

# **Text Representation and Learning**

Information Retrieval and Extraction

---

**Rahul Mishra**  
IIIT-Hyderabad

INTERNATIONAL INSTITUTE OF  
INFORMATION TECHNOLOGY

HYDERABAD

# **Why do we need to represent the text?**

---

**Numerical Input Requirement**

**Feature Extraction**

**Semantic Understanding**

**Generalization**

**Dimensionality Reduction**

**Vectorization**

## Representing words as discrete symbols

---

In traditional NLP, we regard words as discrete symbols:  
hotel, conference, motel – a localist representation

Means one 1, the rest 0s

Such symbols for words can be represented by one-hot vectors:

motel = [0 0 0 0 0 0 0 0 0 1 0 0 0]

hotel = [0 0 0 0 0 0 1 0 0 0 0 0 0]

Vector dimension = number of words in vocabulary (e.g., 500,000)

## Problem with words as discrete symbols

---

**Example:** in web search, if user searches for “Seattle motel”, we would like to match documents containing “Seattle hotel”

But:

$$\begin{aligned}\text{motel} &= [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0] \\ \text{hotel} &= [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]\end{aligned}$$

These two vectors are orthogonal

There is no natural notion of similarity for one-hot vectors!

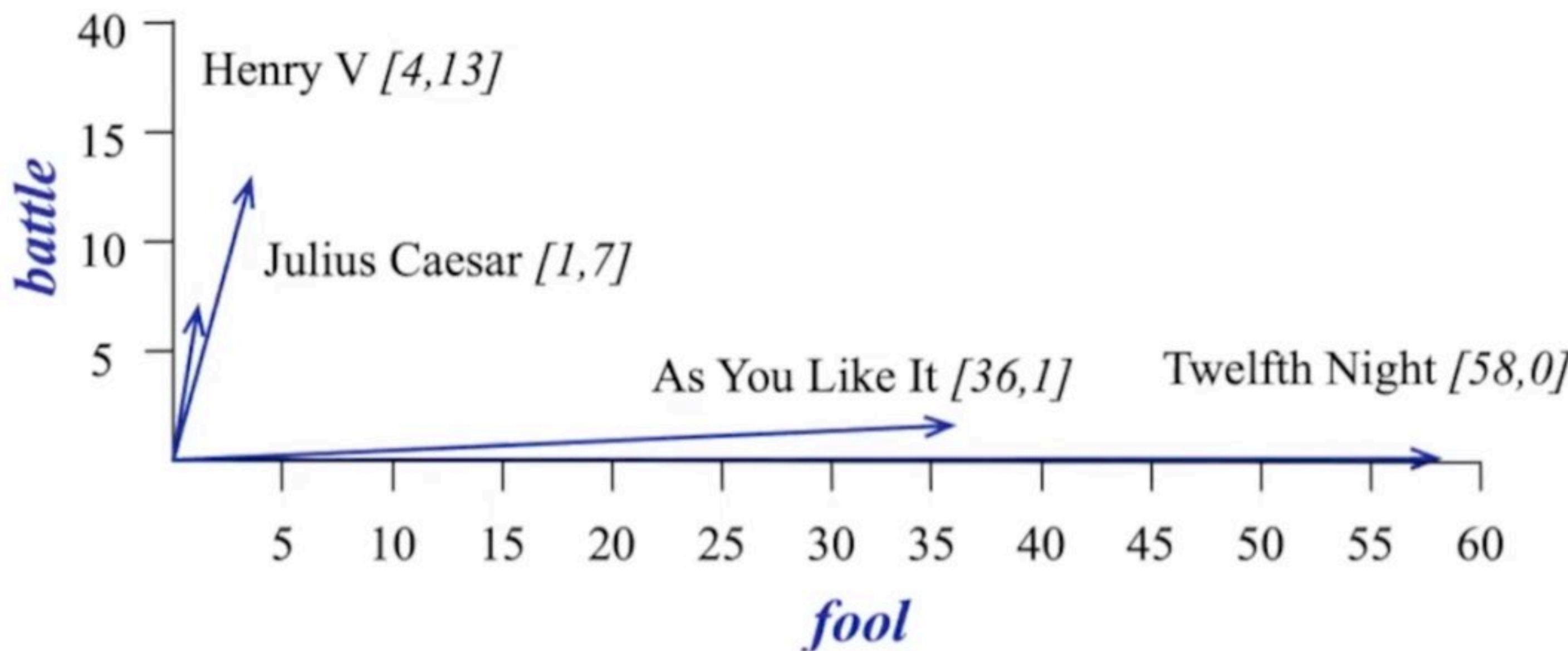
**Solution:**

- Could try to rely on WordNet’s list of synonyms to get similarity?
  - But it is well-known to fail badly: incompleteness, etc.
- **Instead: learn to encode similarity in the vectors themselves**

## Count-based Methods

	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

A term-document matrix



Vectors are similar for the two comedies

But comedies are different than the other two

## What about word vectors?

---

	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

*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"

## Another alternative

---

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

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	...

## Problem with raw count based methods

---

- The co-occurrence matrices we have seen represent each cell by word frequencies.
- Frequency is clearly useful; if *sugar* appears a lot near *apricot*, that's useful information.
- But overly frequent words like *the*, *it*, or *they* are not very informative about the context
- It's a paradox! How can we balance these two conflicting constraints?

## Two common solutions for the word weighting

---

**tf-idf:** tf-idf value for word t in document d:

$$w_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$$

Words like "the" or "it" have very low idf

**PMI:** (Pointwise mutual information)

- $\text{PMI}(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$

See if words like "good" appear more often with "great" than we would expect by chance

## Term Frequency (tf) and Document Frequency (df)

---

Term frequency (tf)

$$tf_{t,d} = \text{count}(t,d)$$

Instead of using raw count, we squash a bit:

$$tf_{t,d} = \log_{10}(\text{count}(t,d)+1)$$

Document frequency (df)

$df_t$  is the number of documents  $t$  occurs in.  
(note this is not collection frequency: total count across all documents)

"Romeo" is very distinctive for one Shakespeare play:

	Collection Frequency	Document Frequency
Romeo	113	1
action	113	31

## Inverse Document Frequency (idf)

---

Inverse document frequency (idf)

$$\text{idf}_t = \log_{10} \left( \frac{N}{\text{df}_t} \right)$$

N is the total number of documents  
in the collection

Word	df	idf
Romeo	1	1.57
salad	2	1.27
Falstaff	4	0.967
forest	12	0.489
battle	21	0.246
wit	34	0.037
fool	36	0.012
good	37	0
sweet	37	0

## Tf-idf weighted representation

---

$$w_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$$

Raw counts:

	As You Like It	Twelfth Night	Julius Caesar
battle	1	0	7
good	114	80	62
fool	36	58	1
wit	20	15	2

tf-idf:

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

## valuation example corpus:

- I like deep learning.
- I like NLP.
- I enjoy flying.

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

Co-occurrence Matrix

valuation

## valuation

Singular Value Decomposition of co-occurrence matrix  $X$

Factorizes  $X$  into  $U\Sigma V^T$ , where  $U$  and  $V$  are orthonormal

$$\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix} = \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{U} \underbrace{\begin{bmatrix} \bullet & & & \\ & \bullet & & \\ & & \bullet & \\ & & & \bullet \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T}$$

Retain only  $k$  singular values, in order to generalize.

