

# Natural Language Processing

AI51701/CSE71001

Lecture 4

09/07/2023

Instructor: Taehwan Kim

# Word (review)

# How should we evaluate word vectors?

- ❑ WordSim353 (Finkelstein et al., 2002)

word pair		similarity
journey	voyage	9.3
king	queen	8.6
computer	software	8.5
law	lawyer	8.4
forest	graveyard	1.9
rooster	voyage	0.6

# How should we evaluate word vectors?

- SimLex-999 (Hill et al., 2014)

word pair		similarity
insane	crazy	9.6
attorney	lawyer	9.4
author	creator	8.0
diet	apple	1.2
new	ancient	0.2

measures **paraphrastic** similarity:

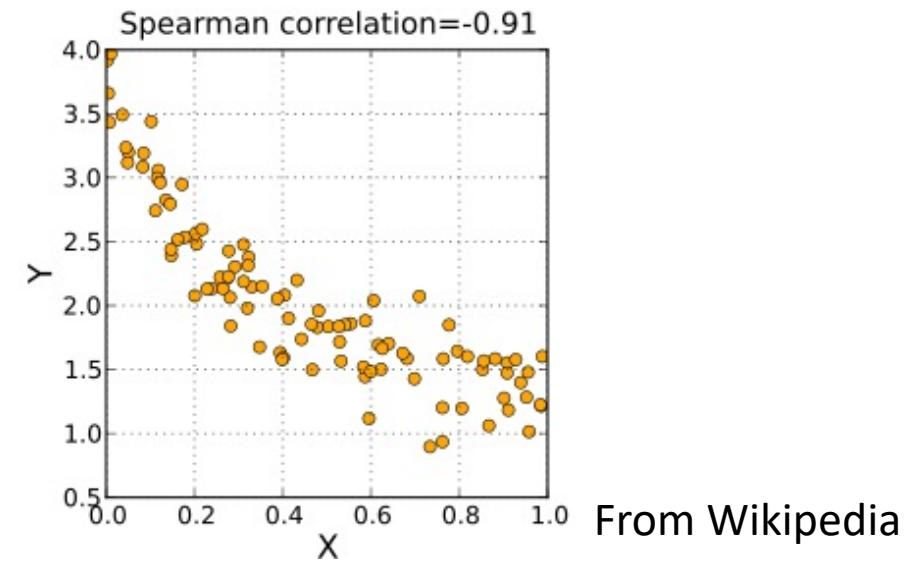
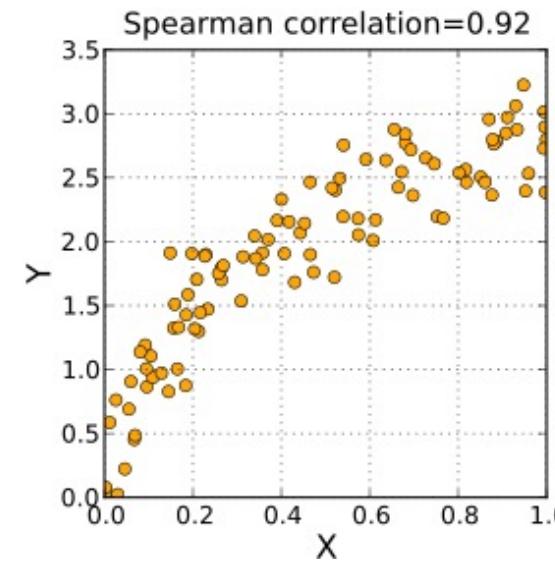
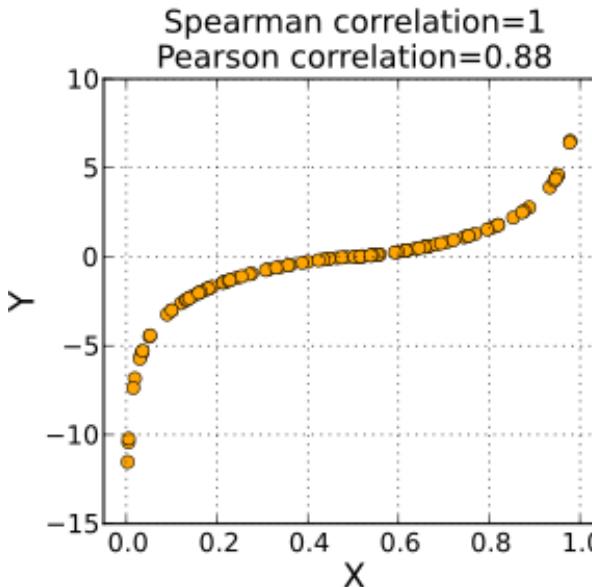
two words are “similar” if they have similar meanings

# How should we evaluate word vectors?

- ❑ There are many word similarity datasets
- ❑ Some focus on topical **relatedness**, others focus on similarity in **meaning**

# Evaluation Metrics for Word Similarity

- ❑ Spearman rank correlation coefficient
- ❑ Measures correlation between two variables:
  - variable 1: human-annotated similarities for word pairs
  - variable 2: cosine similarities computed with your word vectors for the same word pairs



From Wikipedia

# Text Classification

# Text Classification

- ❑ Simplest user-facing NLP application
- ❑ Email (spam, priority, categories):



- ❑ Sentiment:
- ❑ Topic classification
- ❑ Others?

# Text Classification

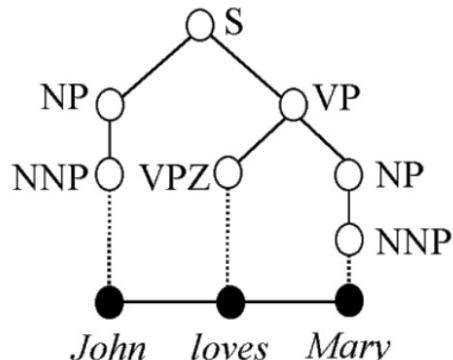
- ❑ Datasets
- ❑ Classification
  - Modeling
  - Inference
  - Learning

# Annotation

- ❑ Supervised machine learning needs labeled datasets, where labels are called **ground truth**
- ❑ In NLP, labels are annotations provided by humans
- ❑ There is always some disagreement among annotators, even for simple tasks
- ❑ These annotations are called a **gold standard**, not ground truth

# How are NLP datasets developed?

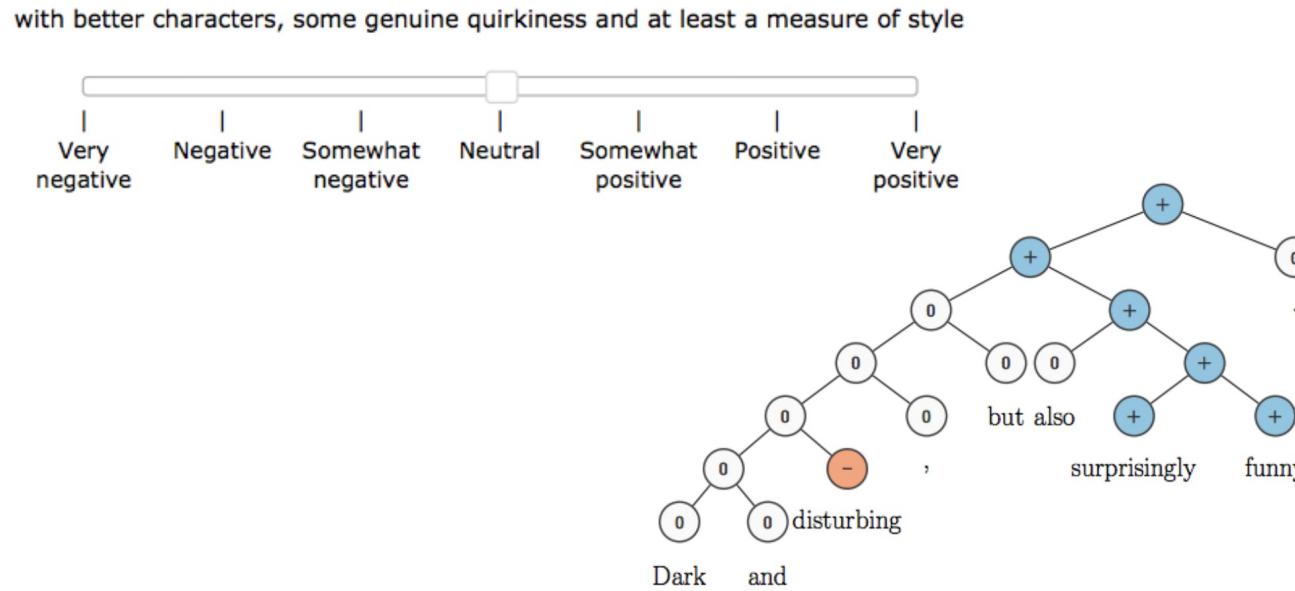
- ❑ Paid, trained human annotators
  - traditional approach
  - researchers write annotation guidelines, recruit & pay annotators (often linguists)
  - more consistent annotations, but costly to scale
  - e.g., Penn Treebank (1993)
    - 1 million words, mostly Wall Street Journal, annotated with part-of-speech tags and syntactic parse trees



# Crowdsourcing

## ❑ Crowdsourcing (e.g., Amazon Mechanical Turk)

- more recent trend
- can't really train annotators, but easier to get multiple annotations for each input (which can then be averaged)
- e.g., Stanford Sentiment Treebank



# Naturally-occurring annotation

- ❑ Long history: used by IBM for speech recognition

There's No Data Like More Data	
• Dick Garwin's correspondence	~2.5M words
• Associated Press	20M words
• Oil company	25M words
• Federal Register	??M words
• American Printing House for the Blind	60M words
• IBM Deposition	100M words
• Canadian Hansard English	100M words

credit: Brown & Mercer, 20 Years of  
Bitext Workshop, 2013

- ❑ How might you find naturally-occurring data for:
  - conversational agents
  - summarization
  - coreference resolution

# Text Classification

- ❑ Datasets
- ❑ **Classification**
  - Modeling
  - Inference
  - Learning

# Modeling, Inference, Learning

modeling: define score function



$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y}{\operatorname{argmax}} \text{ score}(\mathbf{x}, y, \mathbf{w})$$

- ❑ **Modeling:** How do we assign a score to an (x,y) pair using parameters w?

# Modeling, Inference, Learning

**inference:** solve  $\operatorname{argmax}$

**modeling:** define score  $\mathbf{function}$

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_y \text{score}(\mathbf{x}, y, \mathbf{w})$$

- ❑ **Inference:** How do we efficiently search over the space of all labels?

# Modeling, Inference, Learning

**inference:** solve  $\text{argmax}$

**modeling:** define score function

$$\text{classify}(x, w) = \underset{y}{\text{argmax}} \quad \text{score}(x, y, w)$$

**learning:** choose  $w$

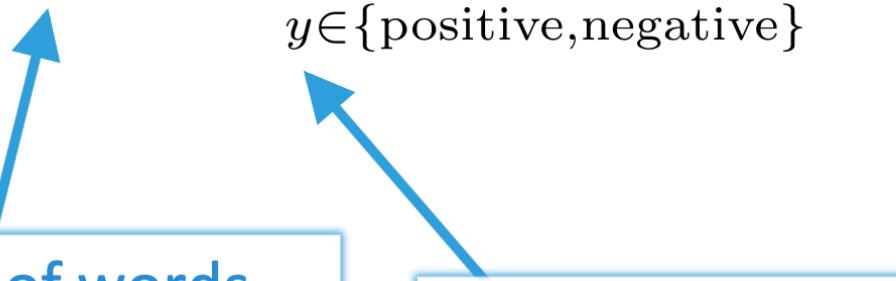
- Learning: How do we choose the weights  $w$ ?

