

TTIC 31210: Advanced Natural Language Processing

Kevin Gimpel
Spring 2017

Lecture 2: Elements of Neural NLP

- Please email me with the following:
 - name
 - email address
 - whether you are taking the course for credit
- I will use the email addresses for the course mailing list

Assignment 1

- Assignment 1 has been posted; due in one week

Roadmap

- review of TTIC 31190 (week 1)
- **deep learning for NLP (weeks 2-4)**
- generative models & Bayesian inference (week 5)
- Bayesian nonparametrics in NLP (week 6)
- EM for unsupervised NLP (week 7)
- syntax/semantics and structure prediction (weeks 8-9)
- applications (week 10)

What is a classifier?

- a function from inputs x to classification labels y
- one simple type of classifier:
 - for any input x , assign a score to each label y , parameterized by vector θ :

$$\text{score}(x, y, \theta)$$

- classify by choosing highest-scoring label:

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

Modeling, Inference, Learning

inference: solve argmax

modeling: define score function

$$\operatorname{classify}(x, \boldsymbol{\theta}) = \operatorname{argmax}_y \operatorname{score}(x, y, \boldsymbol{\theta})$$

learning: choose $\boldsymbol{\theta}$

Notation

- We'll use boldface for vectors:

θ

- Individual entries will use subscripts and no boldface, e.g., for entry i :

θ_i

What is a neural network?

- just think of a neural network as a function
- it has inputs and outputs
- the term “neural” typically means a particular type of functional building block (“neural layers”), but the term has expanded to mean many things
- neural modeling is now better thought of as a modeling strategy (leveraging “distributed representations” or “representation learning”), or a family of related methods

Classifier Framework

$$\text{classify}(\mathbf{x}, \boldsymbol{\theta}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{L}} \text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})$$

- linear model score function:

$$\text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \sum_i \theta_i f_i(\mathbf{x}, \mathbf{y})$$

- we can also use a neural network for the score function!

neural layer = affine transform + nonlinearity

$$\mathbf{z}^{(1)} = g \left(\underbrace{W^{(0)}\mathbf{x} + \mathbf{b}^{(0)}}_{\text{affine transform}} \right)$$

