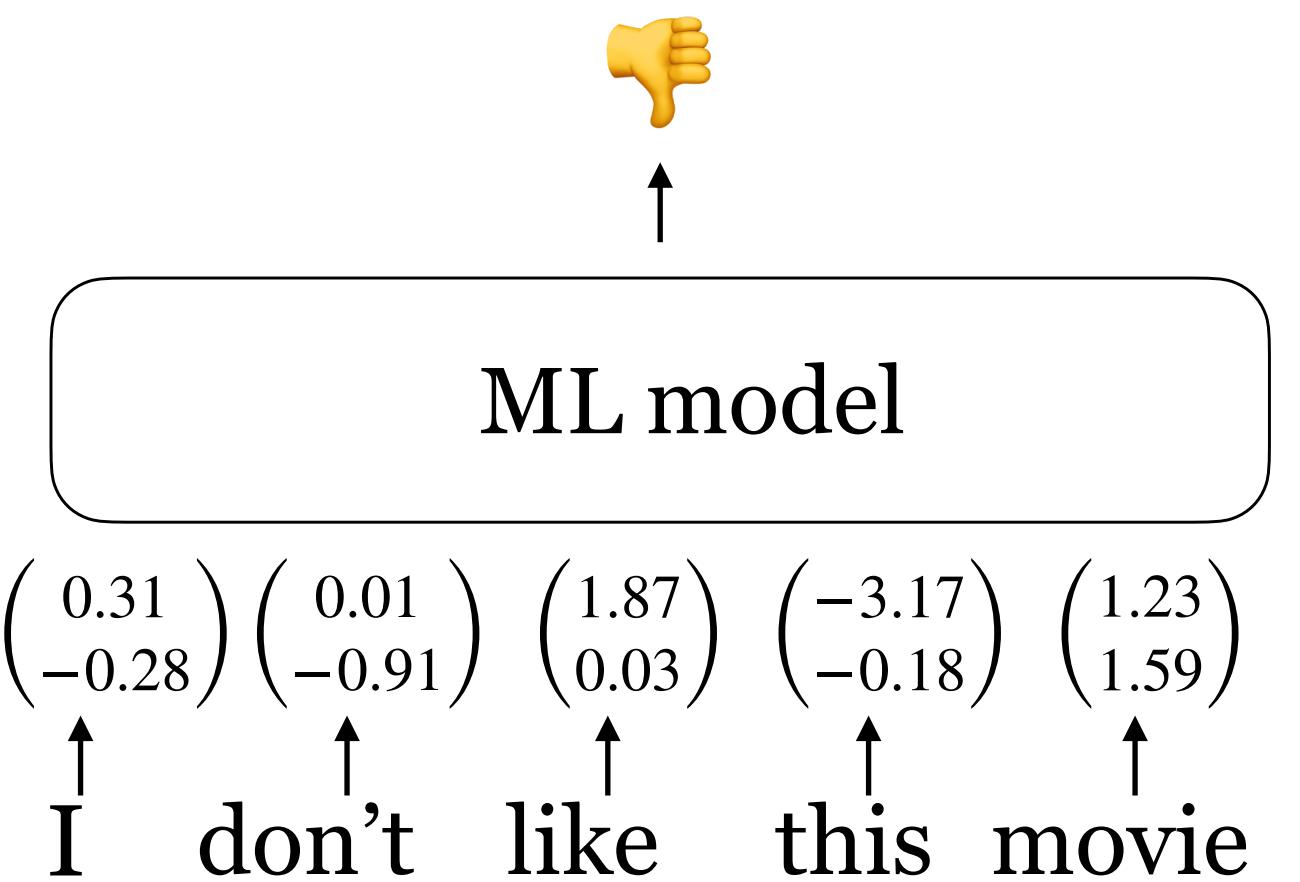
In [19]: model.most_similar('simple')
Out[19]: [('straightforward', 0.7460169196128845),
          ('Simple', 0.7108174562454224),
          ('uncomplicated', 0.6297484636306763),
          ('simplest', 0.6171397566795349),
          ('easy', 0.5990299582481384),
          ('fairly_straightforward', 0.5893306732177734),
          ('deceptively_simple', 0.5743066072463989),
          ('simpler', 0.5537199378013611),
          ('simplistic', 0.5516539216041565),
          ('disarmingly_simple', 0.5365327000617981)]
```

# Evaluating word embeddings

# Extrinsic vs. intrinsic evaluation

- **Extrinsic evaluation**

- Let's plug these word embeddings into a real NLP system and see whether this improves performance
- Could take a long time but still the most important evaluation metric



# Extrinsic vs. intrinsic evaluation



ML model

$$\begin{pmatrix} 0.31 \\ -0.28 \end{pmatrix} \begin{pmatrix} 0.01 \\ -0.91 \end{pmatrix} \begin{pmatrix} 1.87 \\ 0.03 \end{pmatrix} \begin{pmatrix} -3.17 \\ -0.18 \end{pmatrix} \begin{pmatrix} 1.23 \\ 1.59 \end{pmatrix}$$

↑      ↑      ↑      ↑      ↑  
I    don't    like    this    movie

A straightforward solution: given an input sentence  $x_1, x_2, \dots, x_n$

Instead of using a bag-of-words model, we can compute  $\text{vec}(x) = \mathbf{e}(x_1) + \mathbf{e}(x_2) + \dots + \mathbf{e}(x_n)$

And then train a logistic regression classifier on  $\text{vec}(x)$  as we did before!