# Binary Sentiment Classification

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y \in \{\text{positive, negative}\}}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$

a sequence of words  
(a sentence or document)

sentiment label



# Binary Sentiment Classification

modeling: define score function

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y \in \{\text{positive, negative}\}}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$

# Linear Models

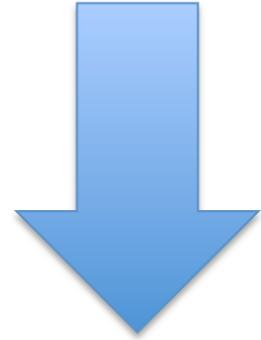
- ❑ Parameters are arranged in a vector  $\mathbf{w}$
- ❑ score function is linear in  $\mathbf{w}$ :

$$\text{score}(\mathbf{x}, y, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) = \sum_i w_i f_i(\mathbf{x}, y)$$

- ❑  $\mathbf{f}$ : vector of feature functions

# Linear Models for Binary Sentiment Classification

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_{y \in \{\text{positive, negative}\}} \text{score}(\mathbf{x}, y, \mathbf{w})$$



$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_{y \in \{\text{positive, negative}\}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y)$$

❑ How do we define  $\mathbf{f}$  ?

# Features for NLP

- ❑ NLP datasets include inputs and outputs
- ❑ Features are usually not included
- ❑ You have to define your own features
- ❑ Contrast this with UCI datasets, which include a fixed-length **dense** feature vector for every instance
- ❑ In (traditional) NLP, features usually **sparse**

# Defining Features

- ❑ This is a large part of NLP
- ❑ Last 30 years: feature engineering
- ❑ Last 10 years: representation learning
  
- ❑ Learning representations doesn't mean that we don't have to look at the data or the output!
- ❑ There's still plenty of engineering required in representation learning

# N-gram

- **N-gram**: a sequence consisting of consecutive n words
- E.g., “Obama is going to”
  - unigram: “Obama”, “is”, “going”, “to”
  - bigram: “Obama is”, “is going”, “going to”
  - trigram: “Obama is going”, “is going to”
- ...

# Unigram Binary Features

- Two example features:

$$f_1(\mathbf{x}, y) = \mathbb{I}[y = \text{positive}] \wedge \mathbb{I}[\mathbf{x} \text{ contains } \text{great}]$$

$$f_2(\mathbf{x}, y) = \mathbb{I}[y = \text{negative}] \wedge \mathbb{I}[\mathbf{x} \text{ contains } \text{great}]$$

where  $\mathbb{I}[S] = 1$  if  $S$  is true, 0 otherwise

- We usually think in terms of feature **templates**
- Unigram binary feature template:

$$f^{\text{u,b}}(\mathbf{x}, y) = \mathbb{I}[y = \text{label}] \wedge \mathbb{I}[\mathbf{x} \text{ contains } \text{word}]$$

- To create features, this feature template is instantiated for particular labels and words

# Higher-Order Binary Feature Templates

- ❑ Unigram binary template:

$$f^{\text{u}, \text{b}}(\boldsymbol{x}, y) = \mathbb{I}[y = \text{label}] \wedge \mathbb{I}[\boldsymbol{x} \text{ contains } word]$$

- ❑ Bigram binary template:

$$f^{\text{b}, \text{b}}(\boldsymbol{x}, y) = \mathbb{I}[y = \text{label}] \wedge \mathbb{I}[\boldsymbol{x} \text{ contains "word1 word2"}]$$

- ❑ Trigram binary template:

...

# Text Classification

- ❑ Datasets
- ❑ Classification
  - Modeling
  - **Inference**
  - Learning

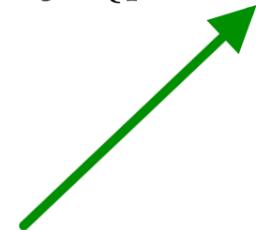
**inference: solve argmax**

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \operatorname{argmax}_y \text{score}(\mathbf{x}, y, \mathbf{w})$$

- ❑ **Inference:** How do we efficiently search over the space of all labels?

# Inference for Text Classification

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y \in \{\text{positive, negative}\}}{\operatorname{argmax}} \text{score}(\mathbf{x}, y, \mathbf{w})$$



❑ trivial (loop over labels)

# Text Classification

- ❑ Datasets
- ❑ Classification
  - Modeling
  - Inference
  - Learning

# Modeling, Inference, Learning

**inference:** solve  $\operatorname{argmax}$

**modeling:** define score function

$$\text{classify}(x, w) = \operatorname{argmax}_y \text{score}(x, y, w)$$

**learning:** choose  $w$

- **Learning:** How should we choose values for the weights  $w$ ?

# Text Classification

- ❑ Modeling
- ❑ Inference
- ❑ Learning
  - **empirical risk minimization**
  - surrogate loss functions
  - gradient-based optimization

# Cost Functions

- ❑ Cost function: scores outputs against a gold standard

$$\text{cost} : \mathcal{L} \times \mathcal{L} \rightarrow \mathbb{R}_{\geq 0}$$

- ❑ Should be as close as possible to the actual evaluation metric for your task
- ❑ for classification, what cost should we use?

$$\text{cost}(y, y') = \mathbb{I}[y \neq y']$$

# Risk Minimization

- given training data:  $\mathcal{T} = \{\langle \mathbf{x}^{(i)}, y^{(i)} \rangle\}_{i=1}^{|\mathcal{T}|}$

where each  $y^{(i)} \in \mathcal{L}$  is a label

- assume data is drawn iid (independently and identically distributed) from (unknown) joint distribution  $P(\mathbf{x}, y)$
- we want to solve the following:

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{P(\mathbf{x}, y)} [\text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))]$$



problem:  $P$  is unknown

# Empirical Risk Minimization (Vapnik et al.)

- ❑ Replace expectation with sum over examples:

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \mathbb{E}_{P(\mathbf{x}, y)} [\text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))]$$



$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} \sum_{i=1}^{|\mathcal{T}|} \text{cost}(y^{(i)}, \text{classify}(\mathbf{x}^{(i)}, \mathbf{w}))$$

problem: NP-hard even for binary classification with linear models

□ Solution: replace “cost loss” (also called “0-1” loss) with a surrogate function that is easier to optimize

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^{|\mathcal{T}|} \text{cost}(y^{(i)}, \text{classify}(\mathbf{x}^{(i)}, \mathbf{w}))$$



generalize to permit any loss function

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^{|\mathcal{T}|} \text{loss}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w})$$

□ Cost loss / 0-1 loss:  $\text{loss}_{\text{cost}}(\mathbf{x}, y, \mathbf{w}) = \text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))$

# Text Classification

- ❑ Modeling
- ❑ Inference
- ❑ Learning
  - empirical risk minimization
  - **surrogate loss functions**
  - gradient-based optimization

# Surrogate Loss Functions

- Cost loss / 0-1 loss:  $\text{loss}_{\text{cost}}(\mathbf{x}, y, \mathbf{w}) = \text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))$ 
  - why is this so difficult to optimize? not necessarily continuous, can't use gradient-based optimization
- max-score loss:  $\text{loss}_{\text{maxscore}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w})$
- perceptron loss:  $\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$ 
  - loss function underlying perceptron algorithm (Rosenblatt, 1957-58)

# Surrogate Loss Functions

- Cost loss / 0-1 loss:  $\text{loss}_{\text{cost}}(\mathbf{x}, y, \mathbf{w}) = \text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))$
- perceptron loss:  $\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$

- hinge loss:

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))$$

- loss function underlying support vector machines

# Log Loss

$$\text{loss}_{\log}(\mathbf{x}, y, \mathbf{w}) = -\log p_{\mathbf{w}}(y \mid \mathbf{x})$$

- Minimize negative log of conditional probability of output given input
  - sometimes called “maximizing conditional likelihood”
- But we don’t have a probabilistic model, we just have a score function

# Score → Probability

- Can turn score into probability by exponentiating (to make it positive) and normalizing:

$$p_{\mathbf{w}}(y \mid \mathbf{x}) \propto \exp\{\text{score}(\mathbf{x}, y, \mathbf{w})\}$$

$$p_{\mathbf{w}}(y \mid \mathbf{x}) = \frac{\exp\{\text{score}(\mathbf{x}, y, \mathbf{w})\}}{\sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \mathbf{w})\}}$$

- This is often called a “softmax” function

# Log Loss

$$\begin{aligned}\text{loss}_{\log}(\mathbf{x}, y, \mathbf{w}) &= -\log p_{\mathbf{w}}(y \mid \mathbf{x}) \\ &= -\log \frac{\exp\{\text{score}(\mathbf{x}, y, \mathbf{w})\}}{\sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \mathbf{w})\}} \\ &= -\text{score}(\mathbf{x}, y, \mathbf{w}) + \log \sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \mathbf{w})\}\end{aligned}$$

- ❑ Similar to perceptron loss!
- ❑ Replace max with “softmax” (a different kind of softmax)

$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$

# Log Loss

$$\begin{aligned}\text{loss}_{\log}(\mathbf{x}, y, \mathbf{w}) &= -\log p_{\mathbf{w}}(y \mid \mathbf{x}) \\ &= -\log \frac{\exp\{\text{score}(\mathbf{x}, y, \mathbf{w})\}}{\sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \mathbf{w})\}} \\ &= -\text{score}(\mathbf{x}, y, \mathbf{w}) + \log \sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \mathbf{w})\}\end{aligned}$$

- Log loss is used in:
  - logistic regression classifiers,
  - conditional random fields,
  - maximum entropy (“maxent”) models

# Regularized Empirical Risk Minimization

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{P(\mathbf{x}, y)} [\text{cost}(y, \text{classify}(\mathbf{x}, \mathbf{w}))]$$

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^{|\mathcal{T}|} \text{loss}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}) + \lambda R(\mathbf{w})$$

**regularization strength**

**regularization term**

# Regularization Terms

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^{|\mathcal{T}|} \text{loss}(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w}) + \lambda R(\mathbf{w})$$

- Most common: penalize large parameter values
- Intuition: large parameters might be instances of overfitting
- Examples:
  - **$L_2$  regularization:** (also called Tikhonov regularization or ridge regression)
  - **$L_1$  regularization:** (also called basis pursuit or LASSO)

# Regularization Terms

□  **$L_2$  regularization:**  $R_{L2}(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_i w_i^2$

- differentiable, widely-used

□  **$L_1$  regularization:**  $R_{L1}(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_i |w_i|$

- not differentiable (but is subdifferentiable)
- leads to sparse solutions (many parameters become zero!)

# Text Classification

- ❑ Modeling
- ❑ Inference
- ❑ Learning
  - empirical risk minimization
  - surrogate loss functions
  - **gradient-based optimization**

# Gradient Descent

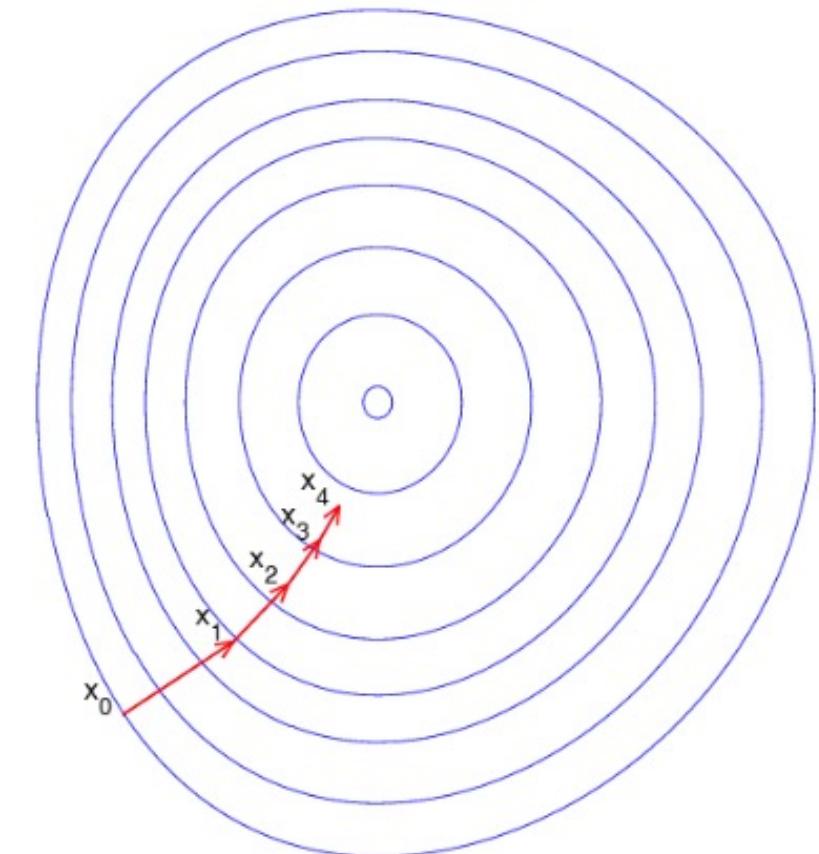
- ❑ Minimizes a function  $F$  by taking steps in proportion to the negative of the gradient:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta^{(t)} \nabla F(\boldsymbol{\theta}^{(t)})$$