nonlinearity

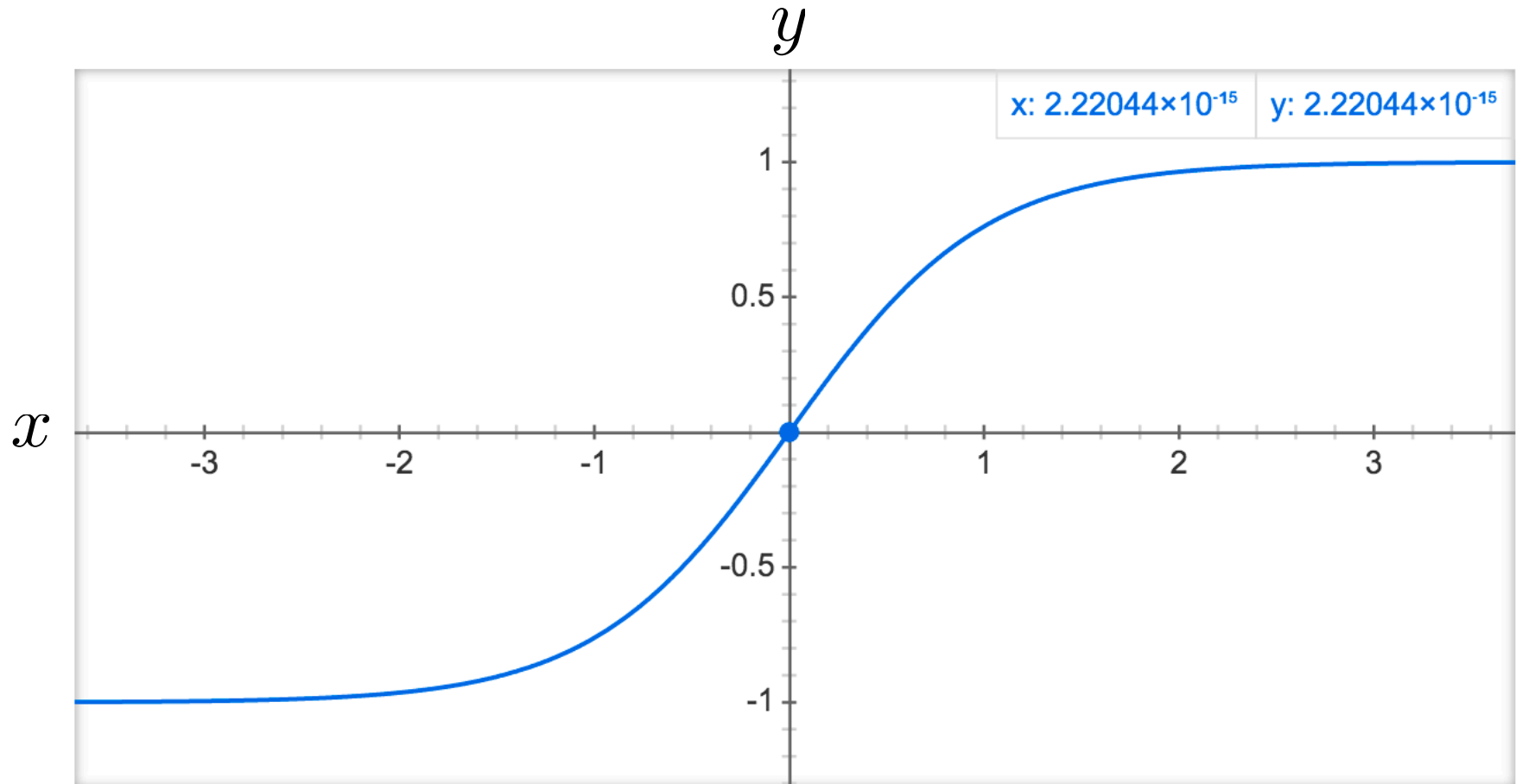
- this is a single “layer” of a neural network
- input vector is \mathbf{x}
- vector of “hidden units” is $\mathbf{z}^{(1)}$

Nonlinearities

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

- most common: elementwise application of g function to each entry in vector
- examples...

