

# Natural Language Processing with Deep Learning

## CS224N/Ling284



Christopher Manning

Lecture 2: Word Vectors, Word Senses, and  
Classifier Review



# Lecture Plan

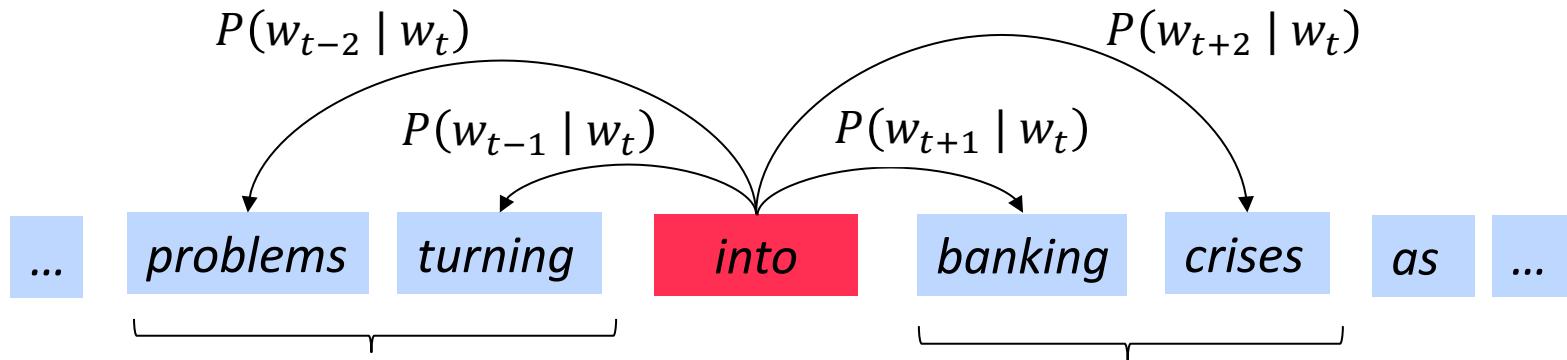
## Lecture 2: Word Vectors and Word Senses

1. Finish looking at word vectors and word2vec (10 mins)
2. Optimization basics (8 mins)
3. Can we capture this essence more effectively by counting? (12m)
4. The GloVe model of word vectors (10 min)
5. Evaluating word vectors (12 mins)
6. Word senses (6 mins)
7. Review of classification and how neural nets differ (10 mins)
8. Course advice (2 mins)

Goal: be able to read word embeddings papers by the end of class

# 1. Review: Main idea of word2vec

- Start with random word vectors
- Iterate through each word in the whole corpus  
遍历
- Try to predict surrounding words using word vectors



- $$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$
- Update vectors so you can predict better
- This algorithm learns word vectors that capture word similarity and meaningful directions in a wordspace

# Word2vec parameters and computations

$$\begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}$$

$U$

outside

$$\begin{bmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix}$$

$V$

center

$$\begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

$U \cdot v_4^T$

dot product

$$\begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

~~softmax~~( $U \cdot v_4^T$ )

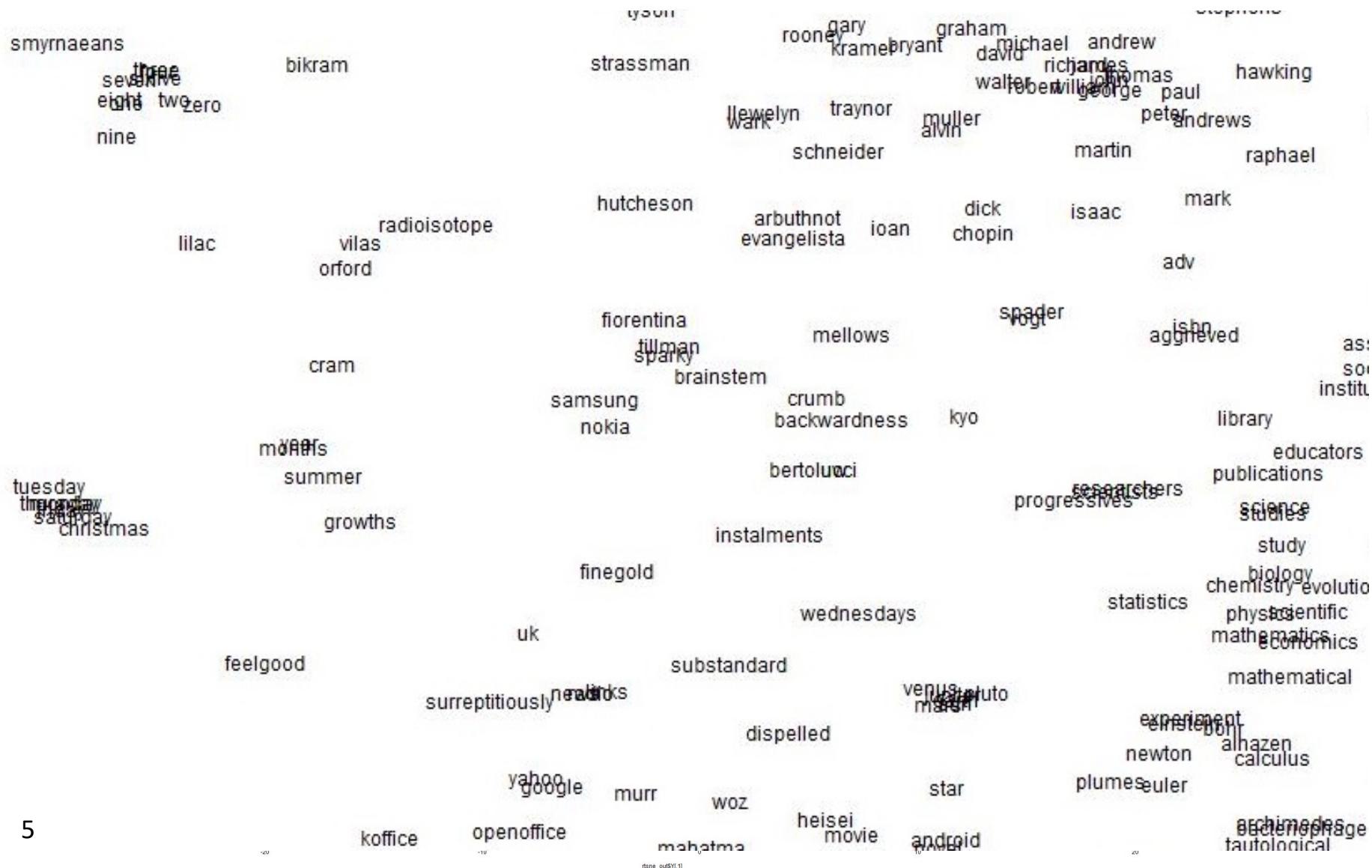
probabilities

!

Same predictions at each position

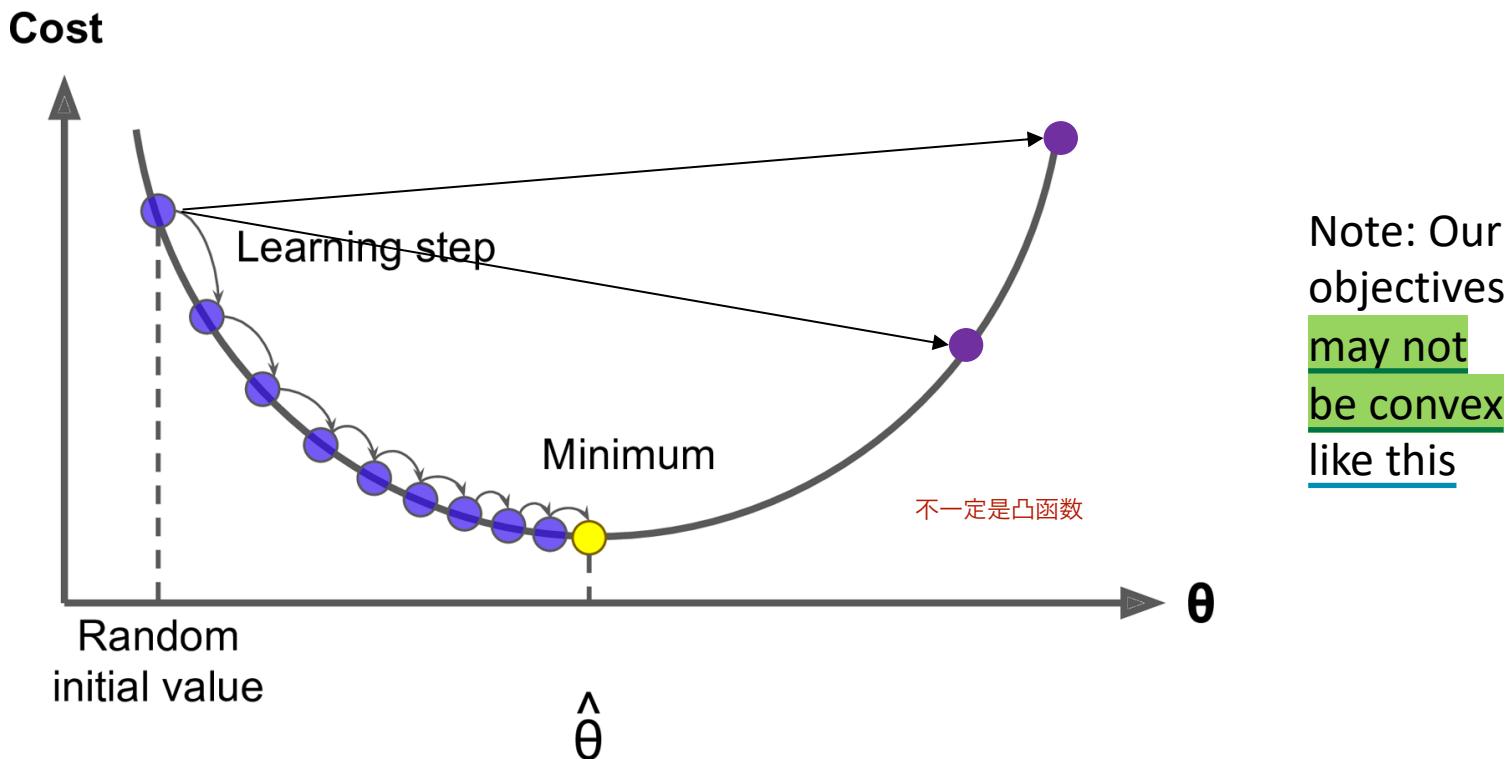
We want a model that gives a reasonably high probability estimate to all words that occur in the context (fairly often)