$\eta^{(t)}$ : stepsize at iteration t

$\nabla F(\boldsymbol{\theta}^{(t)})$ : gradient of objective function

- ❑ With conditions on stepsize and objective function, will converge to local minimum



# Gradient Descent

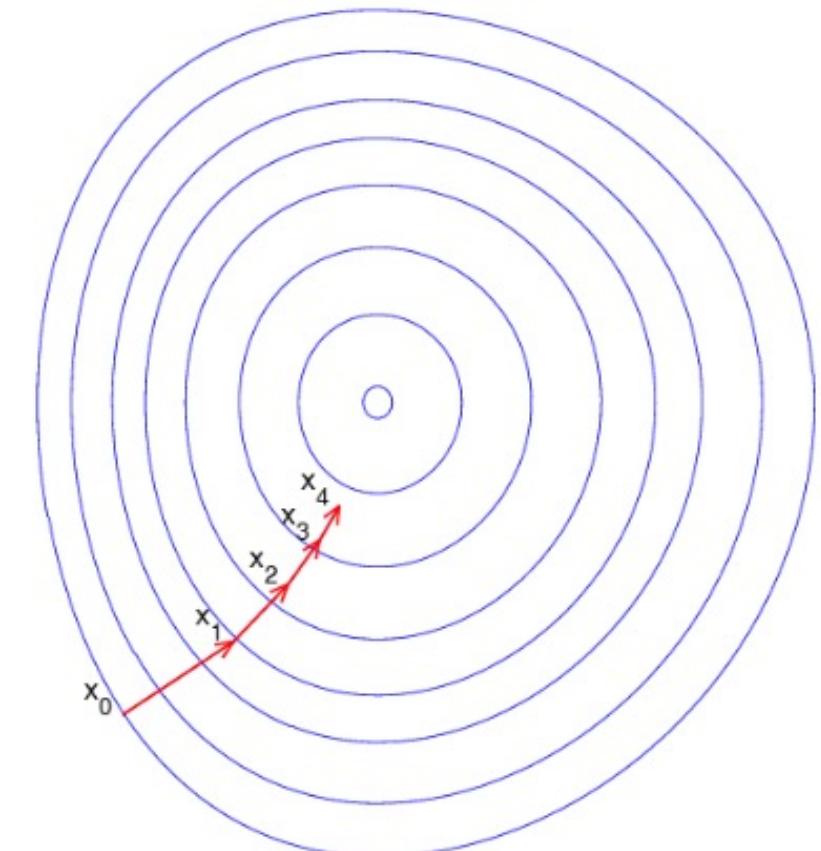
- ❑ Minimizes a function  $F$  by taking steps in proportion to the negative of the gradient:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta^{(t)} \nabla F(\boldsymbol{\theta}^{(t)})$$

$\eta^{(t)}$ : stepsize at iteration t

$\nabla F(\boldsymbol{\theta}^{(t)})$ : gradient of objective function

- ❑ Efficiency concern:  $F$  is a sum over all training examples! every parameter update requires iterating through entire training set. **“batch” algorithm**



# Batch Gradient Descent

## ❑ Batch Gradient Descent Algorithm

Input: data  $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}$  ( $(1 \leq i \leq N)$ ), learning rate  $\alpha_t$ .

- (I)  $t = 0$ ; Make an initial guess  $\theta_0$ ;
- (II) Iterate until termination condition is met

- (i)  $\theta_{t+1} = \theta_t - \alpha_t \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \ell(f(\mathbf{x}^{(i)}; \theta_t), \mathbf{y}^{(i)})$ ;
  - (ii)  $t = t + 1$ ;

# Stochastic Gradient Descent

## ❑ Stochastic Gradient Descent Algorithm

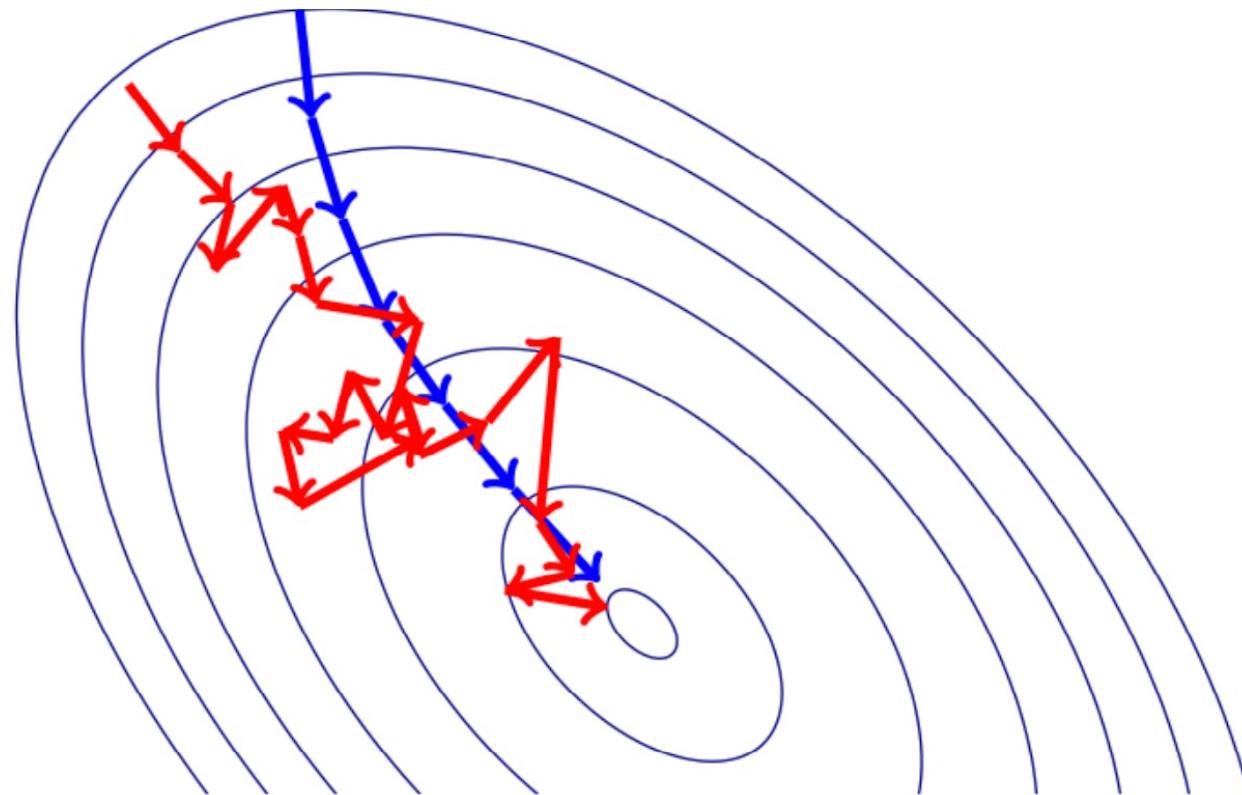
Input: data  $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}$  ( $(1 \leq i \leq N)$ ), learning rate  $\alpha_t$ .

- (I)  $t = 0$ ; Make an initial guess  $\theta_0$ ;
- (II) Iterate until termination condition is met
  - (i) Sample example  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  from training set
  - (ii)  $\theta_{t+1} = \theta_t - \alpha_t \sum_{i=1}^N \nabla_{\theta} \ell(f(\mathbf{x}^{(i)}; \theta_t), \mathbf{y}^{(i)})$ ;
  - (iii)  $t = t + 1$ ;

# Stochastic Gradient Descent

- ❑ Applicable when objective function is a sum
- ❑ Like gradient descent, except calculates gradient on a single example at a time (“online”) or on a small set of examples (“mini-batch”)
- ❑ Converges much faster than (batch) gradient descent
- ❑ With conditions on stepsize and objective function, will converge to local minimum
- ❑ There are many popular variants:
  - SGD+momentum, AdaGrad, AdaDelta, Adam, RMSprop, etc.

# Stochastic Gradient Descent



# What if F is not differentiable?

- Some loss functions are not differentiable:

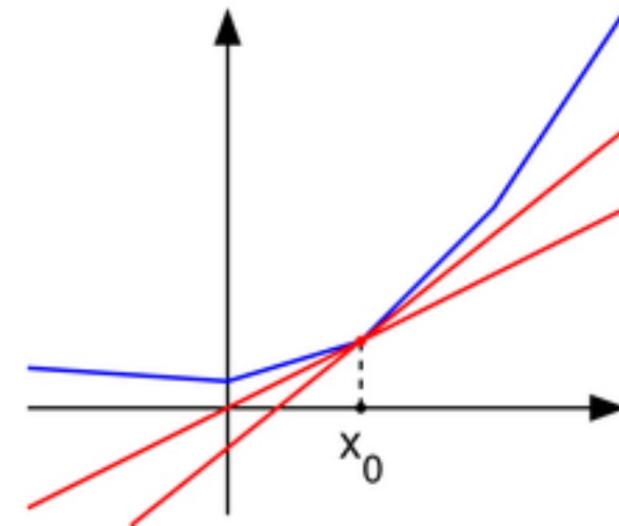
$$\text{loss}_{\text{perc}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \mathbf{w})$$

$$\text{loss}_{\text{hinge}}(\mathbf{x}, y, \mathbf{w}) = -\text{score}(\mathbf{x}, y, \mathbf{w}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \mathbf{w}) + \text{cost}(y, y'))$$

- But they *are subdifferentiable*, so we can compute **subgradients** and use (stochastic) subgradient method (not descent!)
  - So we need to keep track of the best solution so far