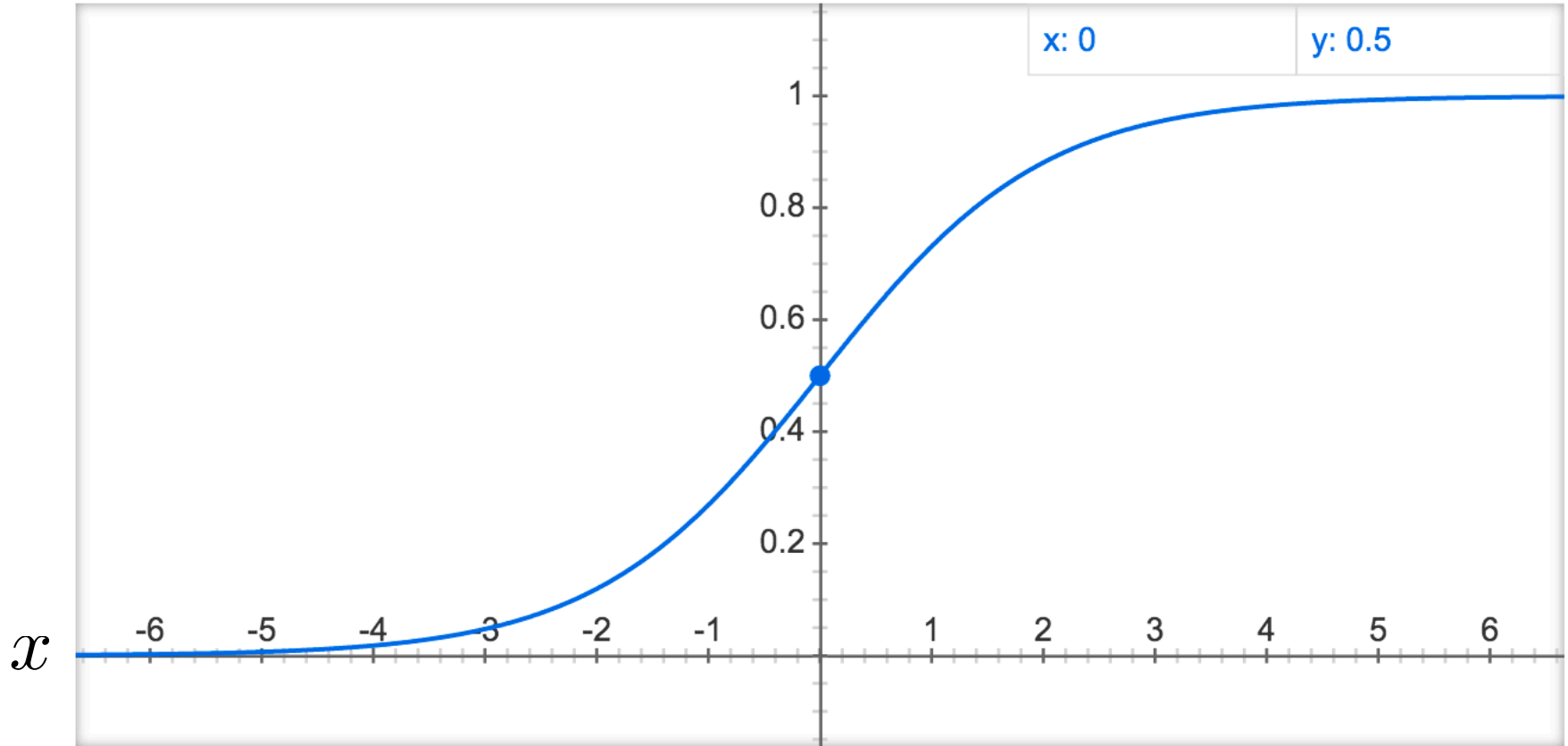
tanh: $y = \tanh(x)$



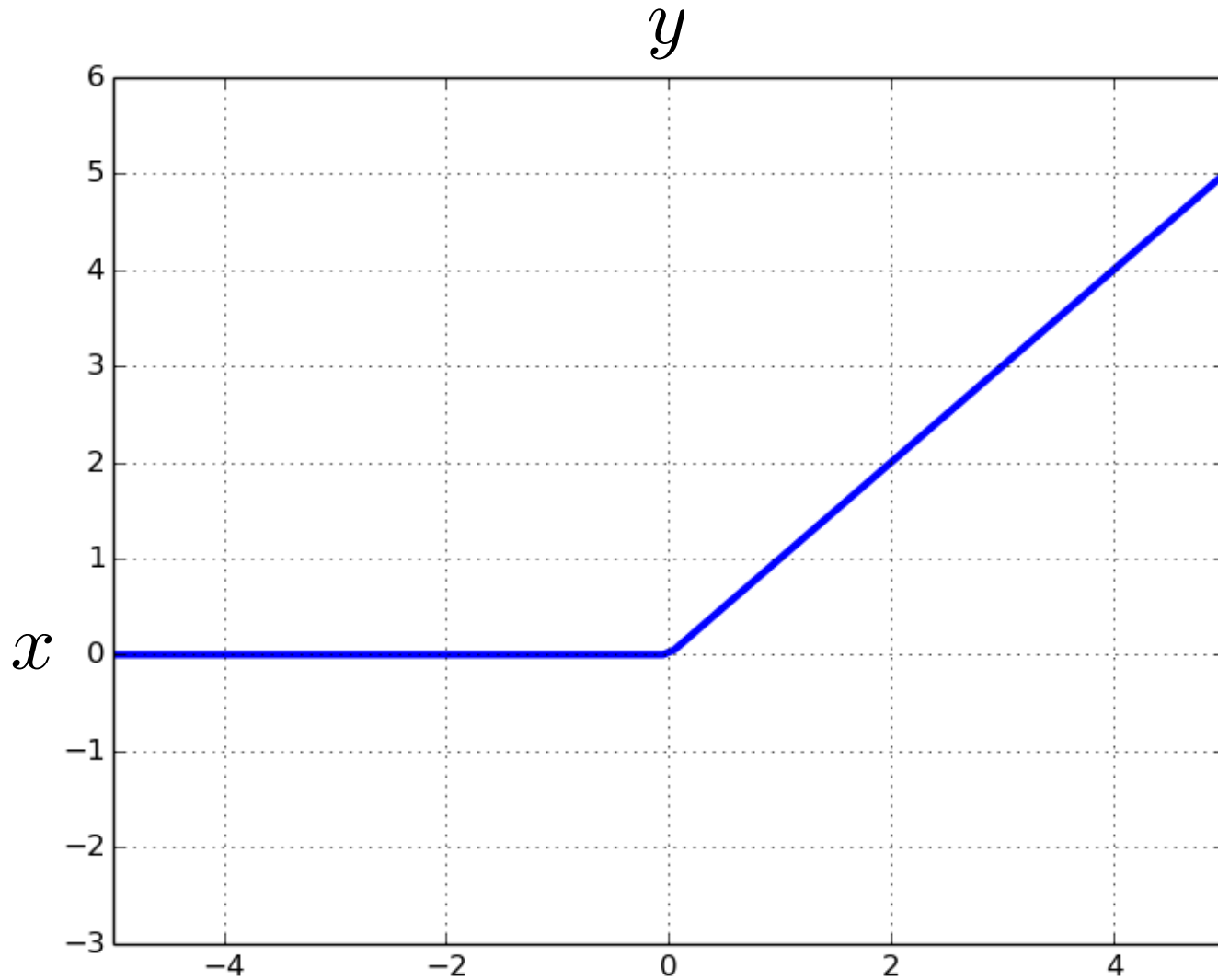
(logistic) sigmoid:

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

y



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



2-layer network

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

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

- this is a 2-layer neural network
- input vector is \mathbf{x}
- output vector is \mathbf{s}

2-layer neural network for sentiment classification

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

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



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

Use softmax function to convert scores into probabilities

$$\text{softmax}(\mathbf{s}) = \begin{bmatrix} \frac{\exp\{s_1\}}{\sum_i \exp\{s_i\}} \\ \dots \\ \frac{\exp\{s_d\}}{\sum_i \exp\{s_i\}} \end{bmatrix}$$

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

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

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

Why nonlinearities?

2-layer network:

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

written in a single equation:

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

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

Understanding the score function

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

entry 2 of bias vector

$$\text{score}(\mathbf{x}, \text{positive}, \boldsymbol{\theta}) = s_1 = g \left(W_{1,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_1^{(1)} \right)$$

$$\text{score}(\mathbf{x}, \text{negative}, \boldsymbol{\theta}) = s_2 = g \left(W_{2,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_2^{(1)} \right)$$

row vector corresponding to row 2 of $W^{(1)}$

Parameter sharing

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

$$\begin{aligned} \text{score}(\mathbf{x}, \text{positive}, \boldsymbol{\theta}) &= s_1 = g \left(W_{1,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_1^{(1)} \right) \\ \text{score}(\mathbf{x}, \text{negative}, \boldsymbol{\theta}) &= s_2 = g \left(W_{2,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_2^{(1)} \right) \end{aligned}$$