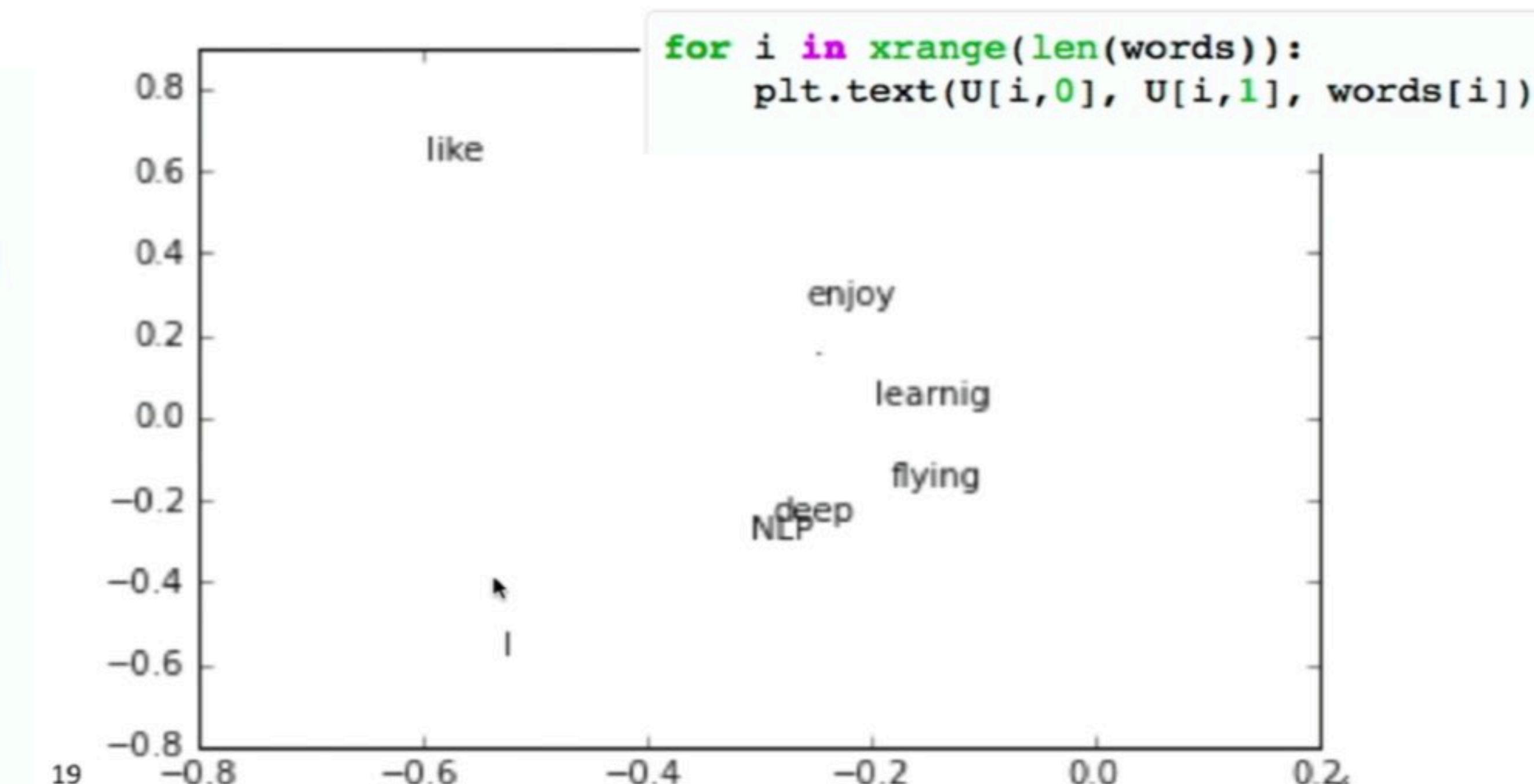
$\hat{X}$  is the best rank  $k$  approximation to  $X$ , in terms of least squares.

Classic linear algebra result. Expensive to compute for large matrices.

Corpus:  
I like deep learning. I like NLP. I enjoy flying.

```
import numpy as np
la = np.linalg
words = ["I", "like", "enjoy",
         "deep", "learnig", "NLP", "flying", ".."]
X = np.array([[0, 2, 1, 0, 0, 0, 0, 0],
              [2, 0, 0, 1, 0, 1, 0, 0],
              [1, 0, 0, 0, 0, 0, 1, 0],
              [0, 1, 0, 0, 1, 0, 0, 0],
              [0, 0, 0, 1, 0, 0, 0, 1],
              [0, 1, 0, 0, 0, 0, 0, 1],
              [0, 0, 1, 0, 0, 0, 0, 1],
              [0, 0, 0, 0, 1, 1, 1, 0]])
U, s, Vh = la.svd(X, full_matrices=False)
```

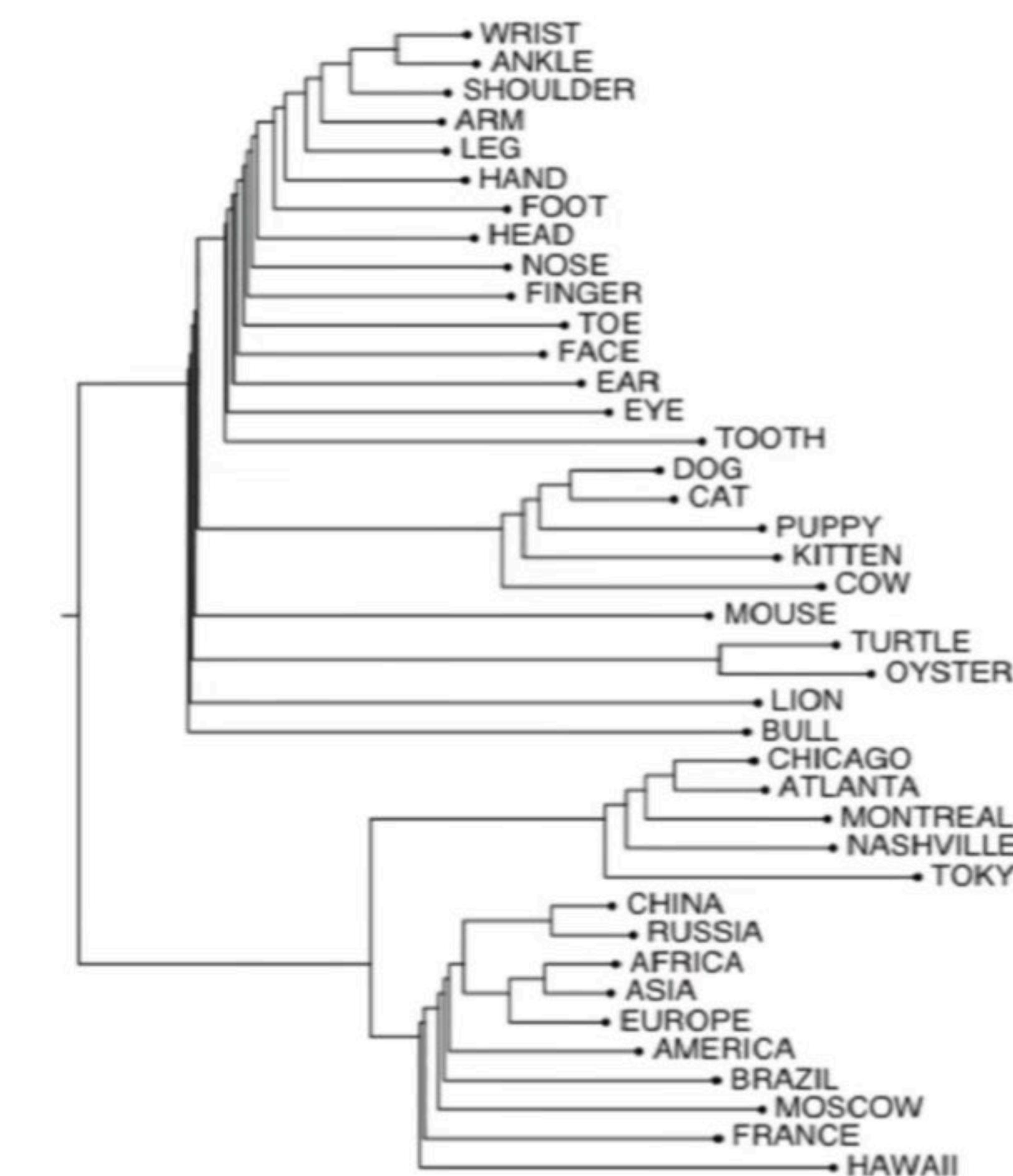
Corpus: I like deep learning. I like NLP. I enjoy flying.  
Printing first two columns of U corresponding to the 2 biggest singular values