# Word2vec maximizes objective function by putting similar words nearby in space



## 2. Optimization: Gradient Descent

- We have a cost function  $J(\theta)$  we want to minimize
- **Gradient Descent** is an algorithm to minimize  $J(\theta)$   $\nabla J(\theta)$
- Idea: for current value of  $\theta$ , calculate gradient of  $J(\theta)$ , then take small step in the direction of negative gradient. Repeat.



# Gradient Descent

- Update equation (in matrix notation):

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

$\alpha$  = step size or learning rate

- Update equation (for a single parameter):

$$\theta_j^{new} = \theta_j^{old} - \alpha \frac{\partial}{\partial \theta_j^{old}} J(\theta)$$

$\propto \nabla J(\theta)$

- Algorithm:

```
while True:  
    theta_grad = evaluate_gradient(J, corpus, theta)  
    theta = theta - alpha * theta_grad
```

# Stochastic Gradient Descent

- Problem:  $J(\theta)$  is a function of **all** windows in the corpus (potentially billions!)
  - So  $\nabla_{\theta} J(\theta)$  is **very expensive to compute**
- You would wait a very long time before making a single update!
- **Very** bad idea for pretty much all neural nets!
- Solution: **Stochastic gradient descent (SGD)**
  - Repeatedly sample windows, and update after each one
- Algorithm:

```
while True:  
    window = sample_window(corpus)  
    theta_grad = evaluate_gradient(J,window,theta)  
    theta = theta - alpha * theta_grad
```

# Stochastic gradients with word vectors!

- Iteratively<sup>反复</sup> take gradients at each such window for SGD
- But in each window, we only have at most  $2m + 1$  words, so  $\nabla_{\theta} J_t(\theta)$  is very sparse!

$$\nabla_{\theta} J_t(\theta) = \begin{bmatrix} 0 \\ \vdots \\ \nabla v_{like} \\ \vdots \\ 0 \\ \nabla u_I \\ \vdots \\ \nabla u_{learning} \\ \vdots \end{bmatrix} \in \mathbb{R}^{2dV}$$

# Stochastic gradients with word vectors!

- We might only update the word vectors that actually appear!
- Solution: either you need sparse matrix update operations to only update certain rows of full embedding matrices  $U$  and  $V$ , or you need to keep around a hash for word vectors

$$|V| \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} & \textcolor{red}{\bullet} \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}^d$$

- If you have millions of word vectors and do distributed computing, it is important to not have to send gigantic 巨大的 updates around!

## 1b. Word2vec: More details

Why two vectors? → Easier optimization. Average both at the end

- But can do algorithm with just one vector per word

Two model variants:

1. Skip-grams (SG)

Predict context (“outside”) words (position independent) given center word

2. Continuous Bag of Words (CBOW)

Predict center word from (bag of) context words

We presented: **Skip-gram model**

Additional efficiency in training:

1. Negative sampling

So far: Focus on naïve softmax (simpler, but expensive, training method)

## The skip-gram model with negative sampling (HW2)

- The normalization factor is too computationally expensive.
- $$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$
- Hence, in standard word2vec and HW2 you implement the skip-gram model with negative sampling
- Main idea: train binary logistic regressions for a true pair (center word and word in its context window) versus several noise pairs (the center word paired with a random word)

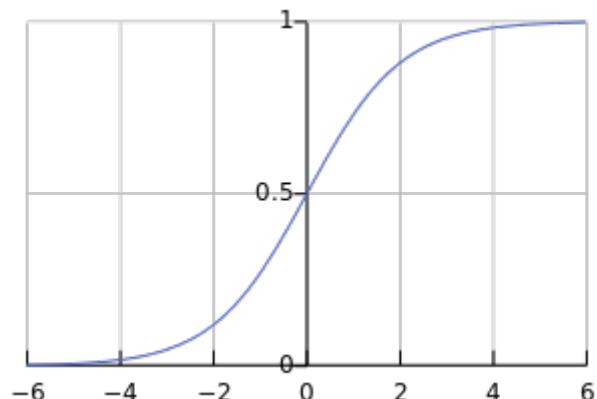
# The skip-gram model with negative sampling (HW2)

- From paper: “Distributed Representations of Words and Phrases and their Compositionality” (Mikolov et al. 2013)
- Overall objective function (they maximize):  $J(\theta) = \frac{1}{T} \sum_{t=1}^T J_t(\theta)$

$$J_t(\theta) = \underbrace{\log \sigma(u_o^T v_c)}_{\text{Positive}} + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

*Sampling negative samples*

- The sigmoid function:  $\sigma(x) = \frac{1}{1+e^{-x}}$   
(we'll become good friends soon)
- So we maximize the probability  
of two words co-occurring in first log  
→



## The skip-gram model with negative sampling (HW2)

- Notation more similar to class and HW2:

$$\begin{aligned} P(O|w, c, \theta) \\ = 1 - P(C|w, c, \theta) \end{aligned}$$

$$J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

- We take  $k$  negative samples (using word probabilities)
- Maximize probability that real outside word appears, minimize prob. that random words appear around center word
- $P(w) = U(w)^{3/4}/Z$ ,  
the unigram distribution  $U(w)$  raised to the  $3/4$  power  
(We provide this function in the starter code).
- The power makes less frequent words be sampled more often

### 3. Why not capture co-occurrence counts directly?

With a co-occurrence matrix  $X$

- 2 options: windows vs. full document
- Window: Similar to word2vec, use window around each word → captures both syntactic (POS) and semantic information
- Word-document co-occurrence matrix will give general topics (all sports terms will have similar entries) leading to “Latent Semantic Analysis” 潜语义分析

# Example: Window based co-occurrence matrix

- Window length 1 (more common: 5–10)
- Symmetric (对称的) (irrelevant whether left or right context) (不相关)
- Example corpus:
  - I like deep learning.
  - I like NLP.
  - I enjoy flying.

# Window based co-occurrence matrix

- 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

# Problems with simple co-occurrence vectors

Increase in size with vocabulary

Very high dimensional: requires a lot of storage

Subsequent classification models have sparsity issues

→ Models are less robust

# Solution: Low dimensional vectors

- Idea: store “most” of the important information in a fixed, small number of dimensions: a dense vector
- Usually 25–1000 dimensions, similar to word2vec
- How to reduce the dimensionality?

# Method: Dimensionality Reduction on X (HW1)

Singular Value Decomposition of co-occurrence matrix  $X$

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

$$\underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{X^k} = \underbrace{\begin{bmatrix} \star & \star \\ \star & \star \\ \star & \star \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} \star & \star & \star & \star & \star \\ \star & \star & \star & \star & \star \\ \star & \star & \star & \star & \star \\ \star & \star & \star & \star & \star \\ \star & \star & \star & \star & \star \end{bmatrix}}_{V^T}$$

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.