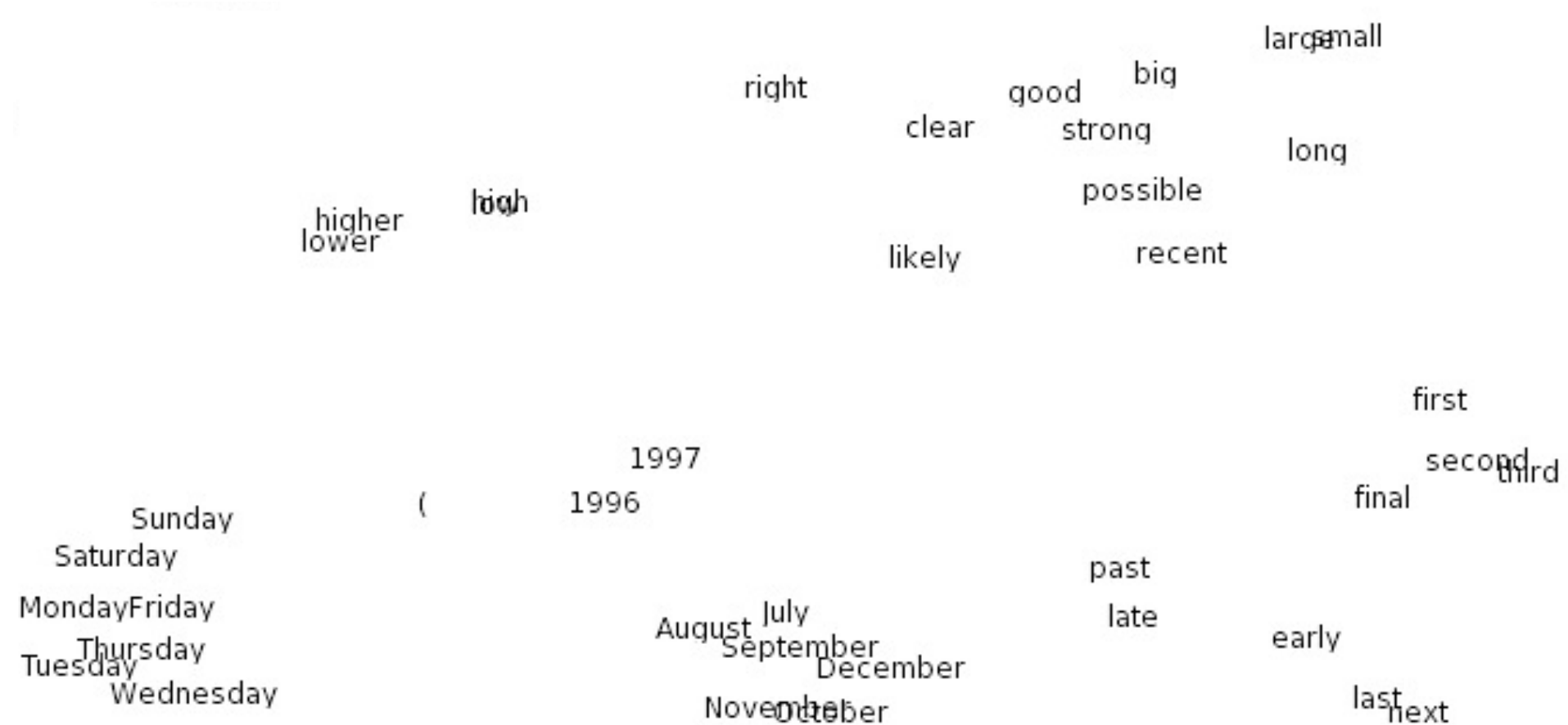
**parameters NOT
shared between labels**

**parameters
shared between
labels**

Word Embeddings

- “Basic unit” of neural NLP
- We’ll talk about ways to learn word embeddings next week
- Today: we’ll assume we have word embeddings as a black box

Word Embeddings



Turian et al. (2010)