Note: There are much better ways to do this e.g., take word embeddings as input of neural networks

# Extrinsic vs. intrinsic evaluation

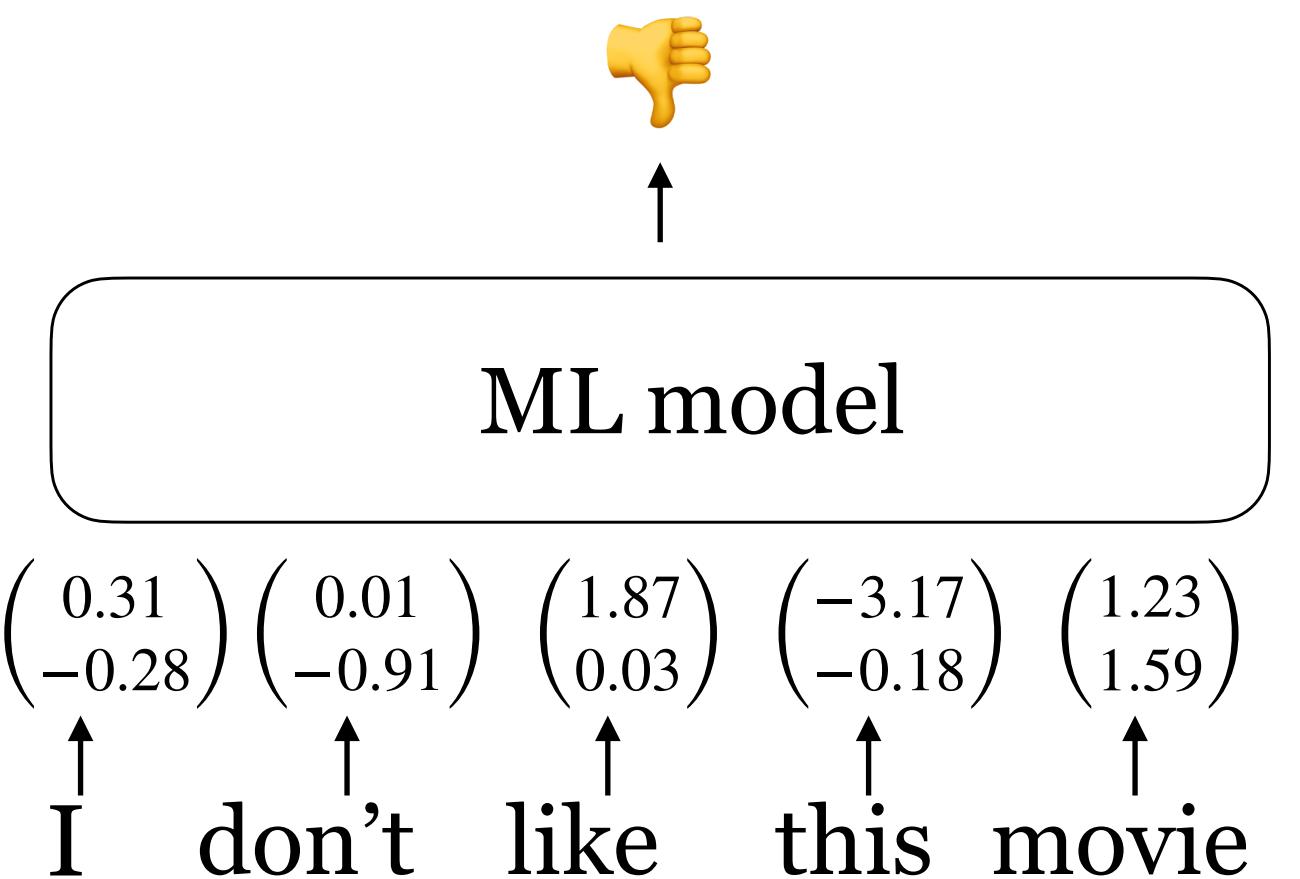
- **Extrinsic evaluation**

- Let's plug these word embeddings into a real NLP system and see whether this improves performance

- Could take a long time but still the most important evaluation metric

- **Intrinsic evaluation**

- Evaluate on a specific/intermediate subtask
- Fast to compute
- Not clear if it really helps downstream tasks



# Intrinsic evaluation: word similarity

## Word similarity

Example dataset: wordsim-353: 353 pairs of words with human judgement

<http://www.cs.technion.ac.il/~gabr/resources/data/wordsim353/>

Word 1	Word 2	Human (mean)
tiger	cat	7.35
tiger	tiger	10
book	paper	7.46
computer	internet	7.58
plane	car	5.77
professor	doctor	6.62
stock	phone	1.62
stock	CD	1.31
stock	jaguar	0.92

$$\cos(\mathbf{u}_i, \mathbf{u}_j) = \frac{\mathbf{u}_i \cdot \mathbf{u}_j}{\|\mathbf{u}_i\|_2 \times \|\mathbf{u}_j\|_2}.$$

Metric: Spearman rank correlation

# Intrinsic evaluation: word similarity

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW <sup>†</sup>	6B	57.2	65.6	68.2	57.0	32.5
SG <sup>†</sup>	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

SG: Skip-gram

# Intrinsic evaluation: word analogy

Word analogy test:  $a : a^* :: b : b^*$

$$b^* = \arg \max_{w \in V} \cos(e(w), e(a^*) - e(a) + e(b))$$

semantic

syntactic

Chicago:Illinois $\approx$ Philadelphia: ?

bad:worst  $\approx$  cool: ?

More examples at

<http://download.tensorflow.org/data/questions-words.txt>

Metric: accuracy

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW <sup>†</sup>	300	6B	63.6	<u>67.4</u>	65.7
SG <sup>†</sup>	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

# OK, and how do we actually use these?

Next lecture!