19

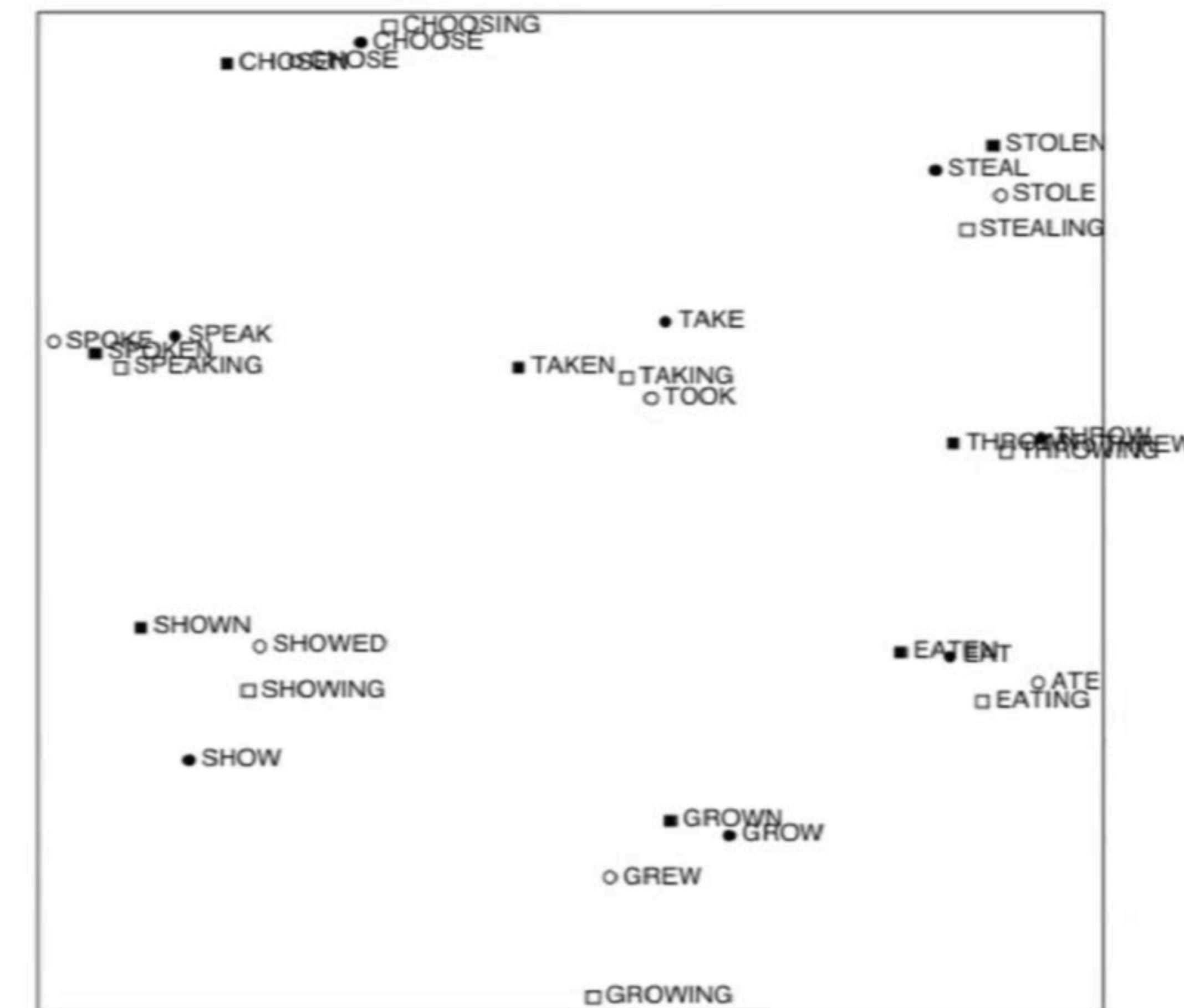
Computational cost scales quadratically for n x m matrix:  
 $O(mn^2)$  flops (when n < m)  
 → Bad for millions of words or documents

## valuatic



An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence  
Rohde et al. 2005

## valua



An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence  
Rohde et al. 2005

# **Distributed word representation**

---

INTERNATIONAL INSTITUTE OF  
INFORMATION TECHNOLOGY

---

H Y D E R A B A D

## Sparse vs Dense Vectors

---

tf-idf (or PMI) vectors are

- **long** (length  $|V| = 20,000$  to  $50,000$ )
- **sparse** (most elements are zero)

Alternative: learn vectors which are

- **short** (length  $50-1000$ )
- **dense** (most elements are non-zero)

## Sparse vs Dense Vectors

---

Why dense vectors?

- Short vectors may be easier to use as **features** in machine learning (fewer weights to tune)
- Dense vectors may **generalize** better than explicit counts
- Dense vectors may do better at capturing synonymy:
  - *car* and *automobile* are synonyms; but are distinct dimensions
  - a word with *car* as a neighbor and a word with *automobile* as a neighbor should be similar, but aren't
- **In practice, they work better**

# Distributional Hypothesis

---

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

---

Ludwig Wittgenstein

PI #43:  
"The meaning of a word is its use in the language"

## Represent words by their “usage”

---

One way to define "usage":

words are defined by their environments (the words around them)

Zellig Harris (1954):

If A and B have almost identical environments we say that they are synonyms.