# Subderivatives

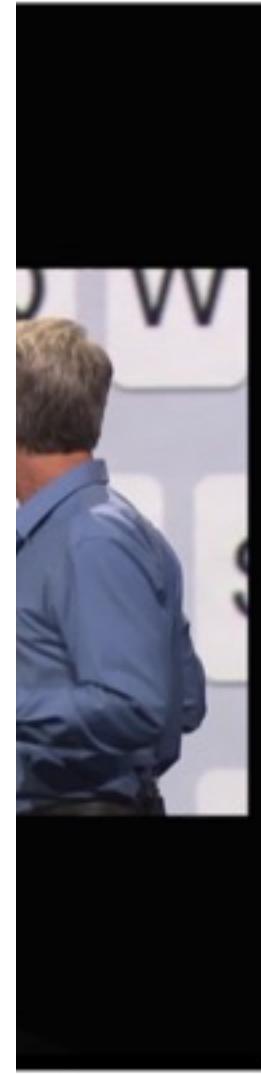
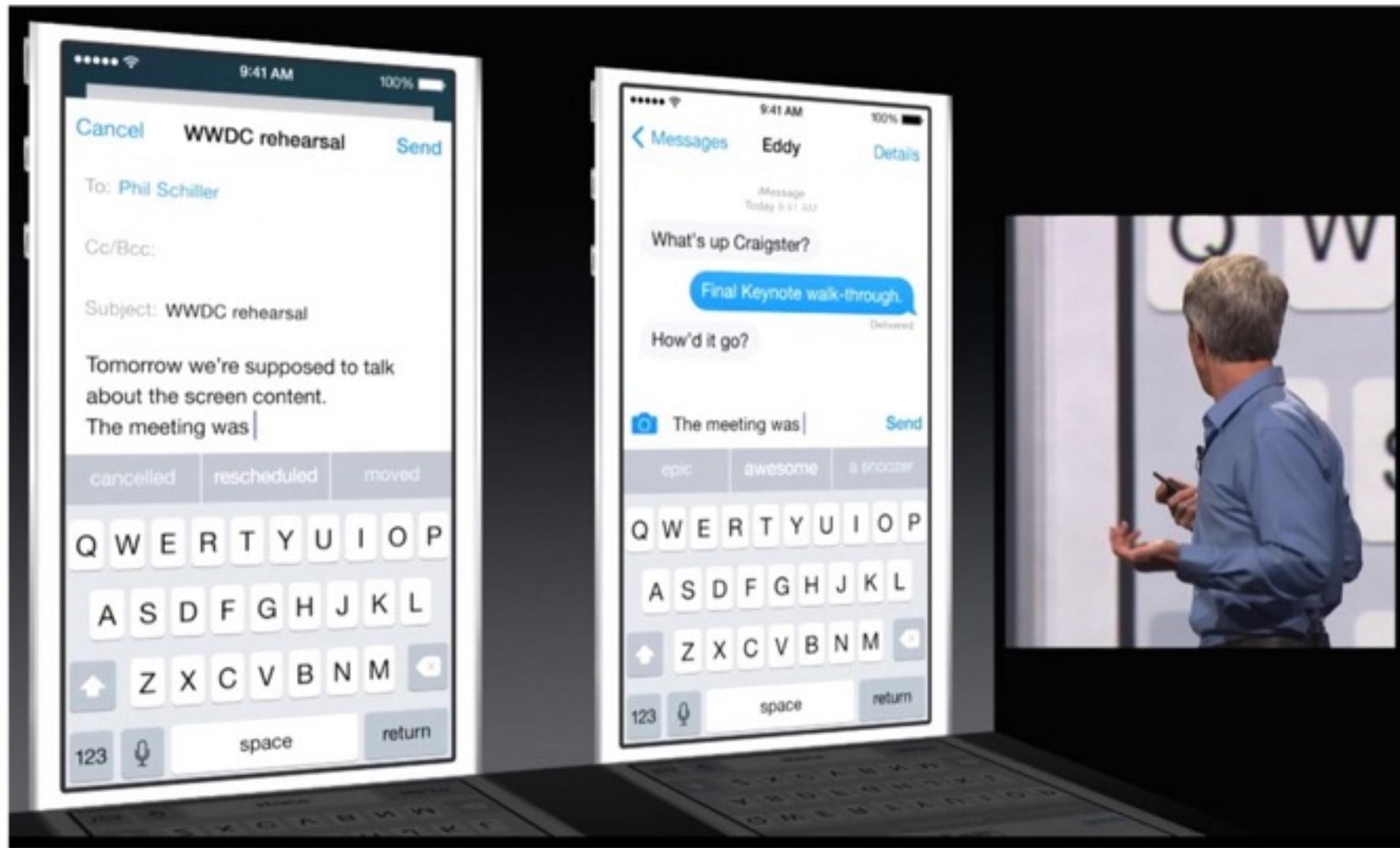
- **Subderivative:** generalization of derivative for nondifferentiable, convex functions
- There may be multiple subderivatives at a point (red lines)
- This set is called the **subdifferential**
- A convex function  $g$  is differentiable at point  $x_0$  if and only if the subdifferential of  $g$  at  $x_0$  contains only the derivative of  $g$  at  $x_0$



# Language Modeling

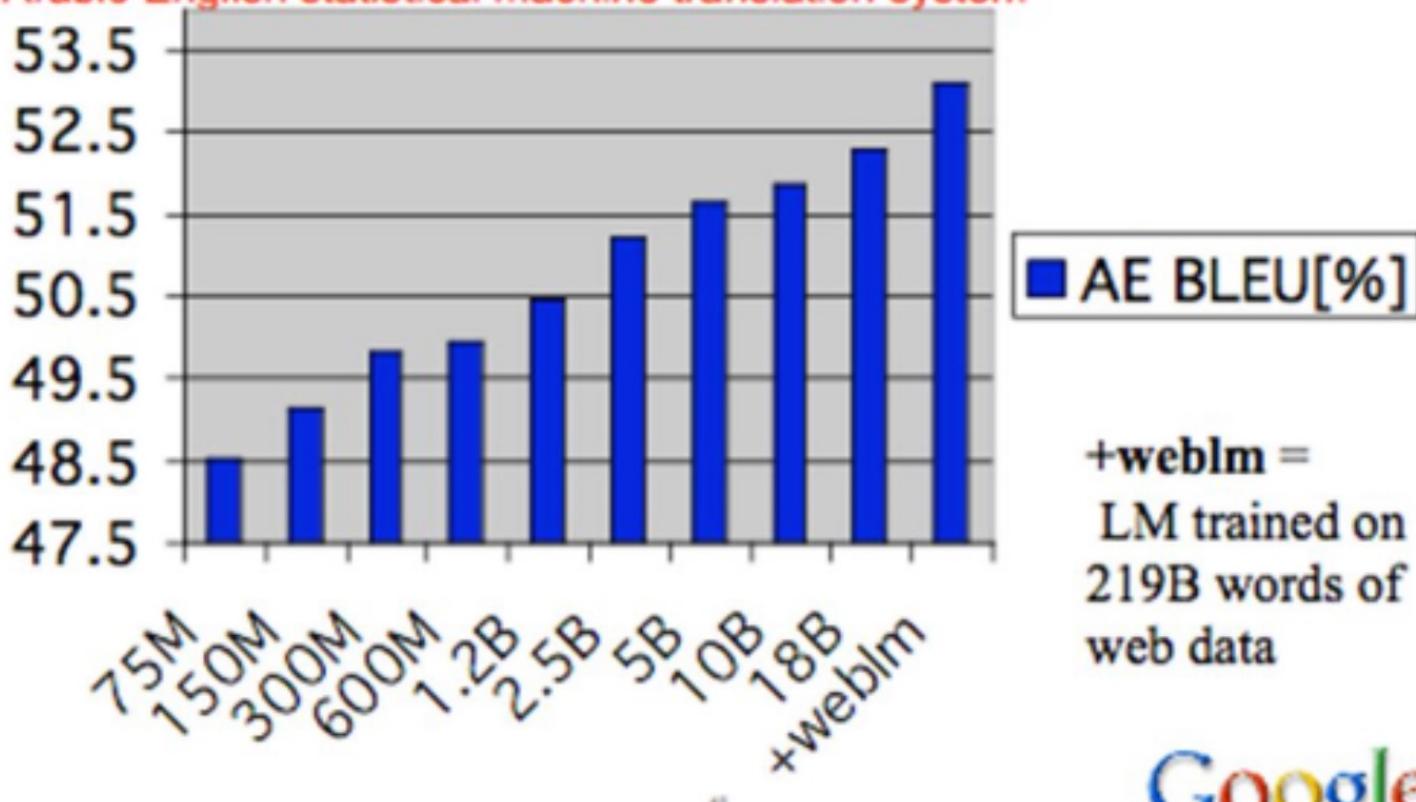
# Probabilistic Language Models

- Language modeling: assign probabilities to sentences
- Why?
  - machine translation:
    - $P(\text{high winds tonite}) > P(\text{large winds tonite})$
  - spelling correction:
    - The office is about fifteen **minuets** from my house
  - speech recognition:
    - $P(\text{I saw a van}) \gg P(\text{eyes awe of an})$
  - summarization, question answering, etc.!



# Language Modeling for Machine Translation

Impact on size of language model training data (in words) on quality of Arabic-English statistical machine translation system



+webLM =  
LM trained on  
219B words of  
web data

Google

# Probabilistic Language Modeling

- ❑ Goal: compute the probability of a sequence of words:

$$P(\mathbf{w}) = P(w_1, w_2, \dots, w_n)$$

- ❑ Related task: probability of next word:

$$P(w_4 \mid w_1, w_2, w_3)$$

- ❑ A model that computes either of these:

$$P(\mathbf{w}) \text{ or } P(w_k \mid w_1, w_2, \dots, w_{k-1})$$

is called a **language model (LM)**

# How to compute $P(w)$

- ❑ How to compute this joint probability:
  - $P(\text{its, water, is, so, transparent, that})$
- ❑ Intuition: let's rely on the Chain Rule of Probability

# Reminder: Chain Rule

- ❑ Factor joint probability into product of conditional probabilities:

$$P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1, w_2) \dots P(w_n \mid w_1, w_2, \dots, w_{n-1})$$

- ❑ We have not yet made any independence assumptions

# Chain Rule for computing joint probability of words in sentence

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i | w_1, w_2, \dots, w_{i-1})$$

$P(\text{"its water is so transparent"}) =$

$P(\text{its}) \times P(\text{water} | \text{its}) \times P(\text{is} | \text{its water})$

$\times P(\text{so} | \text{its water is}) \times P(\text{transparent} | \text{its water is so})$

# How to estimate these probabilities

- ❑ Could we just count and divide?

$P(\text{the lits water is so transparent that}) =$

$\frac{\text{Count}(\text{its water is so transparent that the})}{\text{Count}(\text{its water is so transparent that})}$

- ❑ No! too many possible sentences!

- ❑ We'll never see enough data for estimating these

# Markov Assumption

- ❑ Simplifying assumption:

$$P(\text{the lit water is so transparent that}) \approx P(\text{the l that})$$

- ❑ Or maybe:

$$P(\text{the lit water is so transparent that}) \approx P(\text{the l transparent that})$$

- ❑ I.e., we approximate each component in the product:

$$P(w_i | w_1, \dots, w_{i-2}, w_{i-1}) \approx P(w_i | w_{i-k}, \dots, w_{i-2}, w_{i-1})$$

# Simplest case: Unigram model

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i)$$

- Automatically generated sentences from a unigram model:

fifth an of futures the an incorporated a a the  
inflation most dollars quarter in is mass

thrift did eighty said hard 'm july bullish

that or limited the

# Bigram model

- ❑ condition on the previous word:

$$P(w_1, w_2, \dots, w_n) = \prod_i P(w_i | w_{i-1})$$

- ❑ Automatically generated sentences from a unigram model:

texaco rose one in this issue is pursuing growth in a boiler  
house said mr. gurria mexico 's motion control proposal  
without permission from five hundred fifty five yen

outside new car parking lot of the agreement reached  
this would be a record november

# n-gram models

- ❑ We can extend to trigrams, 4-grams, 5-grams
- ❑ In general this is an insufficient model of language
  - because language has **long-distance dependencies**:  
*“The computer which I had just put into the machine room on the  
fi6h floor crashed.”*
- ❑ But we can often get away with n-gram models

# Estimating bigram probabilities

- The maximum likelihood estimate (MLE)

$$P(w_i \mid w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

# An example

$$P(w_i \mid w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(\text{I} \mid \text{<s>}) = \frac{2}{3} = .67 \quad P(\text{Sam} \mid \text{<s>}) = \frac{1}{3} = .33 \quad P(\text{am} \mid \text{I}) = \frac{2}{3} = .67$$

$$P(\text{</s>} \mid \text{Sam}) = \frac{1}{2} = 0.5 \quad P(\text{Sam} \mid \text{am}) = \frac{1}{2} = .5 \quad P(\text{do} \mid \text{I}) = \frac{1}{3} = .33$$

# More examples:

## Berkeley Restaurant Project sentences

- Can you tell me about any good Cantonese restaurants close by
- Mid priced thai food is what i'm looking for
- Can you give me a listing of the kinds of food that are available

# Raw bigram counts

- ❑ Counts from 9,222 sentences
- ❑ E.g., “*i want*” occurs 827 times

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

# Raw bigram probabilities

- ❑ normalize by unigram counts:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

- ❑ bigram probabilities

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

# Bigram estimates of sentence probabilities

$$P(< s > | I \text{ want } \textit{english} \text{ food} < /s >) =$$

$$P(I | < s >)$$

$$\times P(\textit{want} | I)$$

$$\times P(\textit{english} | \textit{want})$$

$$\times P(\textit{food} | \textit{english})$$

$$\times P(< /s > | \textit{food})$$

$$= .000031$$

# Practical Issues

- ❑ we do everything in log space
  - avoid underflow
  - also adding is faster than multiplying

$$\log(p_1 \times p_2 \times p_3) = \log p_1 + \log p_2 + \log p_3$$

# Language Modeling Toolkits

## ❑ SRILM

- <http://www.speech.sri.com/projects/srilm/>

## ❑ KenLM

- <https://kheafield.com/code/kenlm/>

# Evaluation: How good is our model?

- ❑ Does our language model prefer good sentences to bad ones?
  - assign higher probability to “real” or “frequently observed” sentences
  - than “ungrammatical” or “rarely observed” sentences?