# Simple SVD word vectors in Python

Corpus:

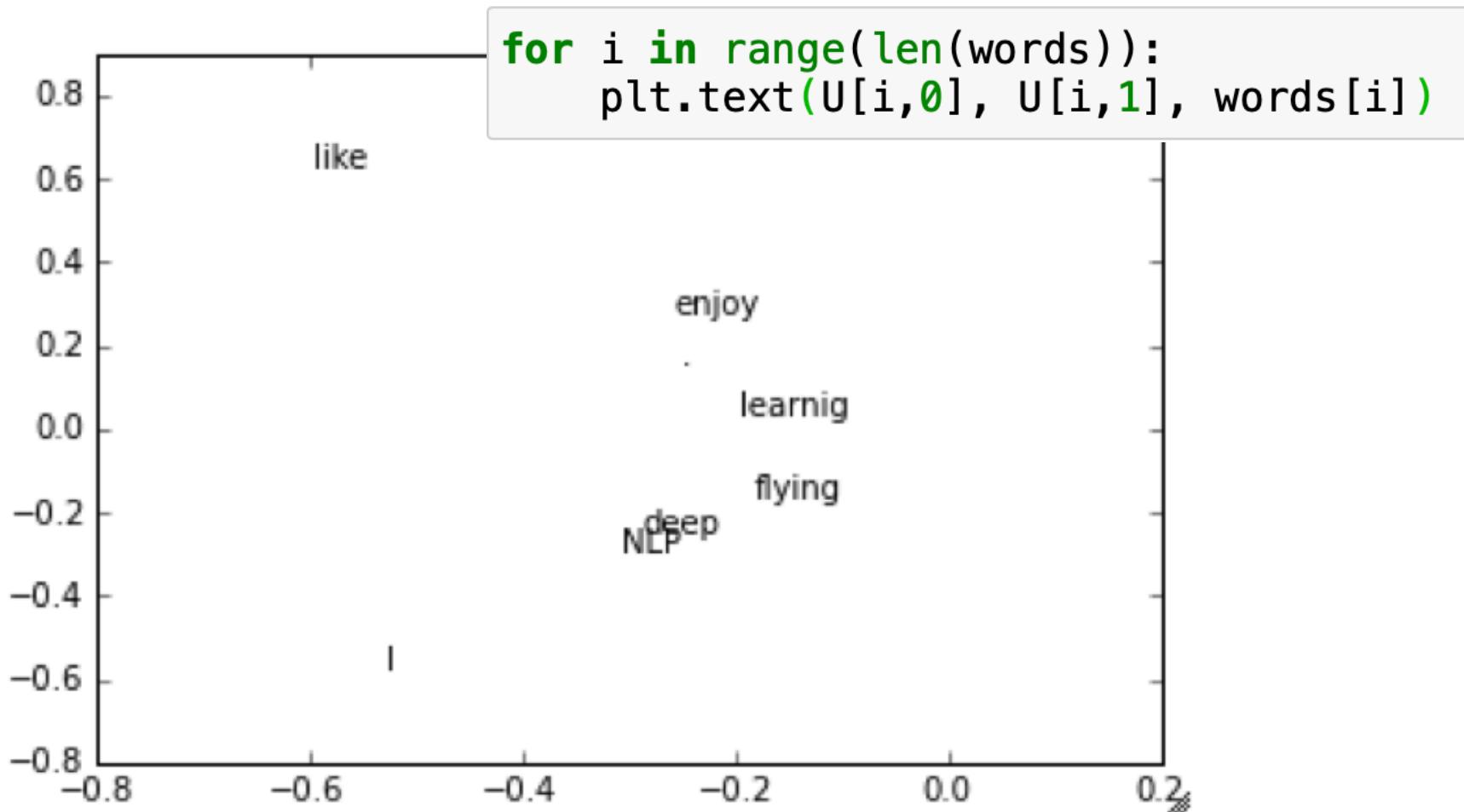
I like deep learning. I like NLP. I enjoy flying.

```
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)
```

# Simple SVD word vectors in Python

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

Printing first two columns of U corresponding to the 2 biggest singular values

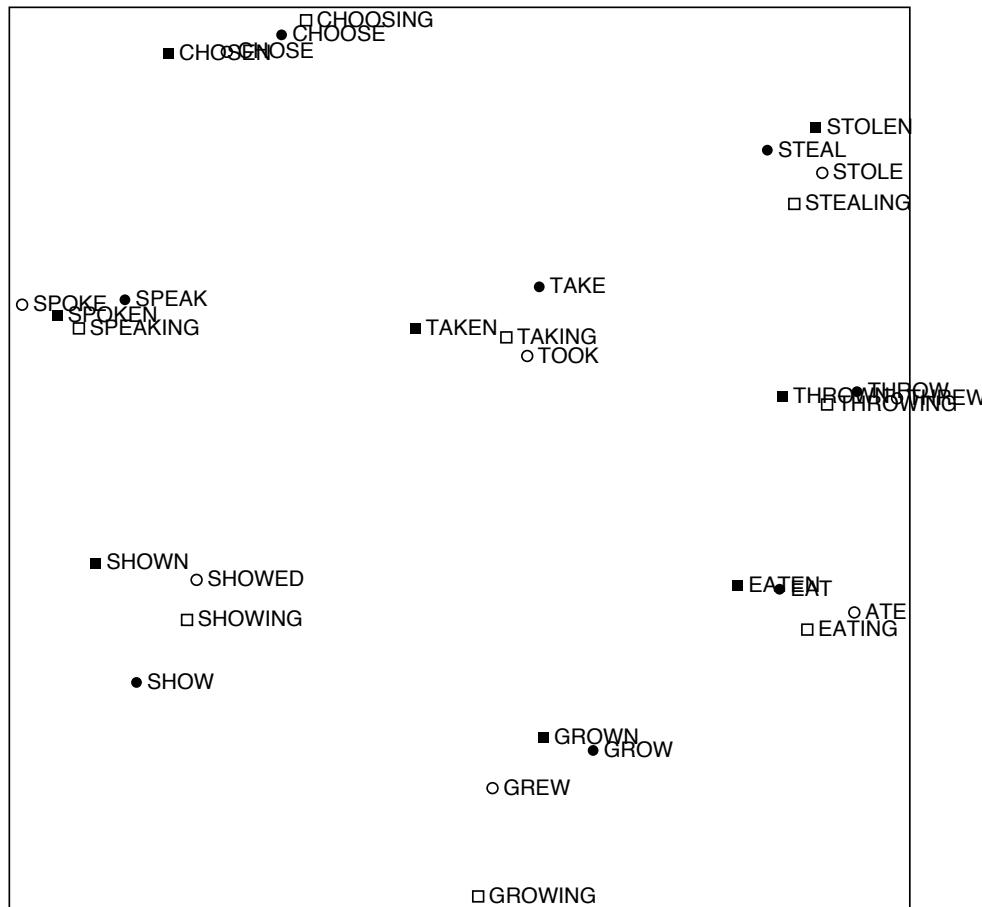


# Hacks to X (several used in Rohde et al. 2005)

Scaling the counts in the cells can help *a lot*

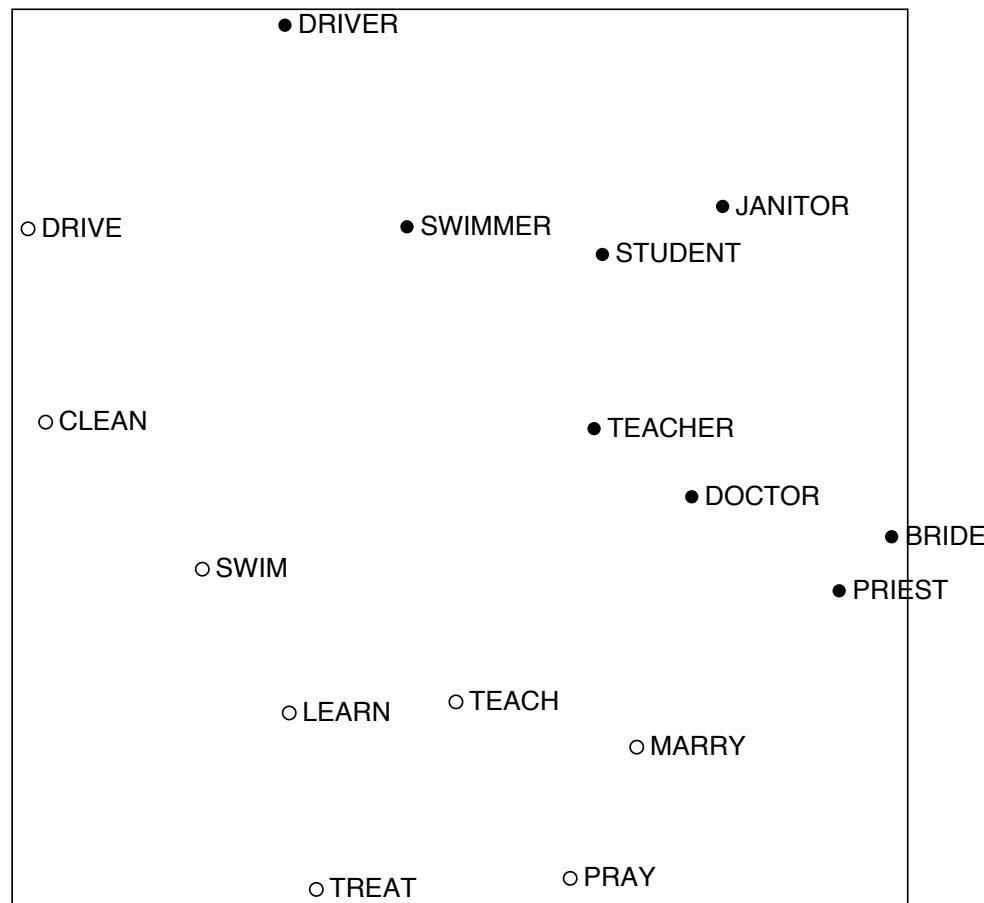
- Problem: function words (*the, he, has*) are too frequent → syntax has too much impact. Some fixes:
  - $\min(X, t)$ , with  $t \approx 100$
  - Ignore them all
- Ramped windows that count closer words more
- Use Pearson correlations instead of counts, then set negative values to 0
- Etc.