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**

## **Represent words by their “usage”**

---

Idea 1: Defining meaning by linguistic distribution

Let's define the meaning of a word by its distribution in language use, meaning its neighboring words or grammatical environments.

## Represent words by their “usage”

---

3 affective dimensions for a word

- **valence:** pleasantness
- **arousal:** intensity of emotion
- **dominance:** the degree of control exerted

	Word	Score		Word	Score	
<b>Valence</b>	love	1.000		toxic	0.008	NRC VAD Lexicon (Mohammad 2018)
	happy	1.000		nightmare	0.005	
<b>Arousal</b>	elated	0.960		mellow	0.069	NRC VAD Lexicon (Mohammad 2018)
	frenzy	0.965		napping	0.046	
<b>Dominance</b>	powerful	0.991		weak	0.045	
	leadership	0.983		empty	0.081	

Idea 1: Defining meaning by linguistic distribution

Idea 2: Meaning as a point in multidimensional space

## Word Vectors

---

Intuition: why vectors?

Consider sentiment analysis:

- With **words**, a feature is a word identity
  - Feature : 'The previous word was "terrible"
  - requires **exact same word** to be in training and test
- With **embeddings**:
  - Feature is a word vector
  - '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!!!**

# Directly Learning Low dimensional dense vectors

---

Old idea. Relevant for this lecture & deep learning:

- Learning representations by back-propagating errors  
(Rumelhart et al., 1986)
- **A neural probabilistic language model** (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Weston, 2008)
- A recent, even simpler and faster model:  
word2vec (Mikolov et al. 2013) → intro now

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

# Getting short dense vectors there!

---

"Neural Language Model"-inspired models

- Word2vec (skipgram, CBOW), GloVe

Singular Value Decomposition (SVD)

- A special case of this is called LSA – Latent Semantic Analysis

Alternative to these "static embeddings":

- Contextual Embeddings (ELMo, BERT)
- Compute distinct embeddings for a word in its context
- Separate embeddings for each token of a word

# Word2vec

---

Predict between every word and its context words!

Two algorithms

1. **Skip-grams (SG)**

Predict context words given target (position independent)

2. **Continuous Bag of Words (CBOW)**

Predict target word from bag-of-words context

Two (moderately efficient) training methods

1. Hierarchical softmax
2. Negative sampling

# Word2vec

## there!

---

Instead of **counting** how often each word  $w$  occurs near "apricot"

- Train a classifier on a binary **prediction** task:
  - Is  $w$  likely to show up near "apricot"?

We don't actually care about this task

- But we'll take the learned classifier weights as the word embeddings

Big idea: **self-supervision**:

- A word  $c$  that occurs near apricot in the corpus acts as the gold "correct answer" for supervised learning
- No need for human labels
- Bengio et al. (2003); Collobert et al. (2011)

## Predict if candidate word $c$ is a “neighbour” there!

---

1. Treat the target word  $t$  and a neighboring context word  $c$  as **positive examples**.
2. Randomly sample other words in the lexicon to get negative examples
3. Use logistic regression to train a classifier to distinguish those two cases
4. Use the learned weights

## Skip-Gram Classifier

~~there!~~

(assuming a +/- 2 word window)

...lemon, a [tablespoon of apricot jam, a] pinch...

c1                    c2 [target]    c3    c4

Goal: train a classifier that is given a candidate (word, context) pair

(apricot, jam)

(apricot, aardvark)

...

And assigns each pair a probability:

$$P(+|w, c)$$

$$P(-|w, c) = 1 - P(+|w, c)$$

## Similarity is computed using dot product there!

---

Remember: two vectors are similar if they have a high dot product

- Cosine is just a normalized dot product

So:

- $\text{Similarity}(w,c) \propto w \cdot c$

We'll need to normalize to get a probability

- (cosine isn't a probability either)

## Turning dot products into probabilities there!

---

$$\text{Sim}(w, c) \approx w \cdot c$$

To turn this into a probability

We'll use the sigmoid from logistic regression:

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

$$\begin{aligned} P(-|w, c) &= 1 - P(+|w, c) \\ &= \sigma(-c \cdot w) = \frac{1}{1 + \exp(c \cdot w)} \end{aligned}$$

---

## Computing $P(+|w, c)$ there!

---

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

This is for one context word, but we have lots of context words.  
We'll assume independence and just multiply them:

$$P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w)$$

$$\log P(+|w, c_{1:L}) = \sum_{i=1}^L \log \sigma(c_i \cdot w)$$

## In nutshell

---

A probabilistic classifier, given

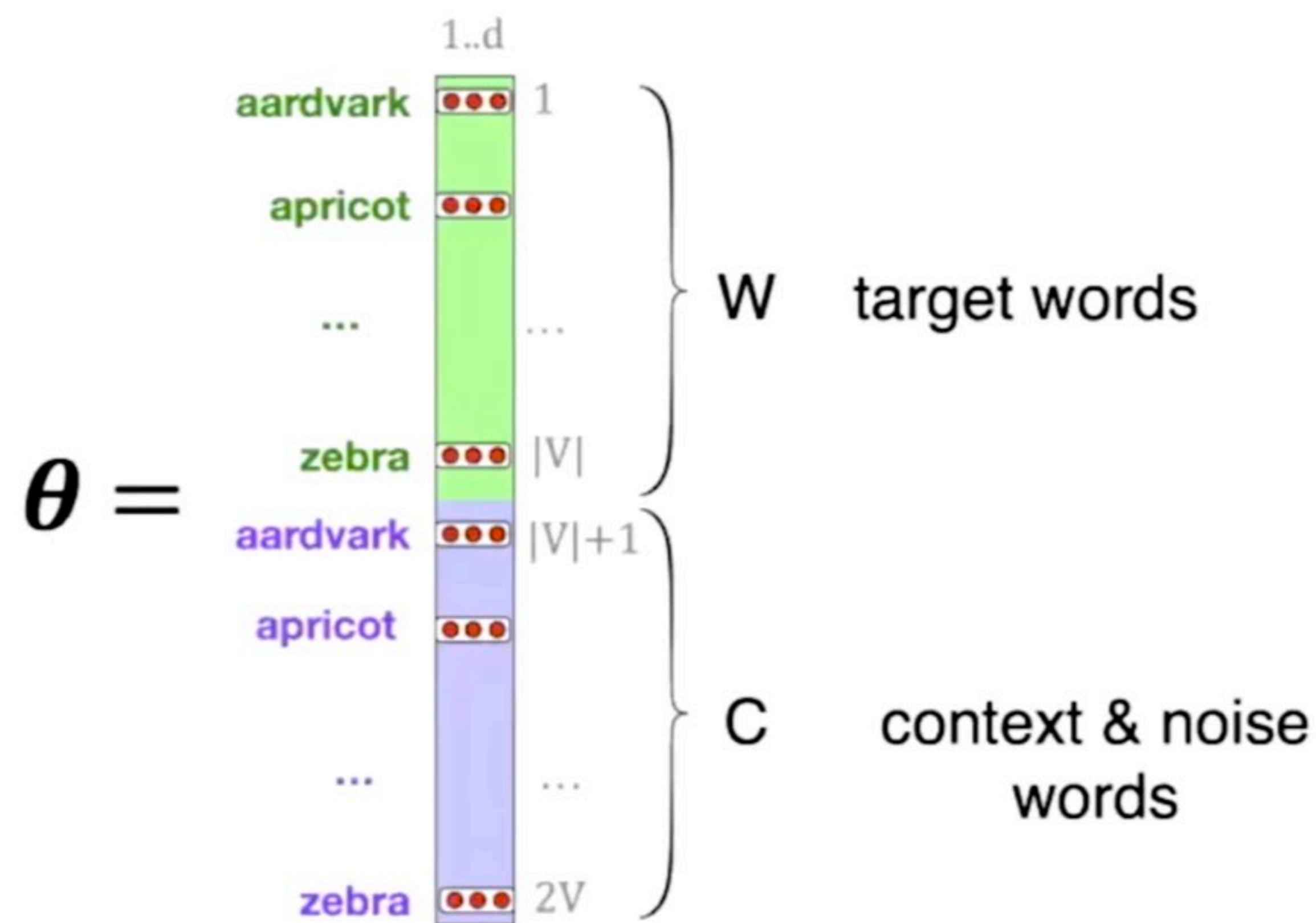
- a test target word  $w$
- its context window of  $L$  words  $c_{1:L}$