# Extrinsic evaluation of N-gram models

- ❑ Best evaluation for comparing models A and B
  - put each model in a task
    - spelling corrector, speech recognizer, MT system
  - run the task, get an accuracy for A and for B
    - how many misspelled words corrected properly
    - how many words translated correctly
  - compare accuracy for A and B

# Intuition of Perplexity

## ❑ The Shannon Game:

- how well can we predict the next word?

*I always order pizza with cheese and \_\_\_\_\_*

*The 33<sup>rd</sup> President of the US was \_\_\_\_\_*

*I saw a \_\_\_\_\_*



*mushrooms* 0.1  
*pepperoni* 0.1  
*anchovies* 0.01  
...  
*fried rice* 0.0001  
...  
*and* 1e-100

- unigrams are terrible at this game (why?)

## ❑ A better model of a text is one which assigns a higher probability to the word that actually occurs

# Probability of Held-out Data

- ❑ probability of held-out sentences:

$$\prod_i P(\mathbf{w}^{(i)})$$

- ❑ let's work with log-probabilities:

$$\log_2 \prod_i P(\mathbf{w}^{(i)}) = \sum_i \log_2 P(\mathbf{w}^{(i)})$$

- ❑ divide by number of words  $M$  in held-out sentences:

$$\frac{1}{M} \sum_i \log_2 P(\mathbf{w}^{(i)})$$

# Probability -> Perplexity

- ❑ average log-probability of held-out words:

$$\ell = \frac{1}{M} \sum_i \log_2 P(\mathbf{w}^{(i)})$$

- ❑ perplexity:

$$PP = 2^{-\ell}$$

# Perplexity as branching factor

- ❑ Given a sentence consisting of random digits
- ❑ Perplexity of this sentence under a model that gives probability 1/10 to each digit?

$$\begin{aligned}\ell &= \frac{1}{M} \log_2 P(w_1, w_2, \dots, w_M) \\ &= \frac{1}{M} \log_2 \prod_{i=1}^M \frac{1}{10} && \text{PP} = 2^{-\ell} = 10 \\ &= \frac{1}{M} \log_2 \left( \frac{1}{10} \right)^M \\ &= \frac{1}{M} M \log_2 \frac{1}{10}\end{aligned}$$

# Lower perplexity = better model

- ❑ Train: 38 million words
- ❑ Test: 1.5 million words

n-gram order:	unigram	bigram	trigram
perplexity:	962	170	109

# Approximating Shakespeare

- |           |  |
|-----------|--|
| 1<br>gram | –To him swallowed confess hear both. Which. Of save on trail for are ay device and<br>rote life have<br>–Hill he late speaks; or! a more to leg less first you enter               |
| 2<br>gram | –Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live<br>king. Follow.<br>–What means, sir. I confess she? then all sorts, he is trim, captain. |
| 3<br>gram | –Fly, and will rid me these news of price. Therefore the sadness of parting, as they say,<br>'tis done.<br>–This shall forbid it should be branded, if renown made it empty.       |
| 4<br>gram | –King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A<br>great banquet serv'd in;<br>–It cannot be but so.   |

## A Neural Probabilistic Language Model

**Yoshua Bengio**

**Réjean Ducharme**

**Pascal Vincent**

**Christian Jauvin**

*Département d’Informatique et Recherche Opérationnelle*

*Centre de Recherche Mathématiques*

*Université de Montréal, Montréal, Québec, Canada*

BENGIOY@IRO.UMONTREAL.CA

DUCHARME@IRO.UMONTREAL.CA

VINCENTP@IRO.UMONTREAL.CA

JAUVINC@IRO.UMONTREAL.CA

- ❑ Idea: use a neural network for n-gram language modeling:

$$P_{\theta}(w_t \mid w_{t-n+1}, \dots, w_{t-2}, w_{t-1})$$

## A Neural Probabilistic Language Model

**Yoshua Bengio**

**Réjean Ducharme**

**Pascal Vincent**

**Christian Jauvin**

*Département d’Informatique et Recherche Opérationnelle*

*Centre de Recherche Mathématiques*

*Université de Montréal, Montréal, Québec, Canada*

BENGIOY@IRO.UMONTREAL.CA

DUCHARME@IRO.UMONTREAL.CA

VINCENTP@IRO.UMONTREAL.CA

JAUVINC@IRO.UMONTREAL.CA

- ❑ This is not the earliest paper on using neural networks for n-gram language models, but it's the most well-known and first to scale up
- ❑ See paper for citations of earlier work

# Neural Probabilistic Language Models (Bengio et al., 2003)

## 1.1 Fighting the Curse of Dimensionality with Distributed Representations

In a nutshell, the idea of the proposed approach can be summarized as follows:

1. associate with each word in the vocabulary a distributed *word feature vector* (a real-valued vector in  $\mathbb{R}^m$ ),
2. express the joint *probability function* of word sequences in terms of the feature vectors of these words in the sequence, and
3. learn simultaneously the *word feature vectors* and the parameters of that *probability function*.

# What is a neural network?

- ❑ Just think of a neural network as a function
- ❑ It has inputs and outputs
- ❑ “neural” typically means one type of functional building block (“neural layers”), but the term has broadened
- ❑ Neural modeling is now better thought of as a modeling strategy (leveraging “distributed representations” or “representation learning”), or a family of related methods

# Notation

**u** = a vector

$u_i$  = entry  $i$  in the vector

**W** = a matrix

$w_{ij}$  = entry  $(i,j)$  in the matrix

**x** = a structured object

$x_i$  = entry  $i$  in the structured object

# neural layer = affine transform + nonlinearity

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

The diagram illustrates the components of the neural layer equation. A blue arrow points from the term  $g$  to the right, labeled "nonlinearity". A blue brace curves under the terms  $\mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)}$ , labeled "affine transform".

- This is a single “layer” of a neural network
- Input vector is  $\mathbf{x}$
- $\mathbf{U}^{(0)}$  and  $\mathbf{b}^{(0)}$  are parameters

# neural layer = affine transform + nonlinearity

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$



"hidden units"

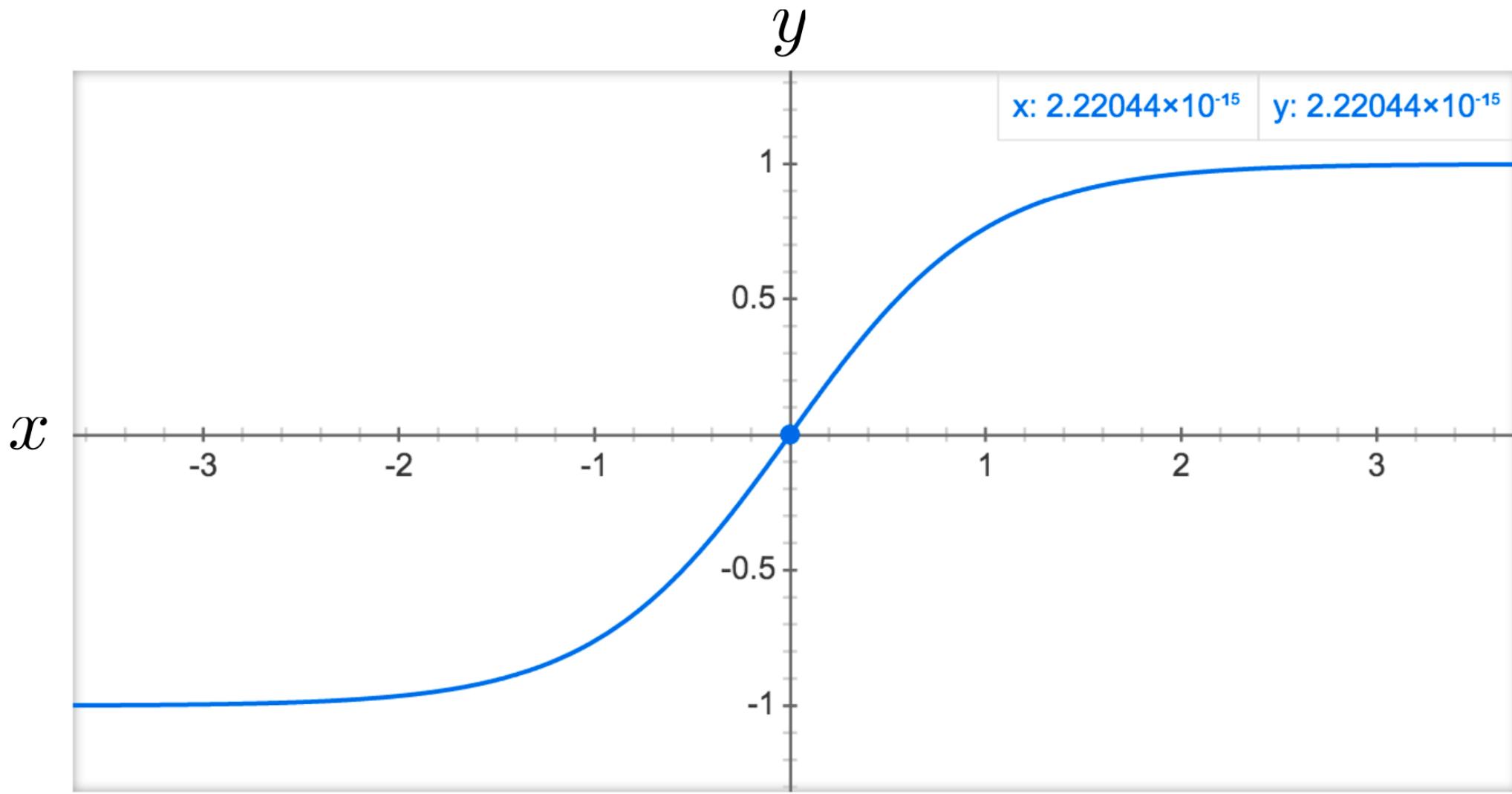
- Vector of “hidden units” is  $\mathbf{z}^{(1)}$
- Think of these as features computed from  $\mathbf{X}$

# Nonlinearities

$$\mathbf{z}^{(1)} = g(\mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)})$$