# Interesting syntactic patterns emerge in the vectors



COALS model from  
 An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence  
 Rohde et al. ms., 2005

# Interesting semantic patterns emerge in the vectors



COALS model from  
An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence  
Rohde et al. ms., 2005

## 4. Towards GloVe: Count based vs. direct prediction

- LSA, HAL (Lund & Burgess),
- COALS, Hellinger-PCA (Rohde et al, Lebret & Collobert)

- Fast training
- Efficient usage of statistics
- Primarily used to capture word similarity
- Disproportionate importance given to large counts  
不成比例的

- Skip-gram/CBOW (Mikolov et al)
- NNLM, HLBL, RNN (Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton)

- Scales with corpus size
- Inefficient usage of statistics
- Generate improved performance on other tasks
- Can capture complex patterns beyond word similarity



# Encoding meaning in vector differences

[Pennington, Socher, and Manning, EMNLP 2014]

Crucial insight: Ratios of co-occurrence probabilities can encode  
meaning components

共现概率表征词意

	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{random}$
$P(x \text{ice})$	large	small	large	small
$P(x \text{steam})$	small	large	large	small
$\frac{P(x \text{ice})}{P(x \text{steam})}$	large	small	~1	~1



# Encoding meaning in vector differences

[Pennington, Socher, and Manning, EMNLP 2014]

**Crucial insight:** Ratios of co-occurrence probabilities can encode meaning components

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



# Encoding meaning in vector differences

Q: How can we capture ratios of co-occurrence probabilities as linear meaning components in a word vector space?

A: Log-bilinear model:  
双线性

with vector differences

$$w_i \cdot w_j = \log P(i|j)$$

$$w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)}$$



# Combining the best of both worlds

## GloVe [Pennington et al., EMNLP 2014]

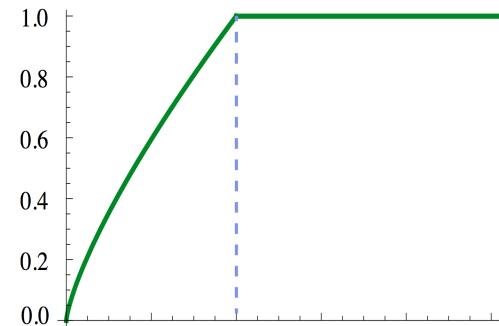


$$w_i \cdot w_j = \log P(i|j)$$

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

- Fast training
- Scalable to huge corpora
- Good performance even with small corpus and small vectors

$f \sim$



# GloVe results

Nearest words to  
[frog](#):

1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



[litoria](#)



[leptodactylidae](#)



[rana](#)



[eleutherodactylus](#)

# 5. How to evaluate word vectors?

- Related to general evaluation in NLP: Intrinsic vs. extrinsic
- Intrinsic: 内在的
  - Evaluation on a specific/intermediate subtask
  - Fast to compute
  - Helps to understand that system
  - Not clear if really helpful unless correlation to real task is established
- Extrinsic: 外在的
  - Evaluation on a real task
  - Can take a long time to compute accuracy
  - Unclear if the subsystem is the problem or its interaction or other subsystems
  - If replacing exactly one subsystem with another improves accuracy → Winning!

# Intrinsic word vector evaluation

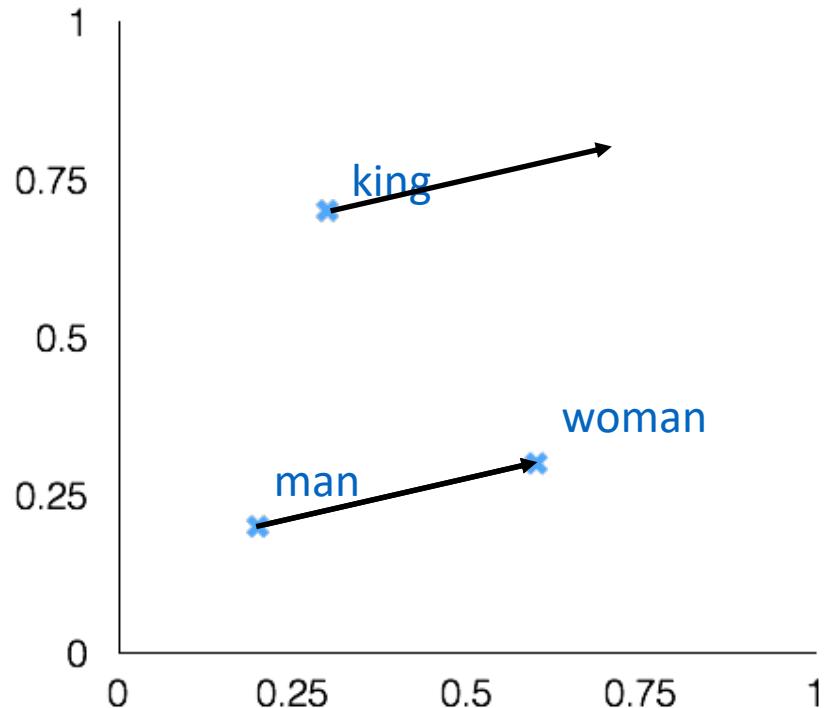
- Word Vector Analogies 类比

a:b :: c:?  
man:woman :: king:?