Estimates probability that  $w$  occurs in this window based on similarity of  $w$  (embeddings) to  $c_{1:L}$  (embeddings).

To compute this, we just need embeddings for all the words.

## Embeddings we need: for w, for c there!

---



## Skip Gram Training there!

---

...lemon, a [tablespoon of apricot jam, a] pinch...

c1

c2 [target] c3 c4



**positive examples +**

t c

---

apricot tablespoon

apricot of

apricot jam

apricot a

For each positive example we'll grab k negative examples, sampling by frequency

# Skip Gram Training there!

---

...lemon, a [tablespoon of apricot jam, a] pinch...

c1

c2 [target]



c3

c4

## positive examples +

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

## negative examples -

t	c	t	c
---	---	---	---

## Learning Vectors there!

---

Given the set of positive and negative training instances,  
and an initial set of embedding vectors

The goal of learning is to adjust those word vectors such  
that we:

- **Maximize** the similarity of the target word, context word pairs ( $w, c_{pos}$ ) drawn from the positive data
- **Minimize** the similarity of the ( $w, c_{neg}$ ) pairs drawn from the negative data.

## Loss function for one $w$ with $C_{pos}, C_{neg} \dots, C_{negk}$ there!

---

Maximize the similarity of the target with the actual context words, and minimize the similarity of the target with the  $k$  negative sampled non-neighbor words.

$$\begin{aligned} L_{CE} &= -\log \left[ P(+|w, c_{pos}) \prod_{i=1}^k P(-|w, c_{neg_i}) \right] \\ &= - \left[ \log P(+|w, c_{pos}) + \sum_{i=1}^k \log P(-|w, c_{neg_i}) \right] \\ &= - \left[ \log P(+|w, c_{pos}) + \sum_{i=1}^k \log (1 - P(+|w, c_{neg_i})) \right] \\ &= - \left[ \log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right] \end{aligned}$$

## Learning the Classifier there!

---

How to learn?

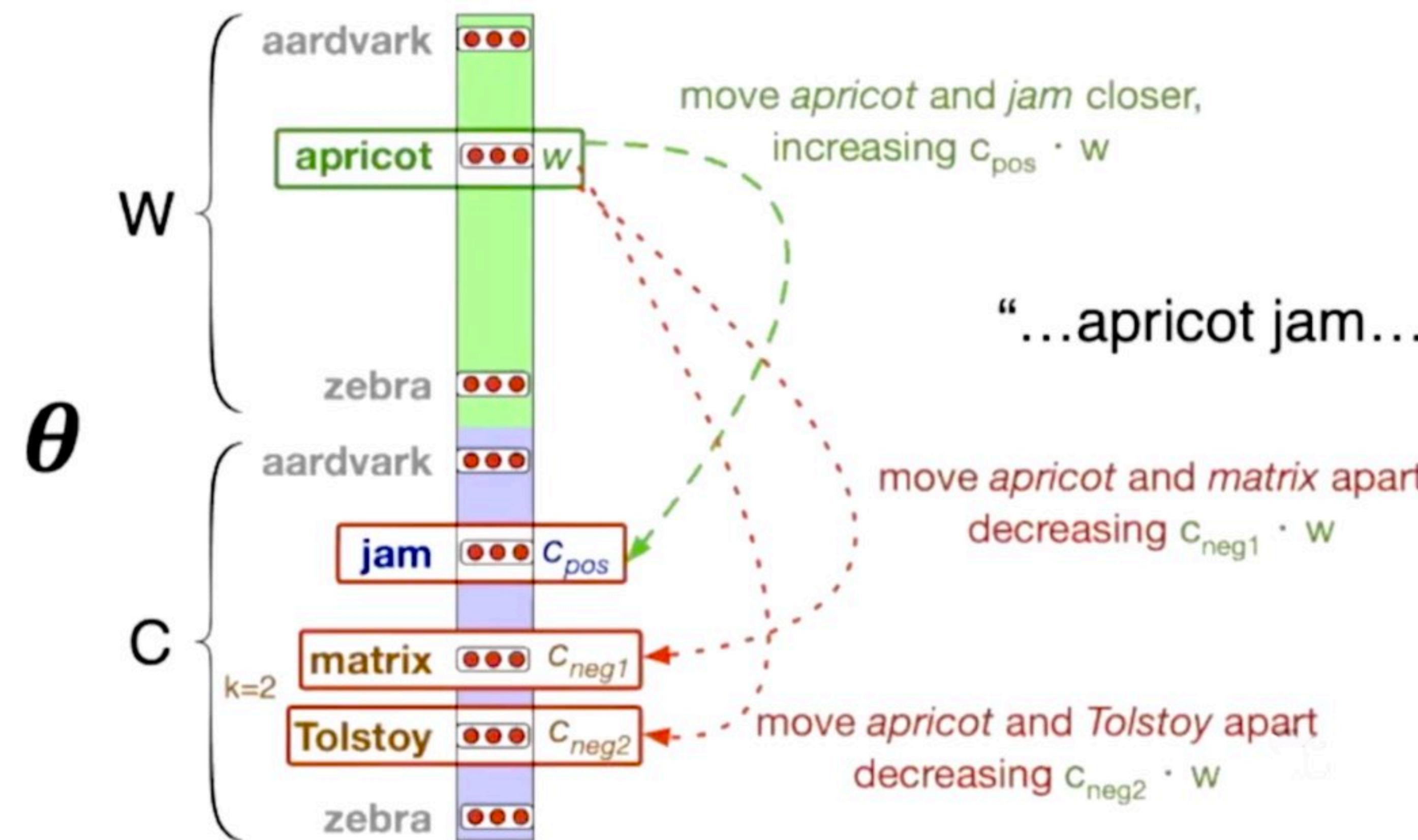
- Stochastic gradient descent!

We'll adjust the word weights to

- make the positive pairs more likely
- and the negative pairs less likely,
- over the entire training set.