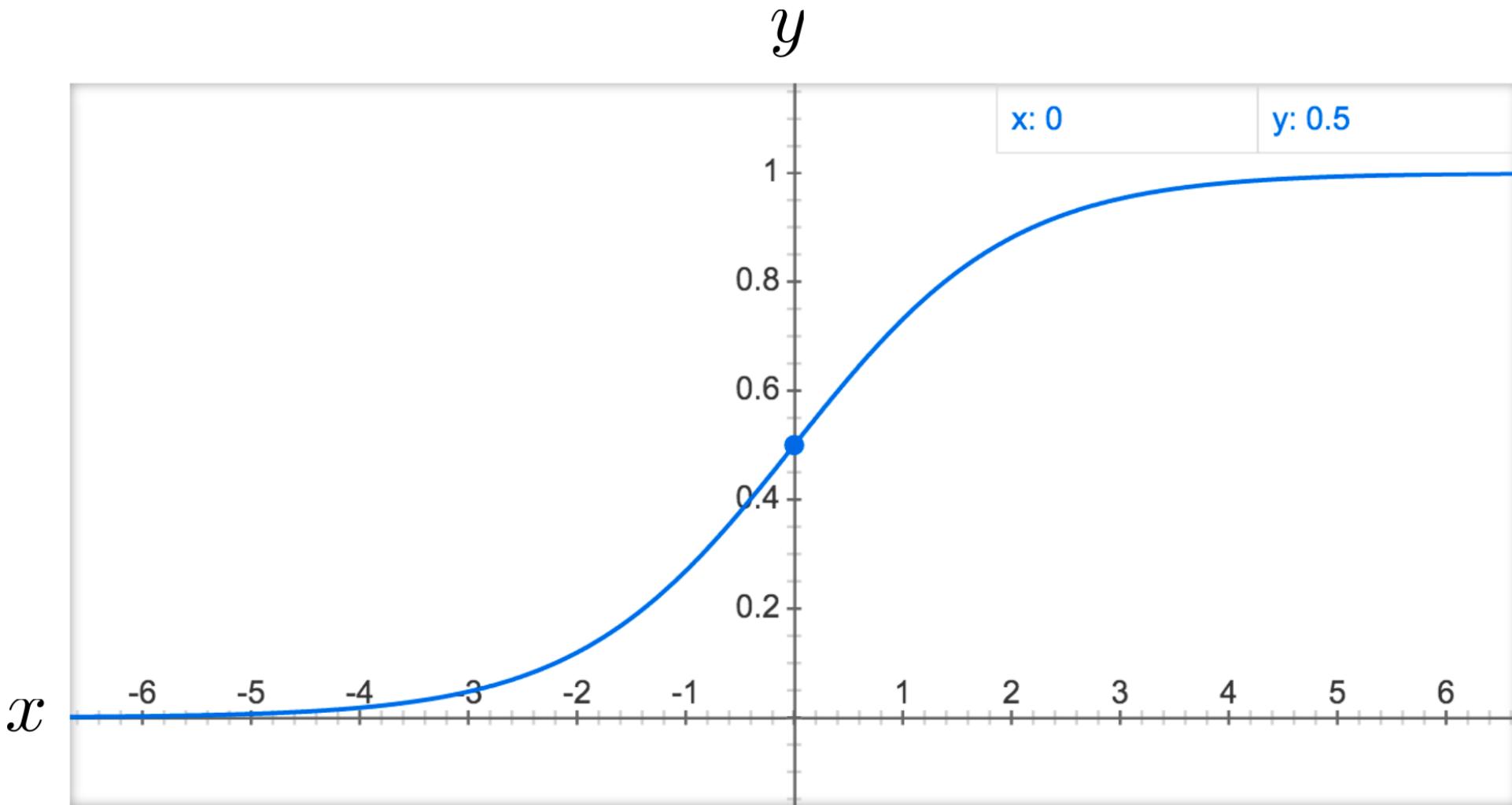
- Most common: elementwise application of  $g$  function to each entry in vector
- Examples...

tanh:  $y = \tanh(x)$

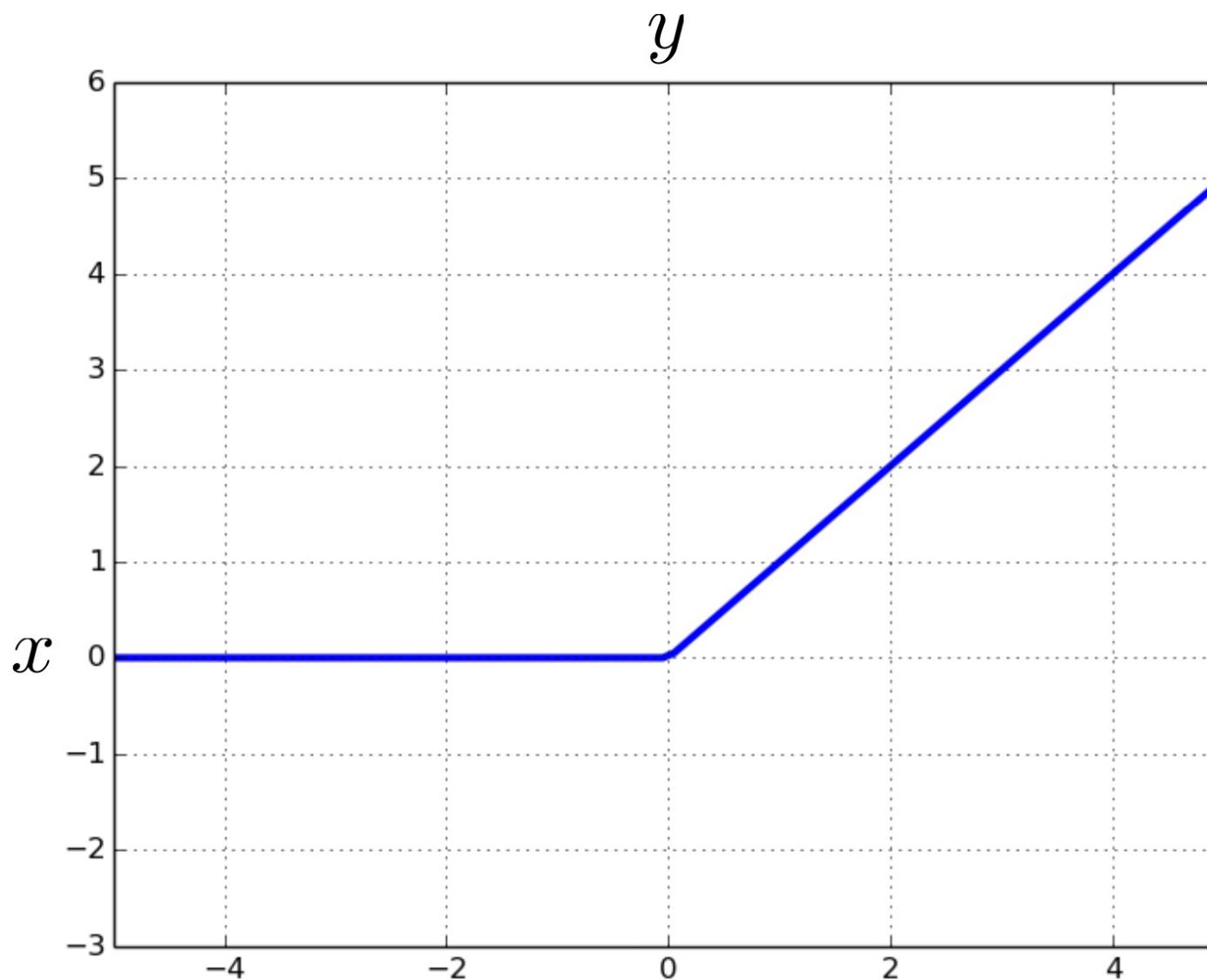


(logistic) sigmoid:

$$y = \frac{1}{1 + \exp\{-x\}}$$



rectified linear unit (ReLU):  $y = \max(0, x)$



# Adding layers...

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

$$\mathbf{z}^{(2)} = g \left( \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)} \right)$$

...

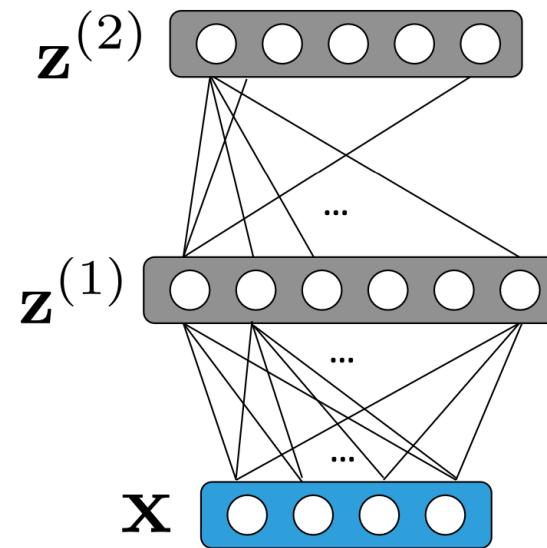
- ❑ Use output of one layer as input to next

# Adding layers...

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

$$\mathbf{z}^{(2)} = g \left( \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)} \right)$$

...



- Use output of one layer as input to next
- “feed-forward” and/or “fully-connected” layers

# Neural Network for Sentiment Classification

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

$$\mathbf{s} = \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$



vector of label scores

# Neural Network for Sentiment Classification

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

$$\mathbf{s} = \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$



$$\mathbf{s} = \begin{bmatrix} \text{score}(\mathbf{x}, \text{positive}, \mathbf{w}) \\ \text{score}(\mathbf{x}, \text{negative}, \mathbf{w}) \end{bmatrix}$$

Use softmax function to convert scores into probabilities

$$\mathbf{s} = \begin{bmatrix} \text{score}(\mathbf{x}, \text{positive}, \mathbf{w}) \\ \text{score}(\mathbf{x}, \text{negative}, \mathbf{w}) \end{bmatrix}$$

$$\mathbf{p} = \text{softmax}(\mathbf{s}) = \begin{bmatrix} \frac{\exp\{\text{score}(\mathbf{x}, \text{positive}, \mathbf{w})\}}{Z} \\ \frac{\exp\{\text{score}(\mathbf{x}, \text{negative}, \mathbf{w})\}}{Z} \end{bmatrix}$$

$$Z = \exp\{\text{score}(\mathbf{x}, \text{positive}, \mathbf{w})\} + \exp\{\text{score}(\mathbf{x}, \text{negative}, \mathbf{w})\}$$

# Why nonlinearities?

network with  
1 hidden layer:

$$\begin{aligned}\mathbf{z}^{(1)} &= g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) \\ \mathbf{s} &= \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)}\end{aligned}$$

- If  $g$  is linear, then we can rewrite the above as a single affine transform
- Can you prove this? (use distributivity of matrix multiplication)

# Learning with Neural Networks

$$\mathbf{z}^{(1)} = g \left( \mathbf{U}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

$$\mathbf{s} = \mathbf{U}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$


$$\mathbf{s} = \begin{bmatrix} \text{score}(\mathbf{x}, \text{positive}, \mathbf{w}) \\ \text{score}(\mathbf{x}, \text{negative}, \mathbf{w}) \end{bmatrix}$$

$$\text{classify}(\mathbf{x}, \mathbf{w}) = \underset{y}{\operatorname{argmax}} \text{ score}(\mathbf{x}, y, \mathbf{w})$$

- We can use any of our loss functions from before, as long as we can compute (sub)gradients
- Algorithm for doing this efficiently: **backpropagation**
- Basically just the chain rule of derivatives

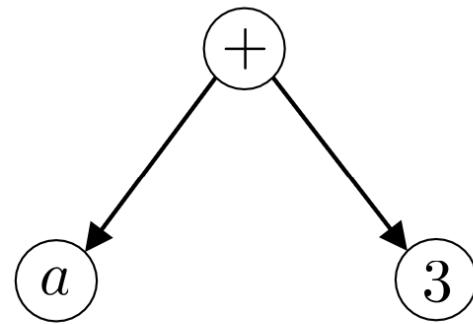
# Computation Graphs

- ❑ A useful way to represent the computations performed by a neural model (or any model!)
- ❑ Why useful? makes it easy to implement automatic differentiation (backpropagation)
- ❑ Many neural net toolkits let you define your model in terms of computation graphs (PyTorch, TensorFlow, DyNet, Theano, etc.)