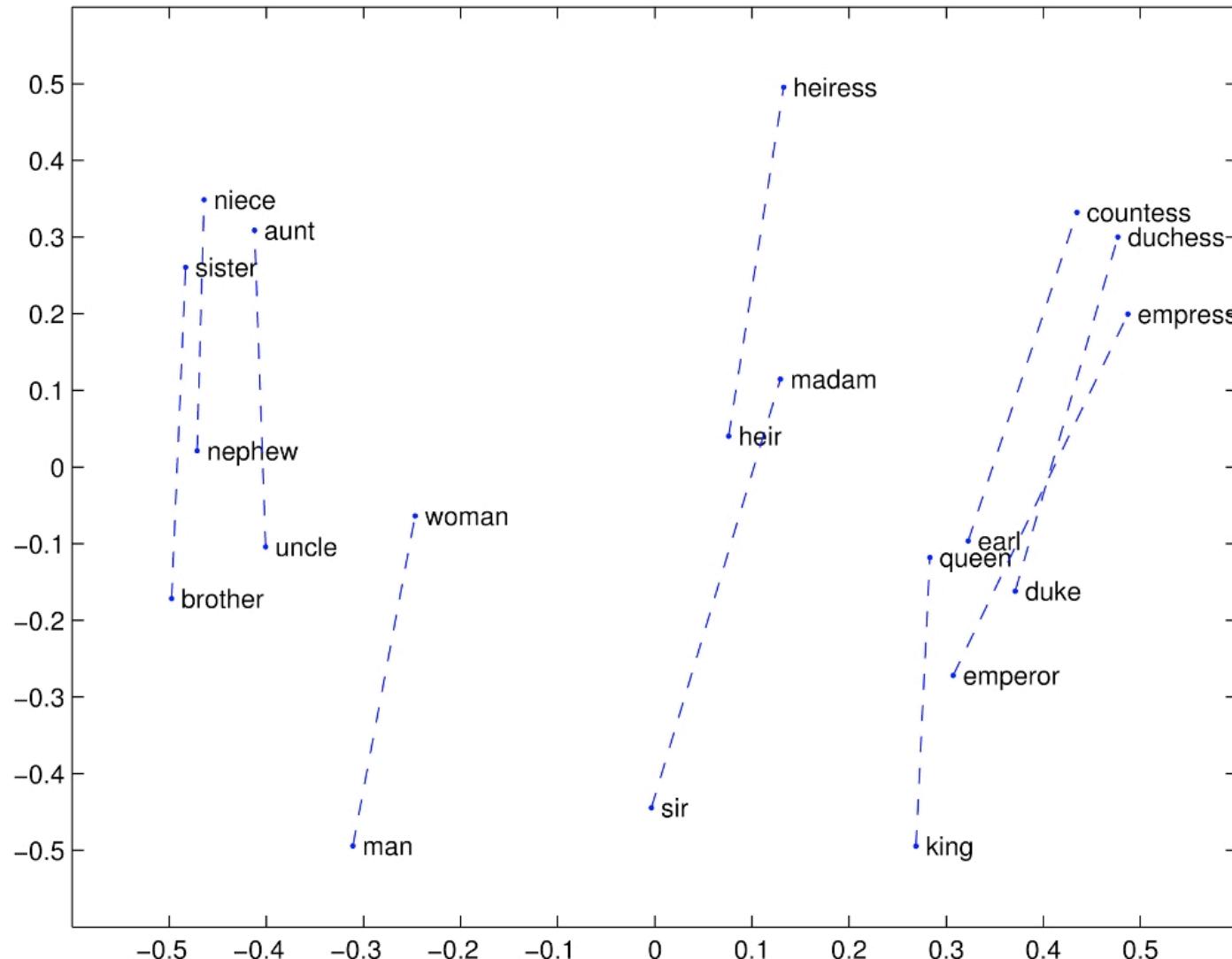


$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

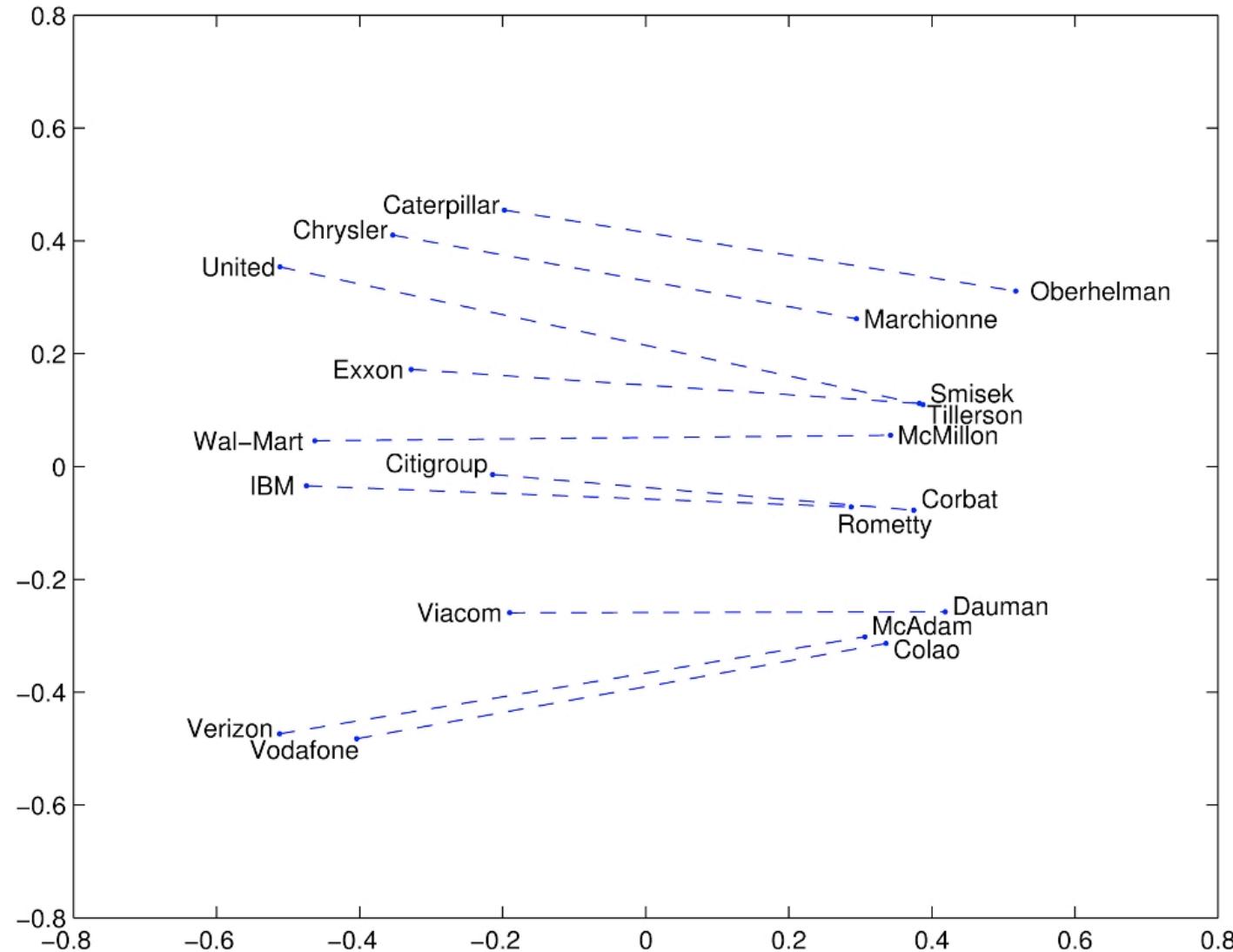
- Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy questions
- Discarding the input words from the search!
- Problem: What if the information is there but not linear?