# Intuition for one step of SGD

## valuation

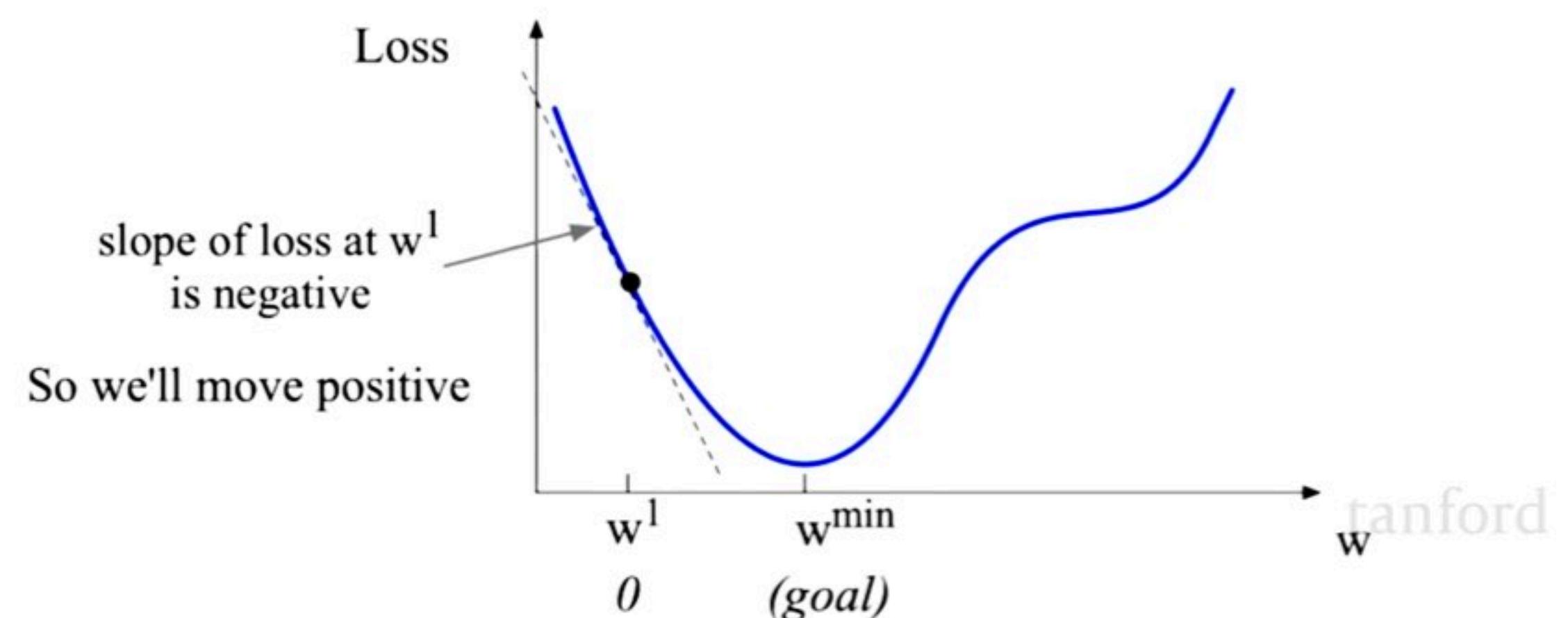
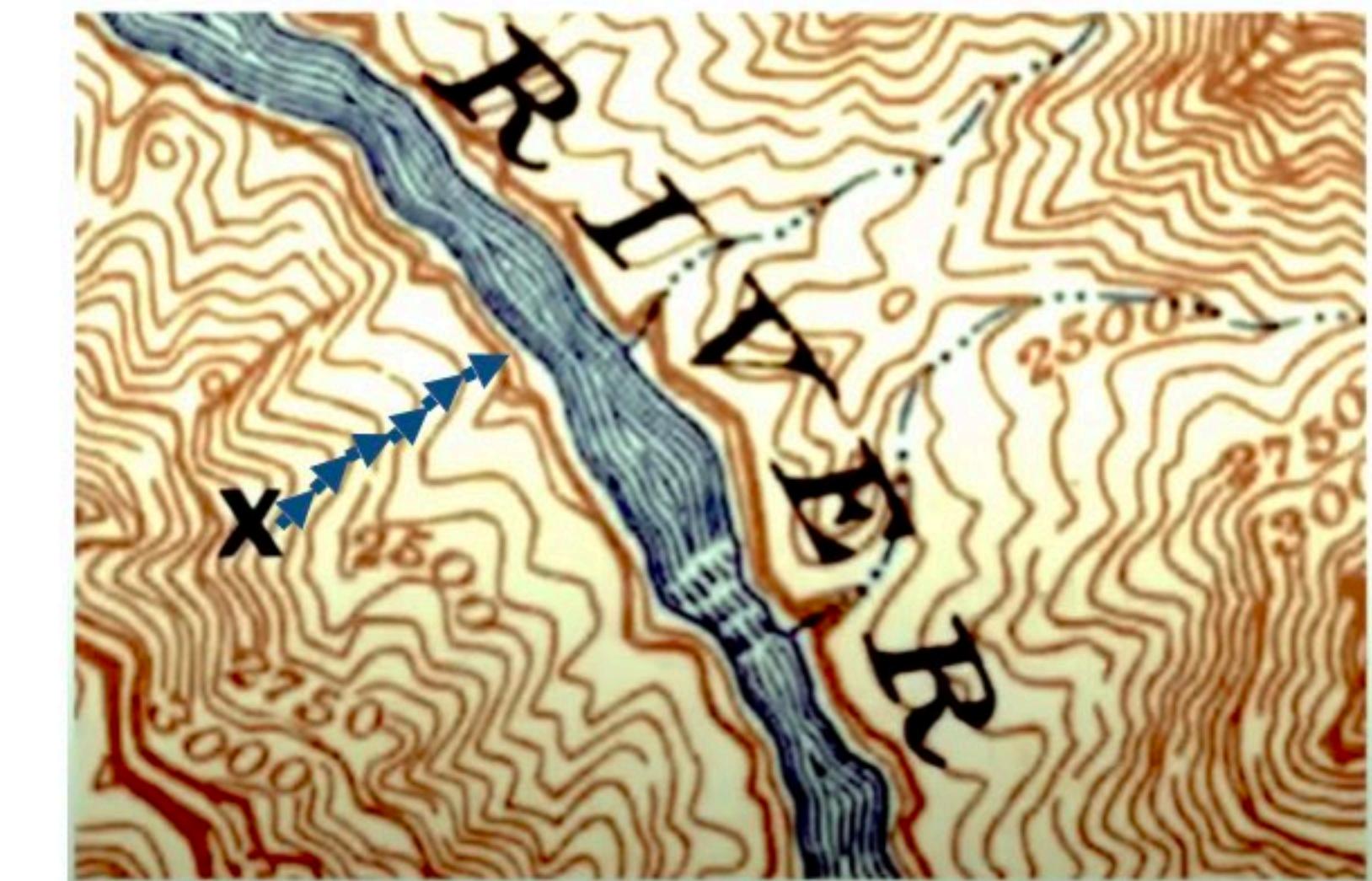