# Questions?

# Acknowledgement

**Princeton COS 484 by Danqi Chen, Tri Dao, Vikram Ramaswamy**

Liang Huang AI 534 (400/401), Machine Learning

Stanford CS224N NLP with Deep Learning by Diyi Yang and Tatsunori Hashimoto (+ older versions by Chris Manning)

New York University NLP by He He

# FAQ

- **Enrollment related**
  - **I was already enrolled but got informed that I was removed because my enrollment was due to a system error. What happened?** Sorry about that. We were informed that you did not meet the eligibility criteria and were enrolled due to a system error, and the system automatically corrected it. We do not have control over this process. Please join the waitlist via the Google Form on the website. Qualified students will be selected using the same criteria.
  - **Does the time I join the waitlist matter?** It does not matter until the end of today. Students added to the waitlist starting tomorrow will not be considered.
  - **When will I hear back? What are my chances?** Sorry, we do not know until the 5th class. The capacity increase is made by the department, not the course staff. Most students have not received enrollment codes yet—we only sent a few to test the system and found that the codes are currently not working. We are waiting for the department to resolve this.
  - **I received an enrollment code today, but I can only waitlist.** Apologies for the confusion. The department needs to increase the course capacity before enrollment codes can be used. We will notify you once codes become active.
- **Ed/Gradescope access related:** Please contact our GSI, Zineng Tang (email address on the website).
- **Slides/recording availability:** Slides will be uploaded within 24 hours after each class. Recordings will be available starting next week—please stay tuned.