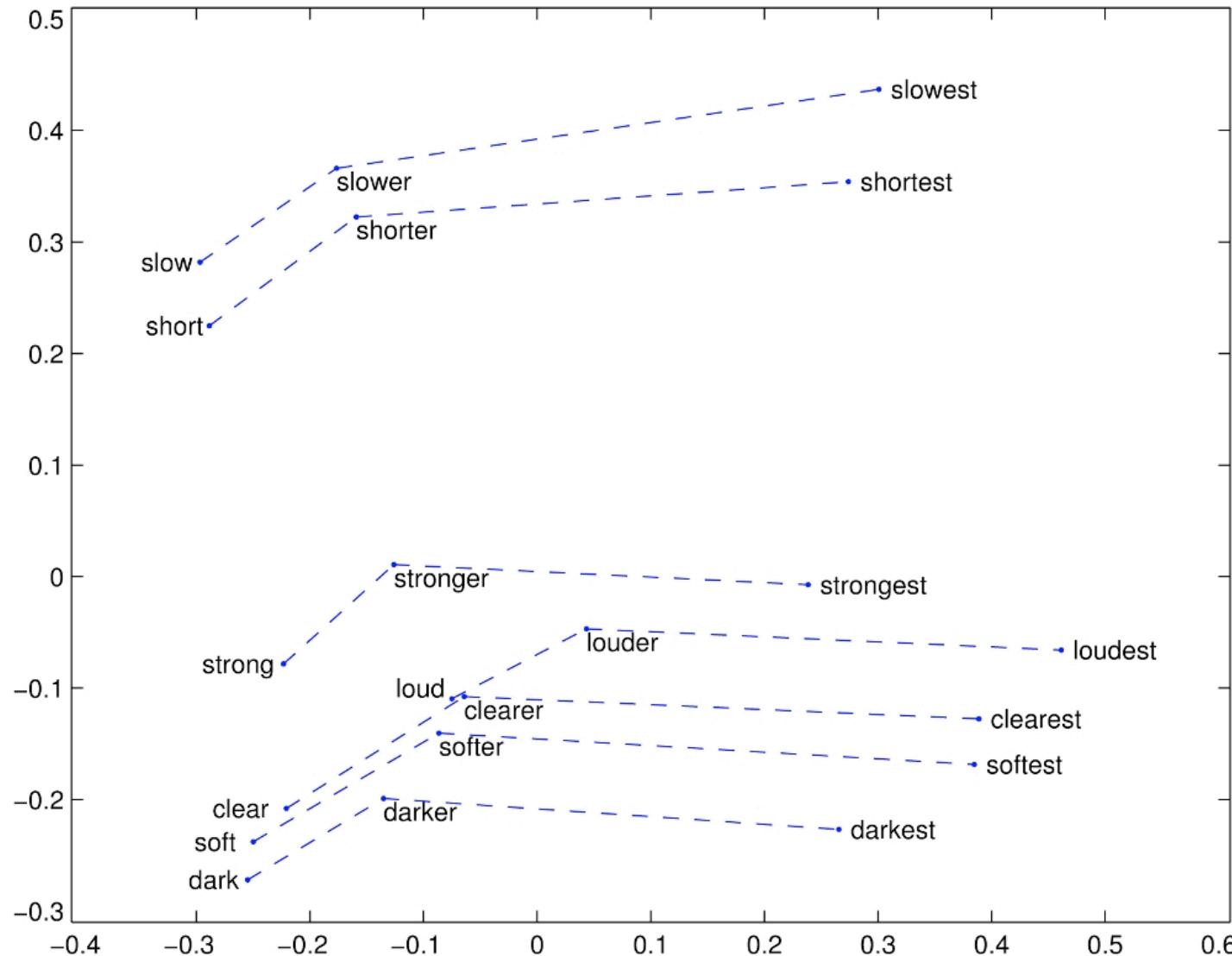
# Glove Visualizations



# Glove Visualizations: Company - CEO



# Glove Visualizations: Superlatives



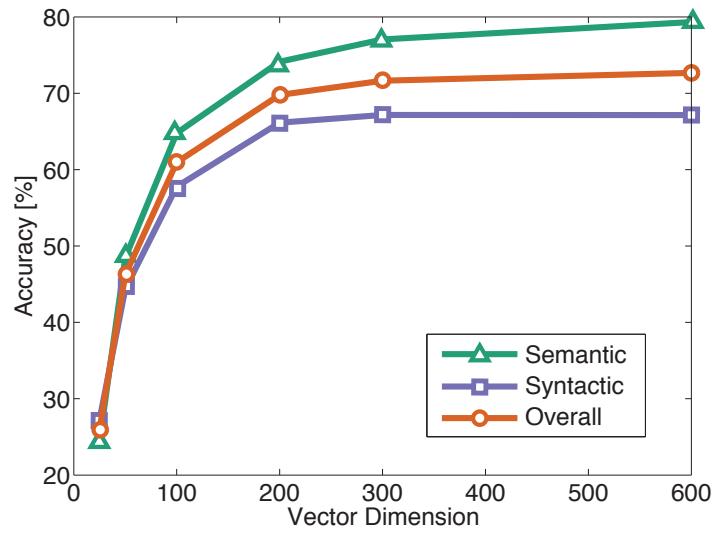
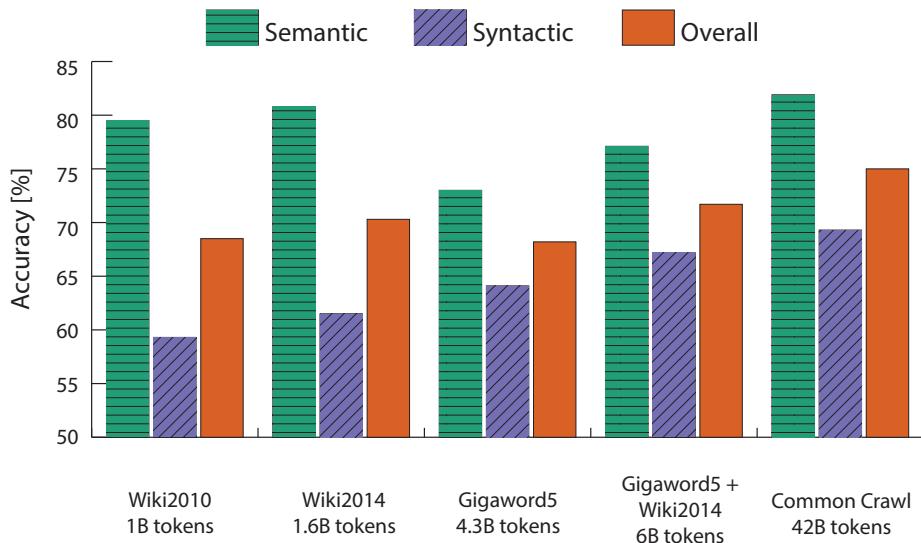
# Analogy evaluation and hyperparameters

## Glove word vectors evaluation

Model	Dim.	Size	Sem.	Syn.	Tot.
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>

# Analogy evaluation and hyperparameters

- More data helps
- Wikipedia is better than news text!
- Dimensionality
- Good dimension is  $\sim 300$



# Another intrinsic word vector evaluation

- Word vector distances and their correlation with human judgments
- Example dataset: WordSim353

<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

# Correlation evaluation

- Word vector distances and their correlation with human judgments

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	<u>65.6</u>	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	<b>75.9</b>	<b>83.6</b>	<b>82.9</b>	<b>59.6</b>	<b>47.8</b>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

- Some ideas from Glove paper have been shown to improve skip-gram (SG) model also (e.g. sum both vectors)

# Extrinsic word vector evaluation

- Extrinsic evaluation of word vectors: All subsequent tasks in this class
- One example where good word vectors should help directly: named entity recognition: finding a person, organization or location

Model	Dev	Test	<u>ACE</u>	<u>MUC7</u>
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	<b>88.7</b>	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	<b>93.2</b>	88.3	<b>82.9</b>	<b>82.2</b>

- Next: How to use word vectors in neural net models!

## 6. Word senses and word sense ambiguity

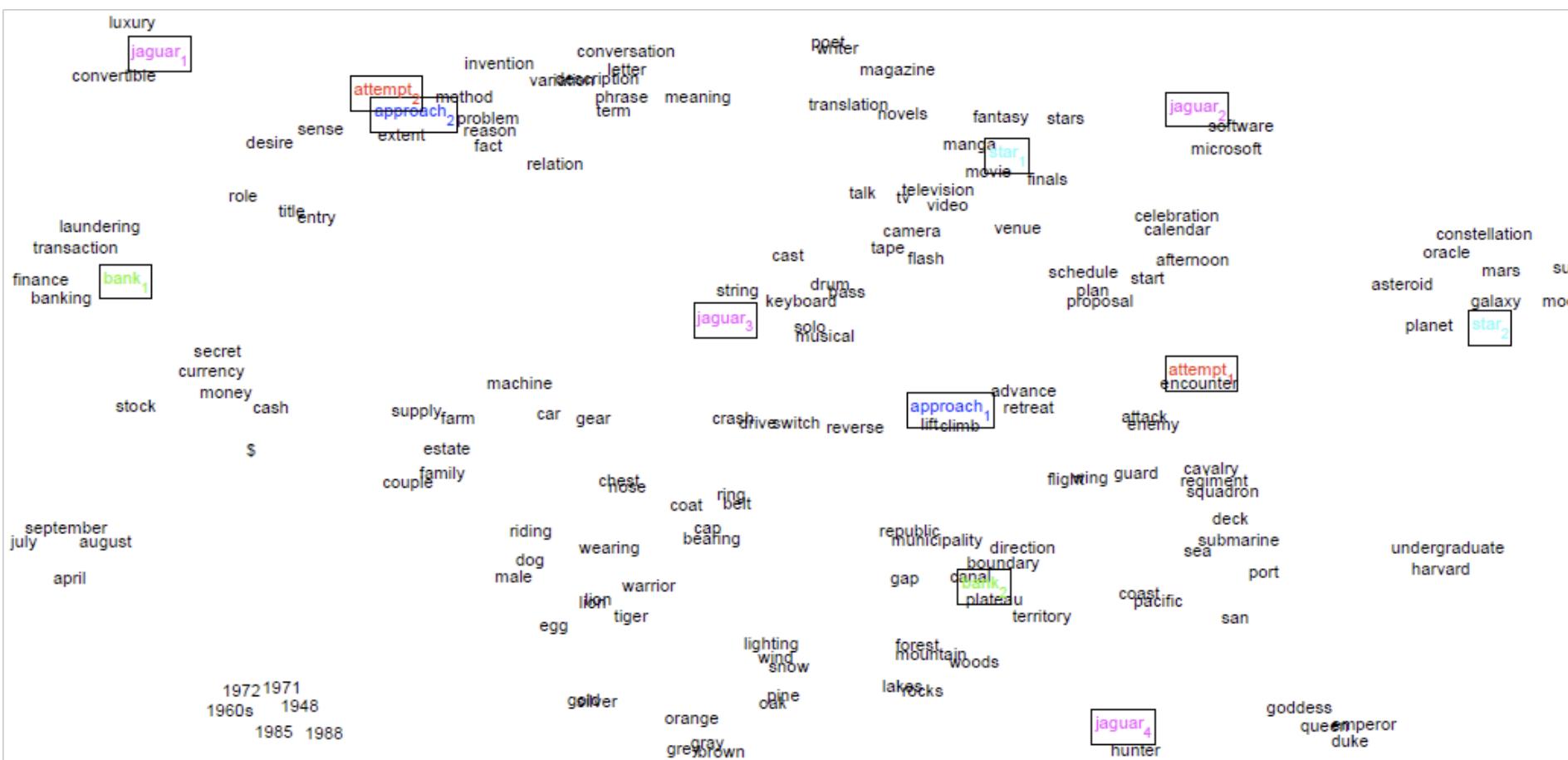
- Most words have lots of meanings!
  - Especially common words
  - Especially words that have existed for a long time
- Example: **pike**
- Does one vector capture all these meanings or do we have a mess?