Two Ways to Represent Word Embeddings

- \mathcal{V} = vocabulary , $|\mathcal{V}|$ = size of vocab
- 1: create $|\mathcal{V}|$ -dimensional “one-hot” vector for each word, multiply by word embedding matrix:

$$emb(x) = W \text{onehot}(\mathcal{V}, x)$$

Two Ways to Represent Word Embeddings

- \mathcal{V} = vocabulary , $|\mathcal{V}|$ = size of vocab
- 1: create $|\mathcal{V}|$ -dimensional “one-hot” vector for each word, multiply by word embedding matrix:

$$emb(x) = W \text{onehot}(\mathcal{V}, x)$$

- 2: store embeddings in a hash/dictionary data structure, do lookup to find embedding for word:

$$emb(x) = \text{lookup}(W, x)$$

Two Ways to Represent Word Embeddings

- \mathcal{V} = vocabulary , $|\mathcal{V}|$ = size of vocab
- 1: create $|\mathcal{V}|$ -dimensional “one-hot” vector for each word, multiply by word embedding matrix:

$$emb(x) = W \text{onehot}(\mathcal{V}, x)$$

- 2: store embeddings in a hash/dictionary data structure, do lookup to find embedding for word:

$$emb(x) = \text{lookup}(W, x)$$

- These are equivalent, second can be much faster (though first can be fast if using sparse operations)

Encoders

- encoder: a function to represent a word sequence as a vector
- simplest: average word embeddings:

$$f_{avg}(\mathbf{x}) = \frac{1}{|\mathbf{x}|} \sum_{i=1}^{|\mathbf{x}|} emb(x_i)$$

- many other functions possible!

Attention

- attention is a useful generic tool
- often used to replace a sum with an attention-weighted sum

Attention

- attention is a useful generic tool
- often used to replace a sum with an attention-weighted sum
- e.g., for a word summing encoder:

$$f_{sum}^{att}(\mathbf{x}) = \sum_{i=1}^{|\mathbf{x}|} att(x_i, i, \mathbf{x}) emb(x_i)$$

$$\sum_{i=1}^{|\mathbf{x}|} att(x_i, i, \mathbf{x}) = 1$$

- many other functions possible!

Ling et al. (EMNLP 2015)

- Not All Contexts Are Created Equal: Better Word Representations with Variable Attention