# Backpropagation

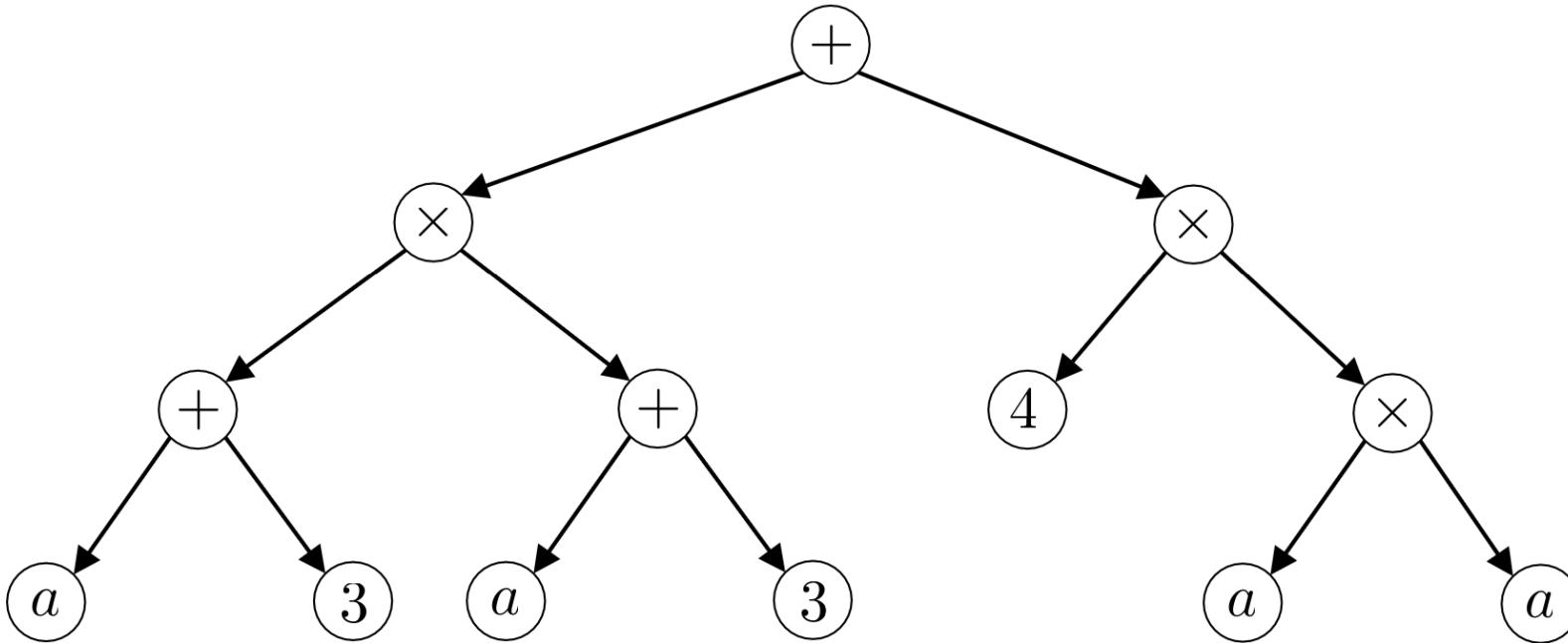
- ❑ Backpropagation has become associated with neural networks, but it's much more general
- ❑ Can be also used to compute gradients in **linear** models for structured prediction

# A simple computation graph:



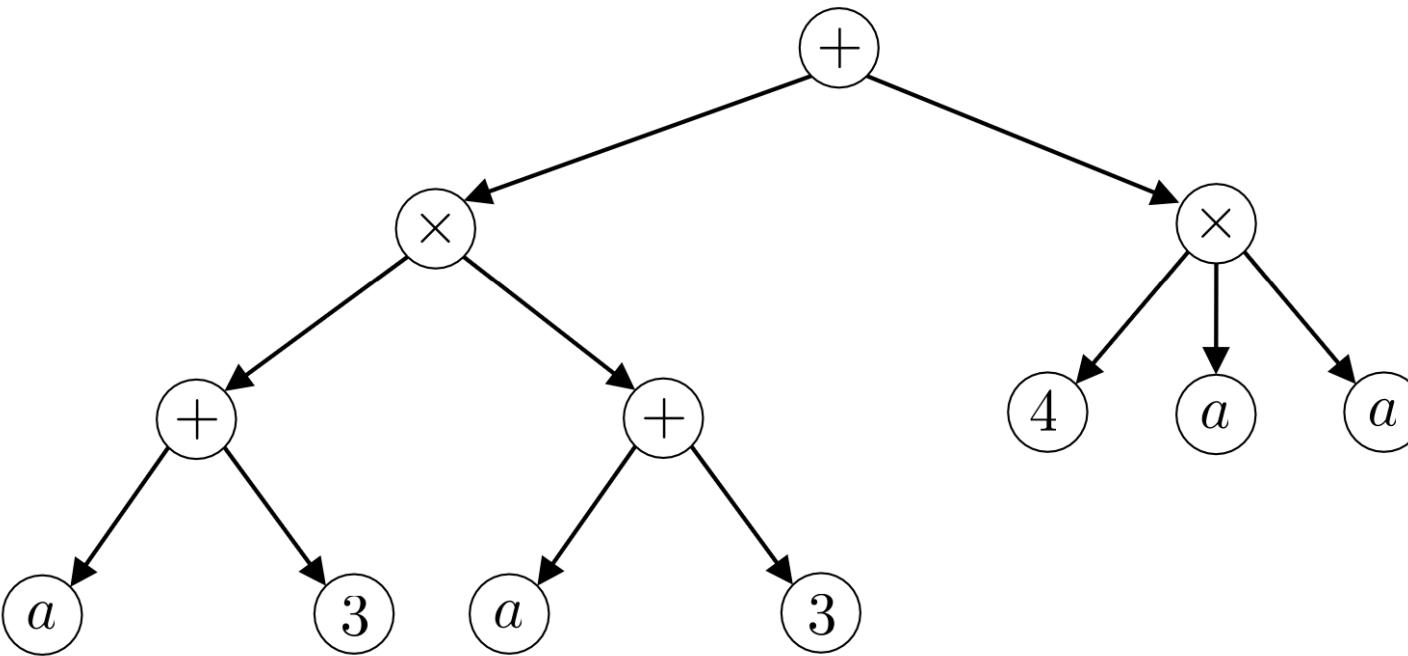
- ❑ Represents expression “ $a + 3$ ”

# A slightly bigger computation graph:

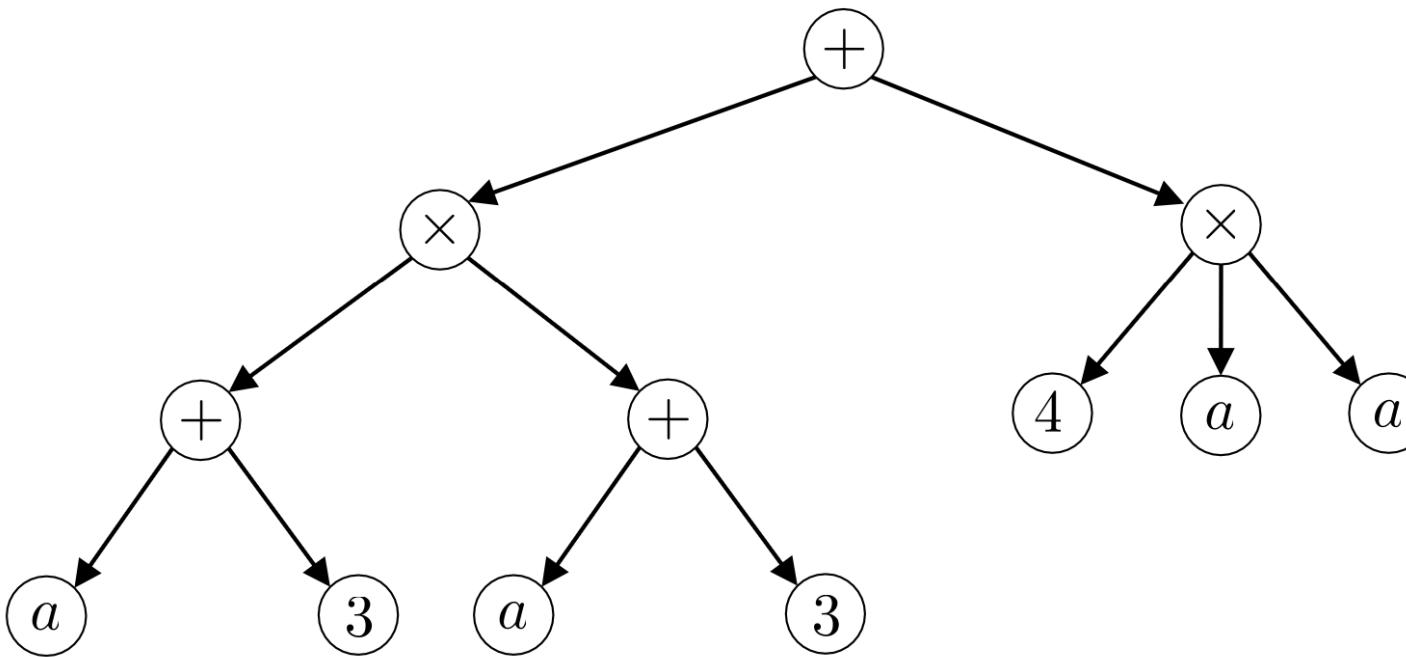


❑ Represents expression “ $(a + 3)^2 + 4a^2$ ”

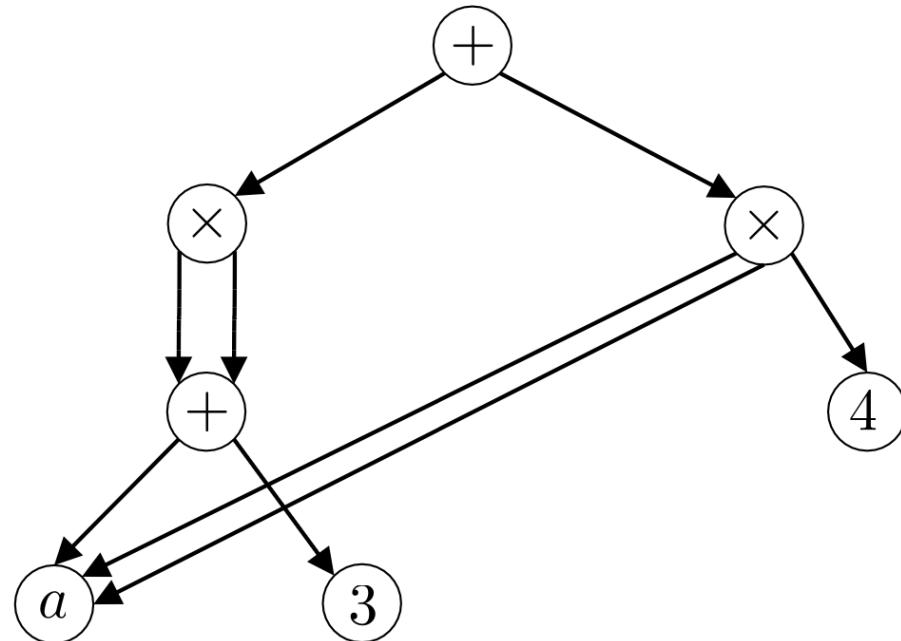
# Operators can have more than 2 operands:



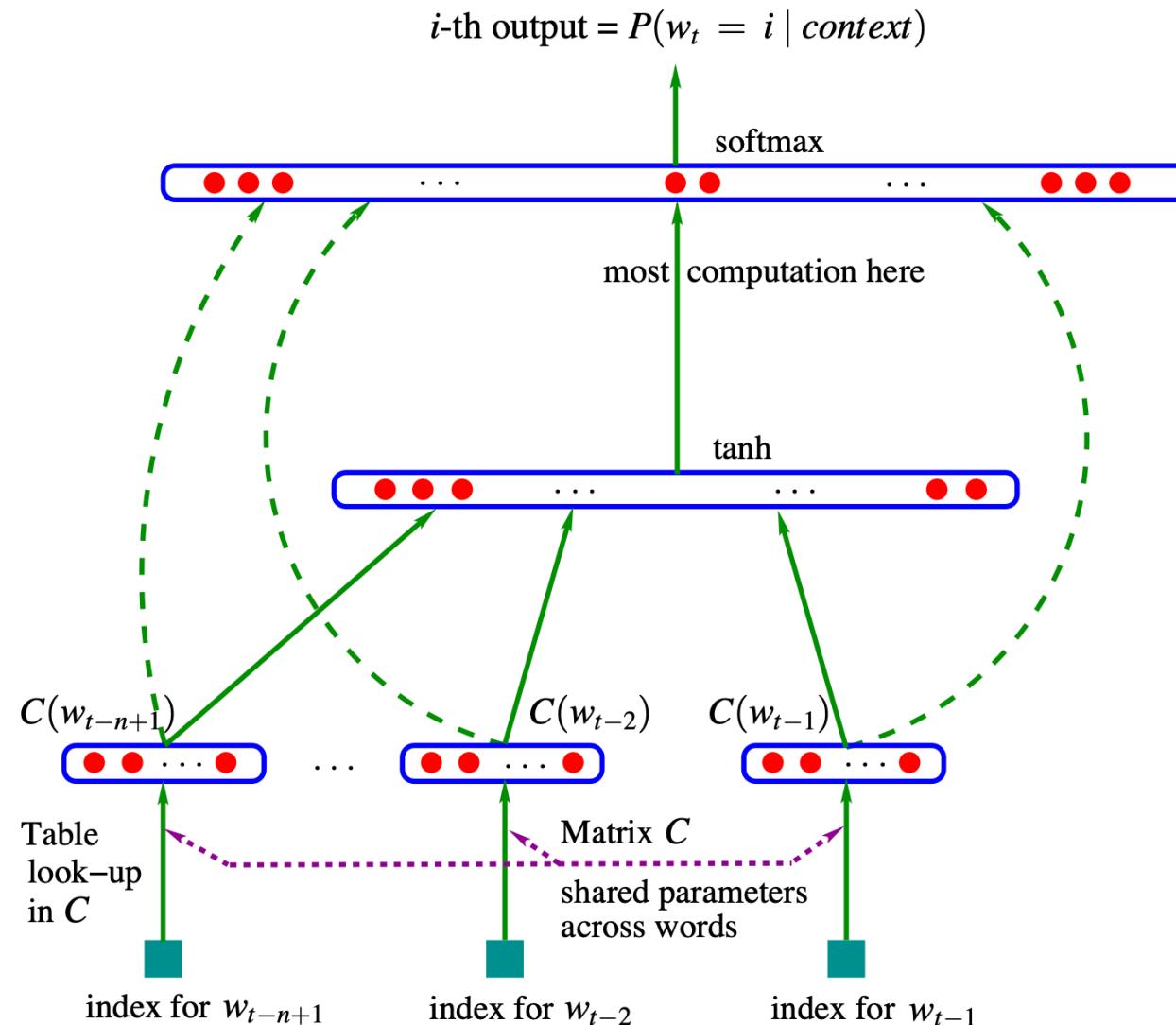
- Still represents expression “ $(a + 3)^2 + 4a^2$ ”



□ More concise:

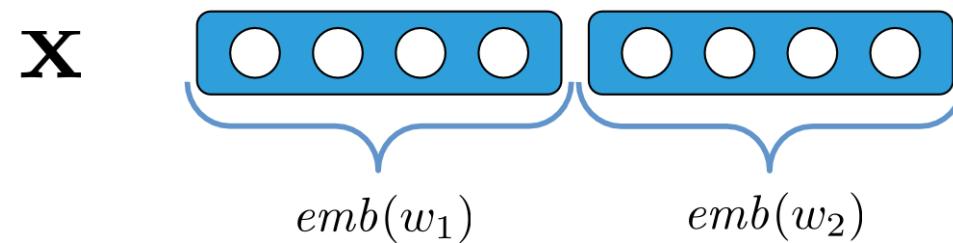


# Model (Bengio et al., 2003)



# A Simple Neural Trigram Language Model

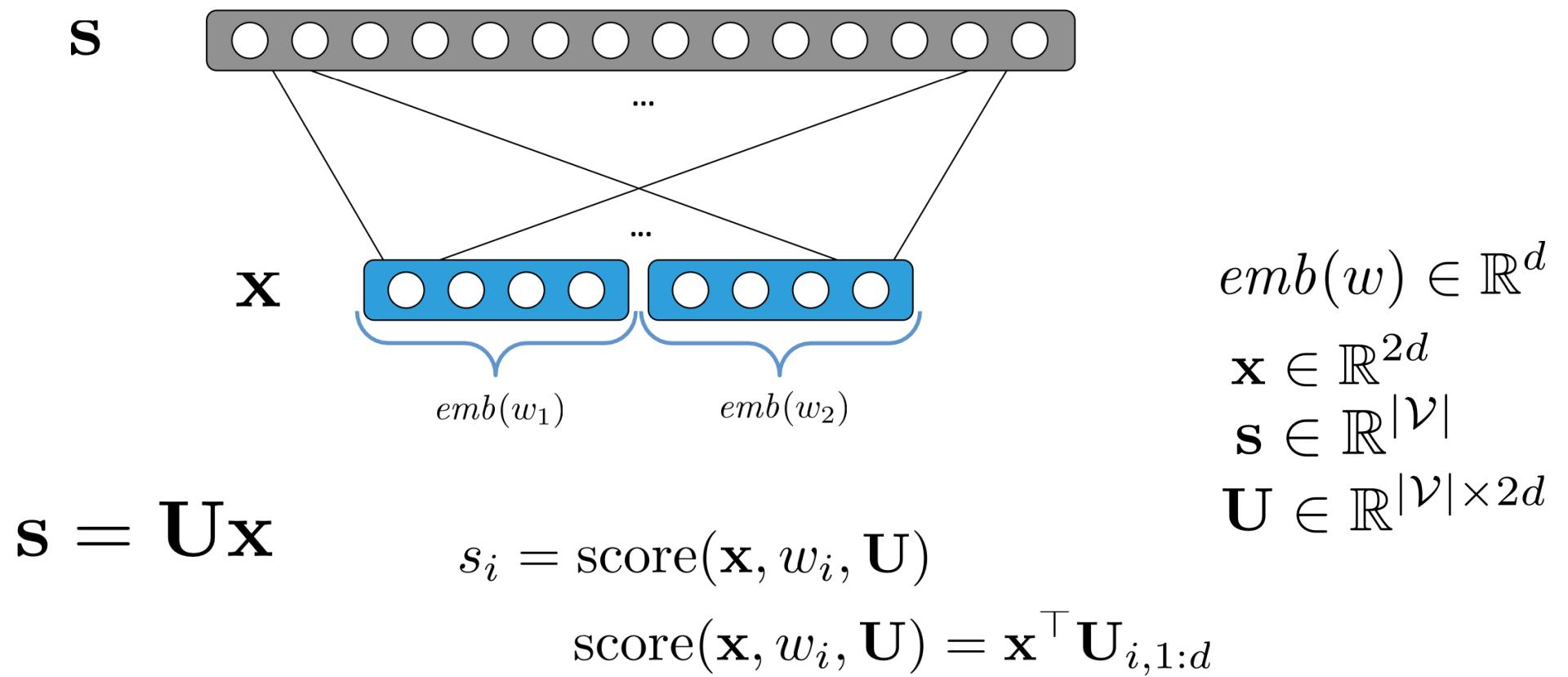
- ❑ Given previous words  $w_1$  and  $w_2$ , predict next word
- ❑ Input is concatenation of vectors (embeddings) of previous words:



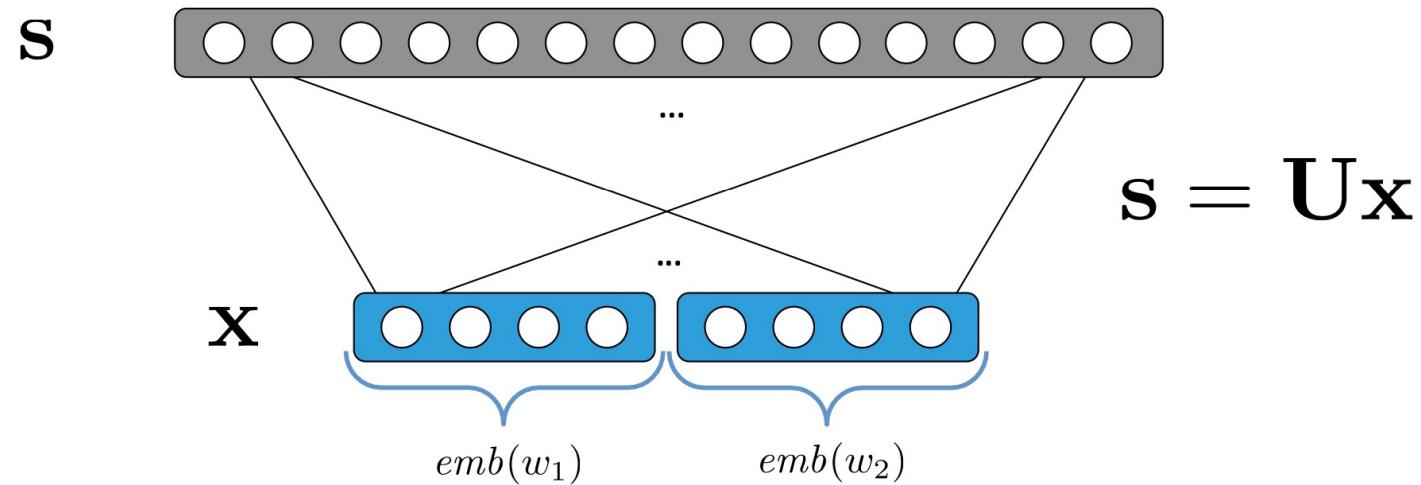
$$\mathbf{x} = \text{cat}(\text{emb}(w_1), \text{emb}(w_2))$$

# A Simple Neural Trigram Language Model

- Output vector contains scores of possible next words:

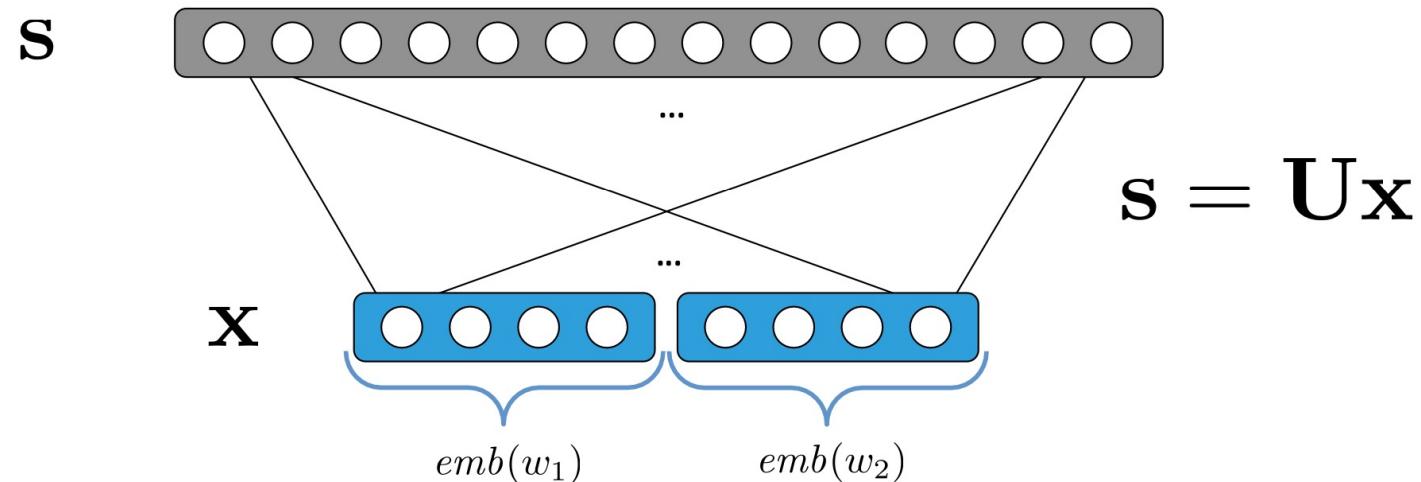


# A Simple Neural Trigram Language Model



❑ How many parameters are in this model?  $|\mathcal{V}| \times 3d$

# A Simple Neural Trigram Language Model



❑ Most common way: log loss

$$\text{loss}_{\log}(\langle w_1, w_2 \rangle, w_3, \boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(w_3 | \langle w_1, w_2 \rangle)$$

$$p_{\boldsymbol{\theta}}(w_3 | \langle w_1, w_2 \rangle) \propto \exp\{\text{score}(\text{cat}(emb(w_1), emb(w_2)), w_3, \mathbf{U})\}$$