## valuation

- At each step
  - Direction: We move in the reverse direction from the gradient of the loss function
  - Magnitude: we move the value of this gradient  $\frac{d}{dw} L(f(x; w), y)$  weighted by a **learning rate**  $\eta$
  - Higher learning rate means move  $w$  faster

$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

slope of loss at  $w^1$   
is negative  
So we'll move positive



## Derivatives of Loss Function

---

$$L_{CE} = - \left[ \log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right]$$

$$\frac{\partial L_{CE}}{\partial c_{pos}} = [\sigma(c_{pos} \cdot w) - 1]w$$

$$\frac{\partial L_{CE}}{\partial c_{neg}} = [\sigma(c_{neg} \cdot w)]w$$

$$\frac{\partial L_{CE}}{\partial w} = [\sigma(c_{pos} \cdot w) - 1]c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i} \cdot w)]c_{neg_i}$$

# Update equation in SGD

---

## valuation

Start with randomly initialized C and W matrices, then incrementally do updates

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

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

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

## valuation

SGNS learns two sets of embeddings

Target embeddings matrix  $W$

Context embedding matrix  $C$

It's common to just add them together,  
representing word  $i$  as the vector  $w_i + c_i$

### valuation

Start with  $V$  random  $d$ -dimensional vectors as initial embeddings

Train a classifier based on embedding similarity

- Take a corpus and take pairs of words that co-occur as positive examples
- Take pairs of words that don't co-occur as negative examples
- Train the classifier to distinguish these by slowly adjusting all the embeddings to improve the classifier performance
- Throw away the classifier code and keep the embeddings.

# Observation and Properties

---

## valuation

**Small windows** ( $C = +/- 2$ ) : nearest words are syntactically similar words in same taxonomy

- *Hogwarts* nearest neighbors are other fictional schools
  - *Sunnydale, Evernight, Blandings*

**Large windows** ( $C = +/- 5$ ) : nearest words are related words in same semantic field

- *Hogwarts* nearest neighbors are Harry Potter world:
  - *Dumbledore, half-blood, Malfoy*

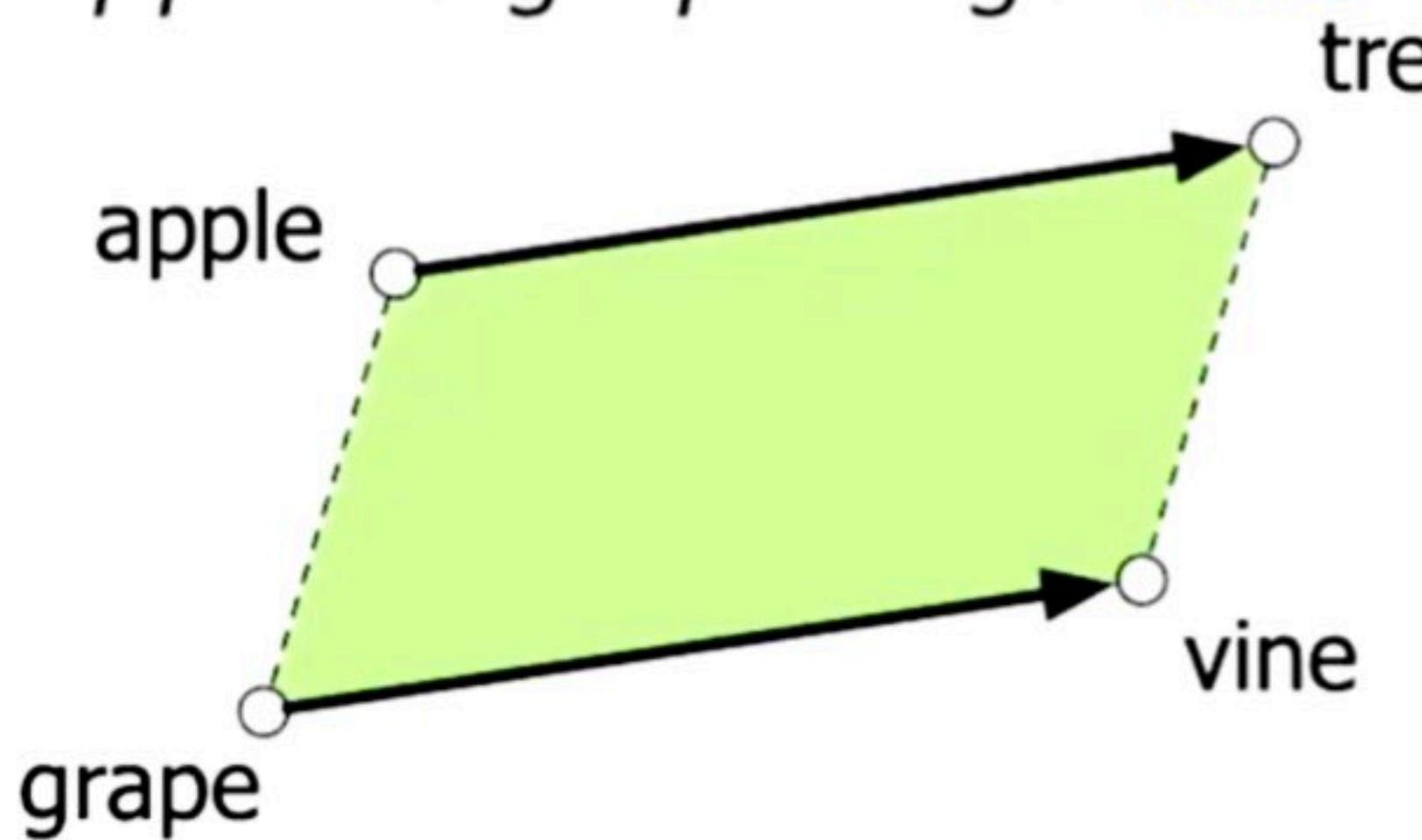
# Analogical Relations

## valuation

The classic parallelogram model of analogical reasoning  
(Rumelhart and Abrahamson 1973)

To solve: "*apple* is to *tree* as *grape* is to \_\_\_\_\_"

Add  $\overrightarrow{\text{tree}} - \overrightarrow{\text{apple}}$  to  $\overrightarrow{\text{grape}}$  to get  $\overrightarrow{\text{vine}}$



## valuation

The parallelogram method can solve analogies with both sparse and dense embeddings (Turney and Littman 2005, Mikolov et al. 2013b)

$$\overrightarrow{\text{king}} - \overrightarrow{\text{man}} + \overrightarrow{\text{woman}} \text{ is close to } \overrightarrow{\text{queen}}$$

$$\overrightarrow{\text{Paris}} - \overrightarrow{\text{France}} + \overrightarrow{\text{Italy}} \text{ is close to } \overrightarrow{\text{Rome}}$$

For a problem  $a:a^*:b:b^*$ , the parallelogram method is:

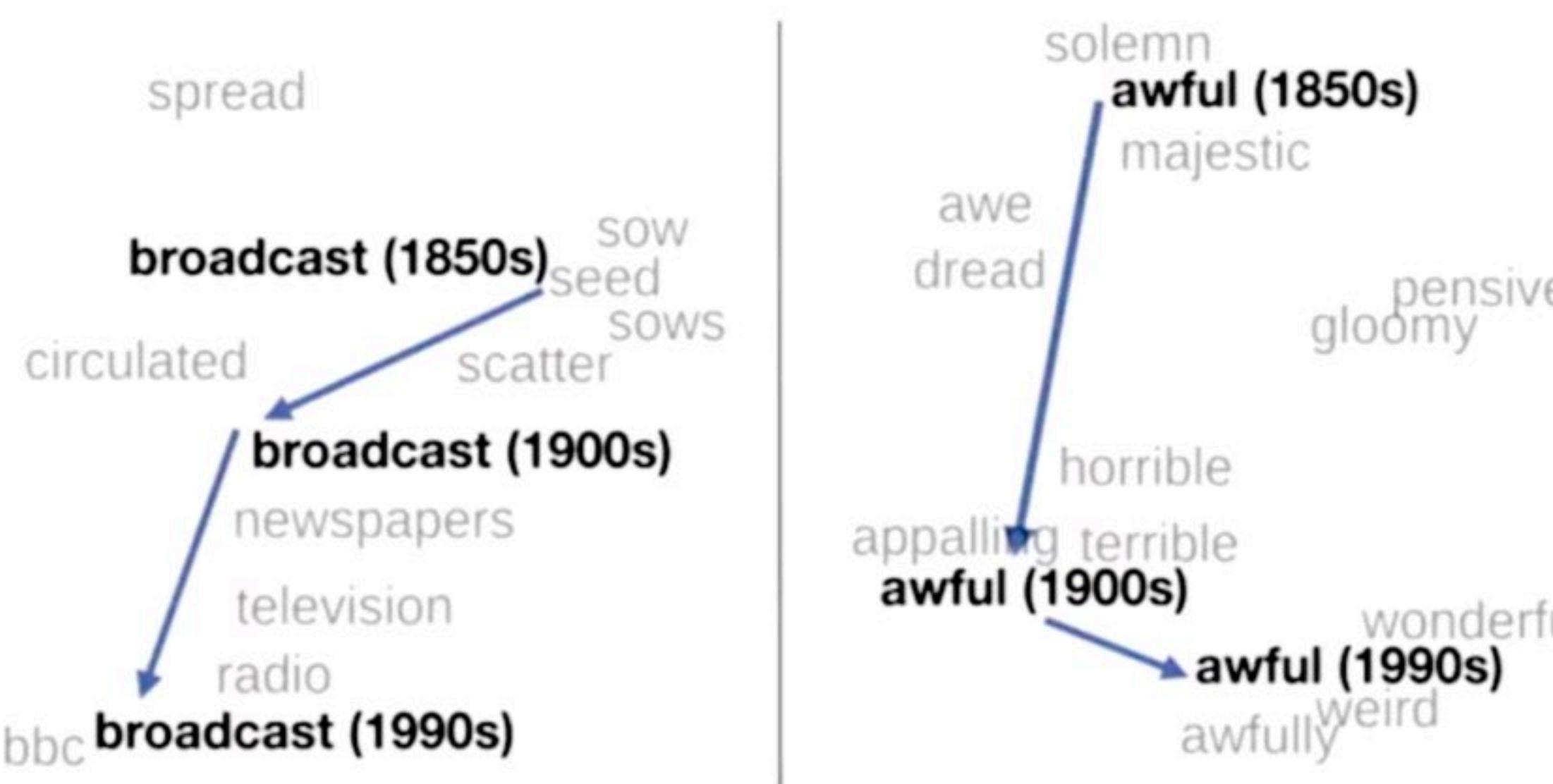
$$\hat{b}^* = \operatorname*{argmax}_x \text{distance}(x, a^* - a + b)$$

# Historical Semantics

## valuation

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

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



## valuation

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

## valuation

- LSA, HAL (Lund & Burgess),
- COALS, Hellinger-PCA (Rohde et al, Lebret & Collobert)
- Skip-gram/CBOW (Mikolov et al)
- NNLM, HLBL, RNN (Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton)

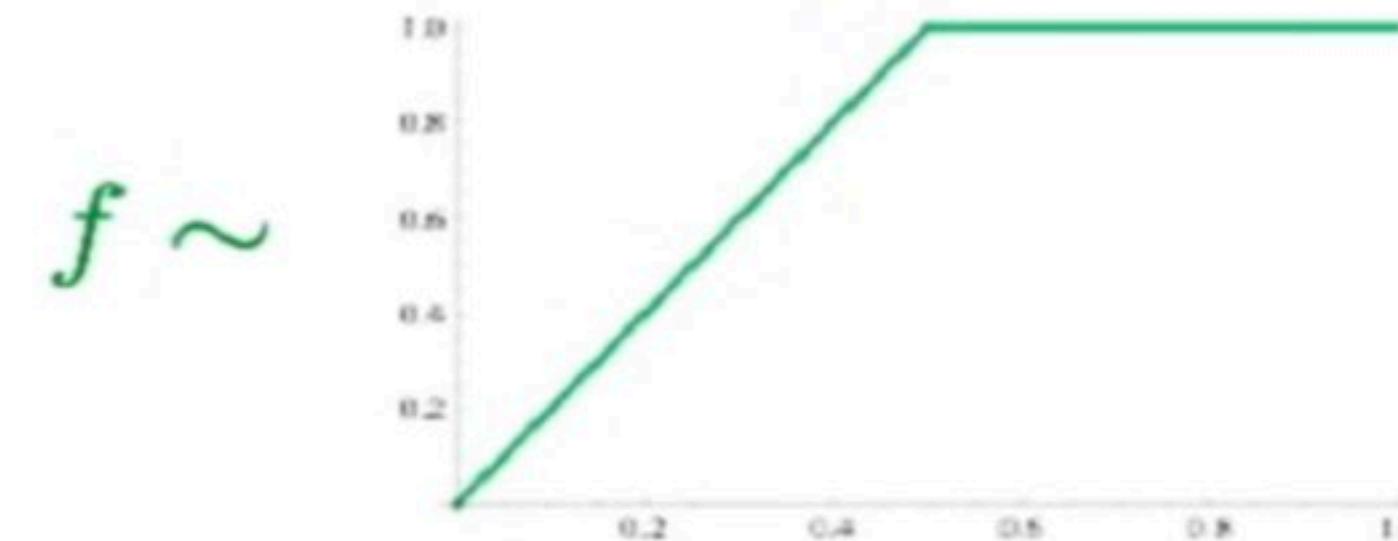
- 
- Fast training
  - Efficient usage of statistics
  - Primarily used to capture word similarity
  - Disproportionate importance given to large counts
  - Scale with corpus size
  - Inefficient usage of statistics
  - Generate improved performance on other tasks
  - Can capture complex patterns beyond word similarity

## GloVe: Combine the best of both worlds

---

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2$$

- Fast training
- Scalable to huge corpora
- Good performance even with small corpus, and small vectors
- By Pennington, Socher, Manning (2014)



## Problem with Word2vec and similar methods

---

These have two problems:

- Always the same representation for a **word type** regardless of the context in which a **word token** occurs
  - We might want very fine-grained word sense disambiguation
- We just have one representation for a word, but words have different **aspects**, including semantics, syntactic behavior, and register/connotations