# pike

- A sharp point or staff
- A type of elongated fish
- A railroad line or system
- A type of road
- The future (coming down the pike)
- A type of body position (as in diving)
- To kill or pierce with a pike
- To make one's way (pike along)
- In Australian English, pike means to pull out from doing something: *I reckon he could have climbed that cliff, but he piked!*

# Improving Word Representations Via Global Context And Multiple Word Prototypes (Huang et al. 2012)

- Idea: Cluster word windows around words, retrain with each word assigned to multiple different clusters  $\text{bank}_1$ ,  $\text{bank}_2$ , etc



# Linear Algebraic Structure of Word Senses, with Applications to Polysemy

多义性

(Arora, ..., Ma, ..., TACL 2018)

- Different senses of a word reside in a linear superposition (weighted sum) in standard word embeddings like word2vec
- $v_{\text{pike}} = \alpha_1 v_{\text{pike}_1} + \alpha_2 v_{\text{pike}_2} + \alpha_3 v_{\text{pike}_3}$
- Where  $\alpha_1 = \frac{f_1}{f_1+f_2+f_3}$ , etc., for frequency  $f$
- Surprising result:
  - Because of ideas from sparse coding you can actually separate out the senses (providing they are relatively common)

tie				
trousers	season	scoreline	wires	operatic
blouse	teams	goalless	cables	soprano
waistcoat	winning	equaliser	wiring	mezzo
skirt	league	clinching	electrical	contralto
sleeved	finished	scoreless	wire	baritone
pants	championship	replay	cable	coloratura

## 7. Classification review and notation

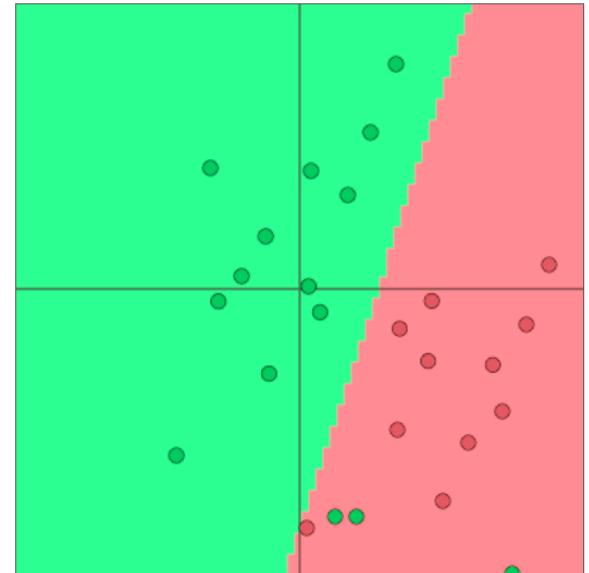
- Generally we have a **training dataset** consisting of **samples**

$$\{x_i, y_i\}_{i=1}^N$$

- $x_i$  are **inputs**, e.g. words (indices or vectors!), sentences, documents, etc.
  - Dimension  $d$
- $y_i$  are **labels** (one of  $C$  classes) we try to predict, for example:
  - classes: sentiment, named entities, buy/sell decision
  - other words
  - later: multi-word sequences

# Classification intuition

- Training data:  $\{x_i, y_i\}_{i=1}^N$
- Simple illustration case:
  - Fixed 2D word vectors to classify
  - Using softmax/logistic regression
  - Linear decision boundary
- **Traditional ML/Stats approach:** assume  $x_i$  are fixed, train (i.e., set) softmax/logistic regression weights  $W \in \mathbb{R}^{C \times d}$  to determine a decision boundary (hyperplane) as in the picture
- **Method:** For each fixed  $x$ , predict:



Visualizations with ConvNetJS by Karpathy!  
<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

$$p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

# Details of the softmax classifier

$$p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

We can tease apart the prediction function into three steps:

1. Take the  $y^{\text{th}}$  row of  $W$  and multiply that row with  $x$ :

$$W_y \cdot x = \sum_{i=1}^d W_{yi} x_i = f_y$$

Compute all  $f_c$  for  $c = 1, \dots, C$

2. Apply softmax function to get normalized probability:

$$p(y|x) = \frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} = \text{softmax}(f_y)$$

3. Choose the  $y$  with maximum probability

# Training with softmax and cross-entropy loss

- For each training example  $(x, y)$ , our objective is to **maximize the probability of the correct class  $y$**
- Or we can **minimize the negative log probability of that class:**

$$-\log p(y|x) = -\log \left( \frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} \right)$$

# Background: What is “cross entropy” loss/error?

- Concept of “cross entropy” is from information theory
- Let the true probability distribution be  $p$
- Let our computed model probability be  $q$
- The cross entropy is:

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c)$$

- Assuming a ground truth (or true or gold or target) probability distribution that is 1 at the right class and 0 everywhere else:  $p = [0, \dots, 0, 1, 0, \dots, 0]$  then:
- **Because of one-hot  $p$ , the only term left is the negative log probability of the true class**

# Classification over a full dataset

- Cross entropy loss function over full dataset  $\{x_i, y_i\}_{i=1}^N$

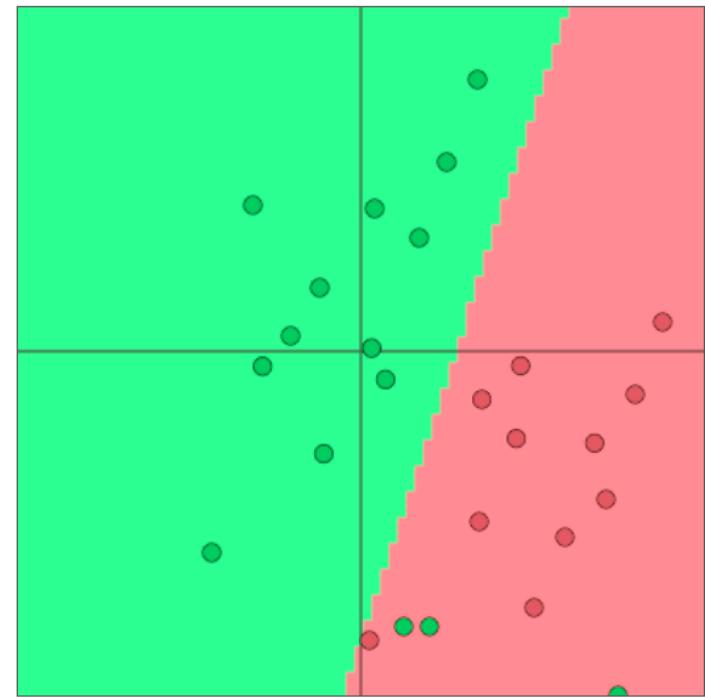
$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left( \frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right)$$

- Instead of

$$f_y = f_y(x) = W_y \cdot x = \sum_{j=1}^d W_{yj} x_j$$

We will write  $f$  in matrix notation:

$$f = Wx$$



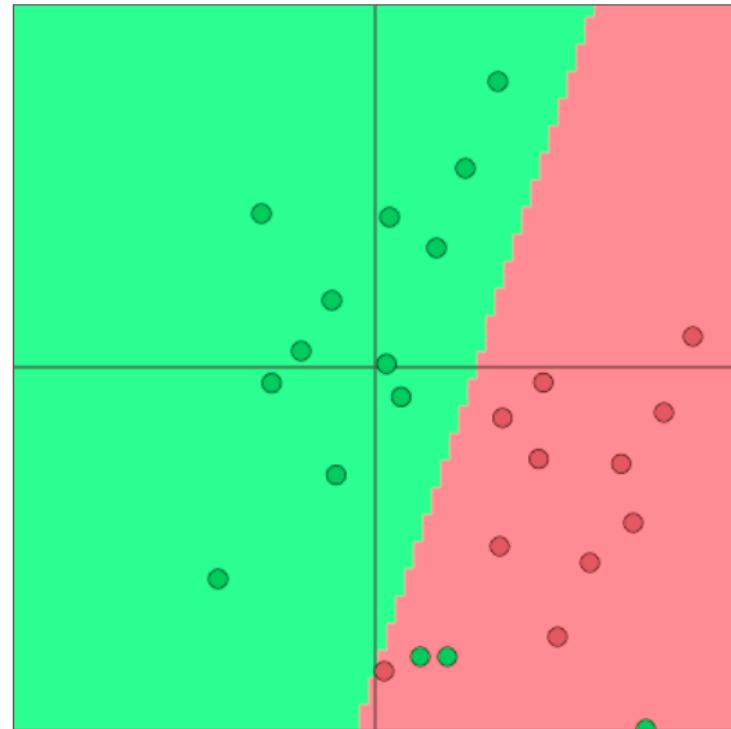
# Traditional ML optimization

- For general machine learning  $\theta$  usually only consists of columns of  $W$ :

$$\theta = \begin{bmatrix} W_{\cdot 1} \\ \vdots \\ W_{\cdot d} \end{bmatrix} = W(:) \in \mathbb{R}^{Cd}$$

- So we only update the decision boundary via

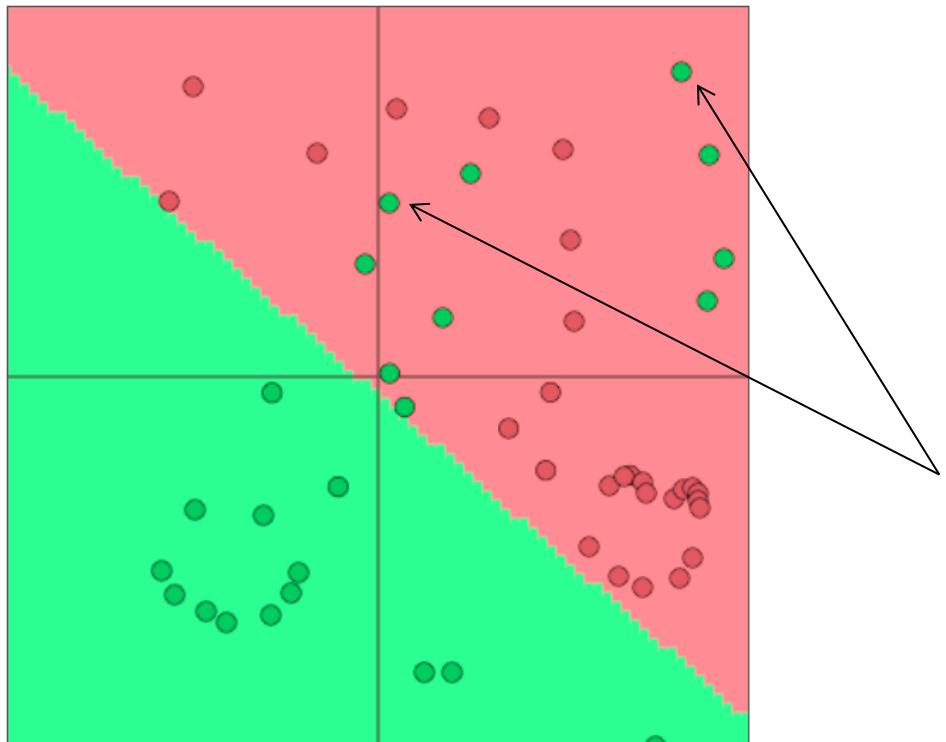
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \end{bmatrix} \in \mathbb{R}^{Cd}$$



Visualizations with ConvNetJS by Karpathy

# Neural Network Classifiers

- Softmax ( $\approx$  logistic regression) alone not very powerful
- Softmax gives only linear decision boundaries



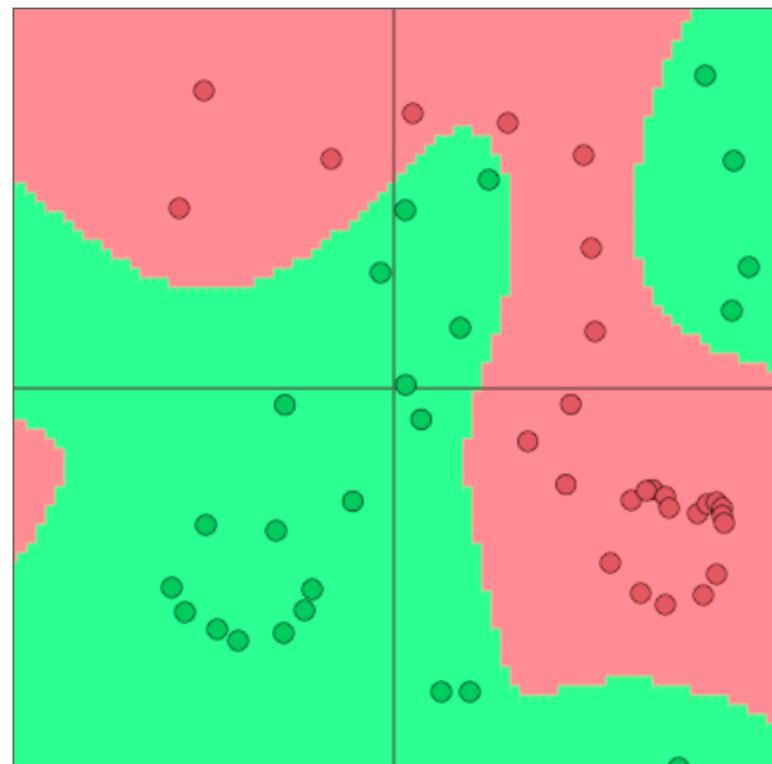
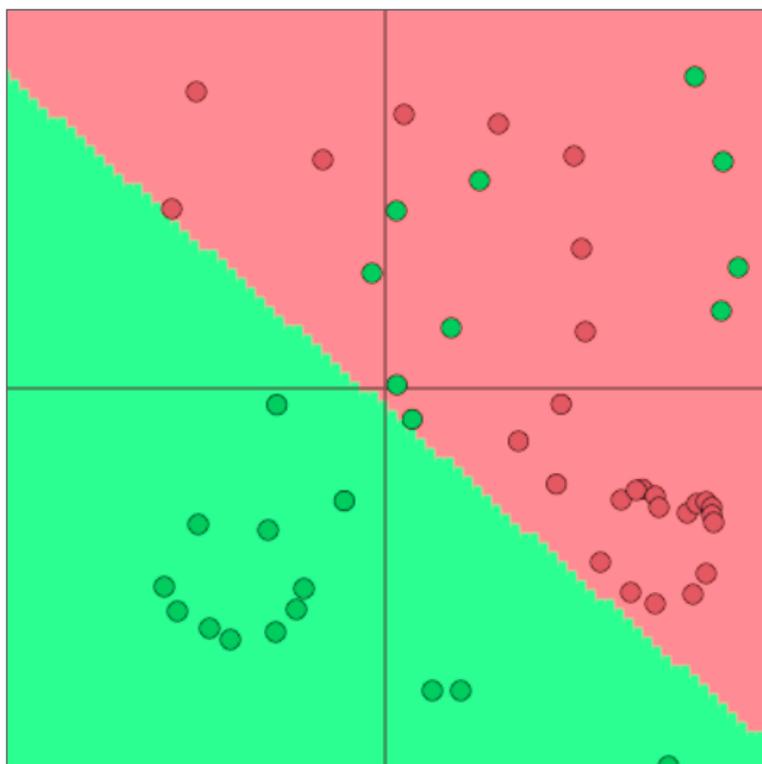
This can be quite limiting

→ Unhelpful when a problem is complex

Wouldn't it be cool to get these correct?

# Neural Nets for the Win!

- Neural networks can learn much more complex functions and nonlinear decision boundaries!
  - In the original space



# Classification difference with word vectors

- Commonly in NLP deep learning:
  - We learn **both**  $W$  and word vectors  $x$
  - We learn **both** conventional parameters and representations
  - The word vectors re-represent one-hot vectors—move them around in an intermediate layer vector space—for easy classification with a (linear) softmax classifier via layer  $x = Le$

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \\ \nabla_{x_{aardvark}} \\ \vdots \\ \nabla_{x_{zebra}} \end{bmatrix} \in \mathbb{R}^{Cd + Vd}$$

Very large number of parameters!

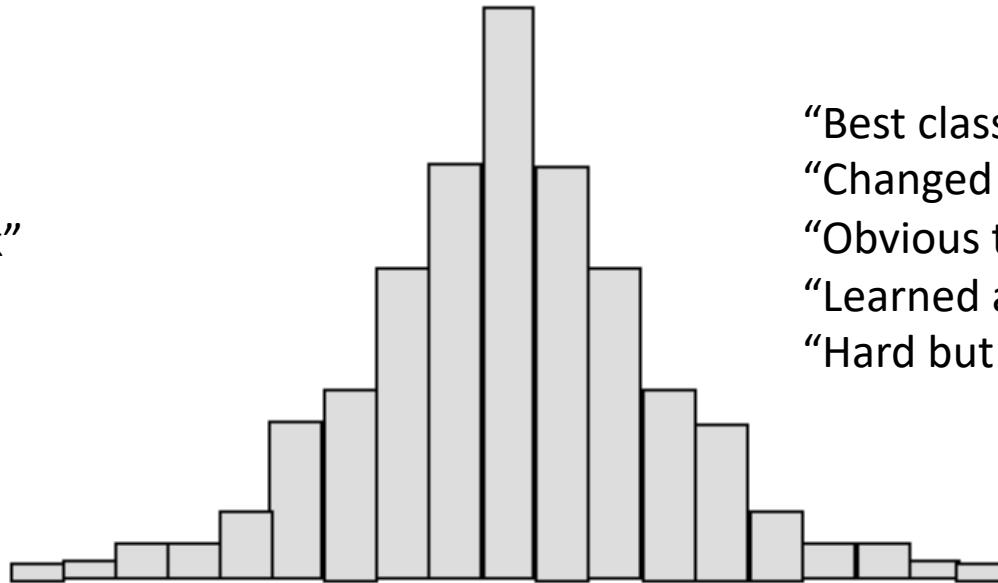


## 8. The course

### A note on your experience 😊

“Terrible class”  
“Don’t take it”  
“Too much work”

“Best class at Stanford”  
“Changed my life”  
“Obvious that instructors care”  
“Learned a ton”  
“Hard but worth it”



- This is a hard, advanced, graduate level class
- I and all the TAs really care about your success in this class
- Give Feedback. Take responsibility for holes in your knowledge
- **Come to office hours/help sessions**

# Office Hours / Help sessions

- **Come to office hours/help sessions!**
  - Come to discuss final project ideas as well as the assignments
  - Try to come early, often and off-cycle
- **Help sessions:** daily, at various times, see calendar
  - Attending in person: Just show up! Our friendly course staff will be on hand to assist you
  - SCPD/remote access: Use queuestatus
- **Chris's office hours:**
  - Mon 4–6pm. Come along next Monday?