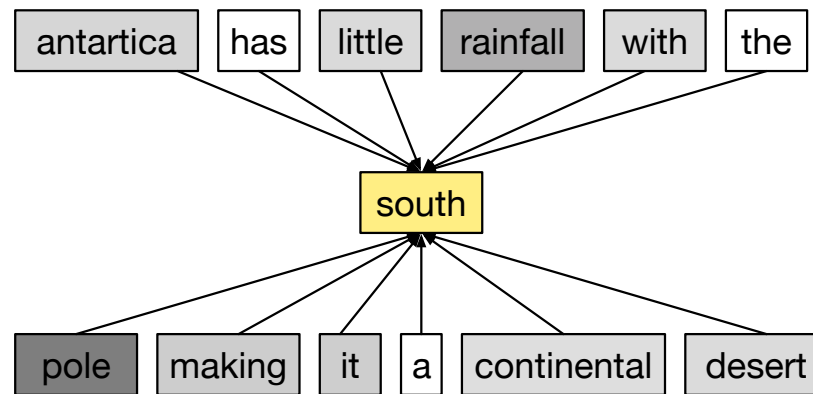
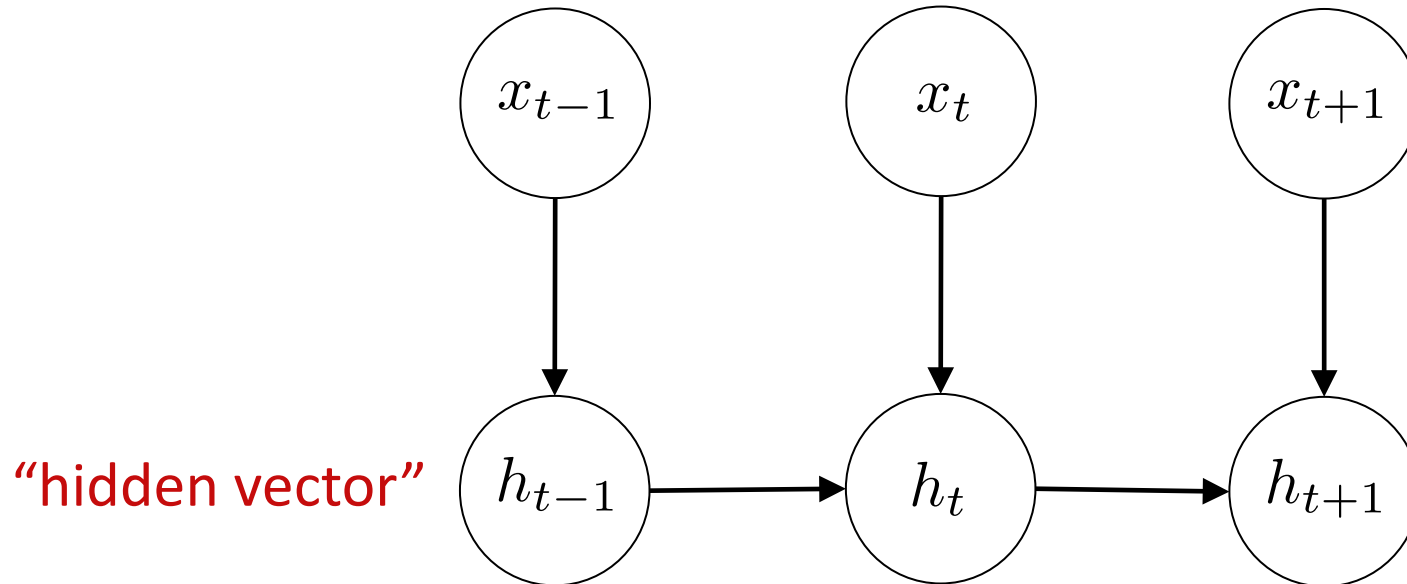


Figure 1: Illustration of the inferred attention parameters for a sentence from our training data when predicting the word *south*; darker cells indicate higher weights.

Recurrent Neural Networks

$$h_t = \tanh \left(W^{(xh)} x_t + W^{(hh)} h_{t-1} + b^{(h)} \right)$$



Neural Similarity Learning

- A common need:
compute similarity/affinity of two things
 - maybe two things of the same type,
 - two things with different types being mapped to same space, or
 - two things with different types being mapped to different spaces, but being compared with a learned similarity function
- Examples?

Synonym Pairs

- Faruqui et al. (NAACL 2014), Wieting et al. (TACL 2014), *inter alia*

contamination	pollution
converged	convergence
captioned	subtitled
outwit	thwart
bad	villain
broad	general
permanent	permanently
bed	sack
carefree	reckless
absolutely	urgently
...	...

Translation Pairs

- Haghighi et al. (ACL 2008), Mikolov et al. (2013), Faruqui and Dyer (EACL 2014)

dog	hund
man	mann
woman	frau
city	stadt
person	man
the	der
the	die
the	das
...	...

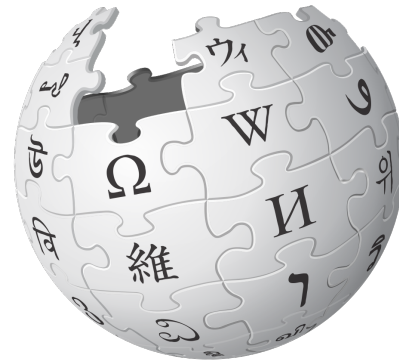
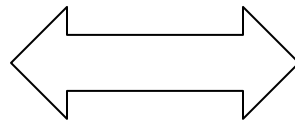
Sentence Pairs

this was also true for pompeii , where the temple of jupiter that was already there was enlarged and made more roman when the romans took over .

this held true for pompeii , where the previously existing temple of jupiter was enlarged and romanized upon conquest .



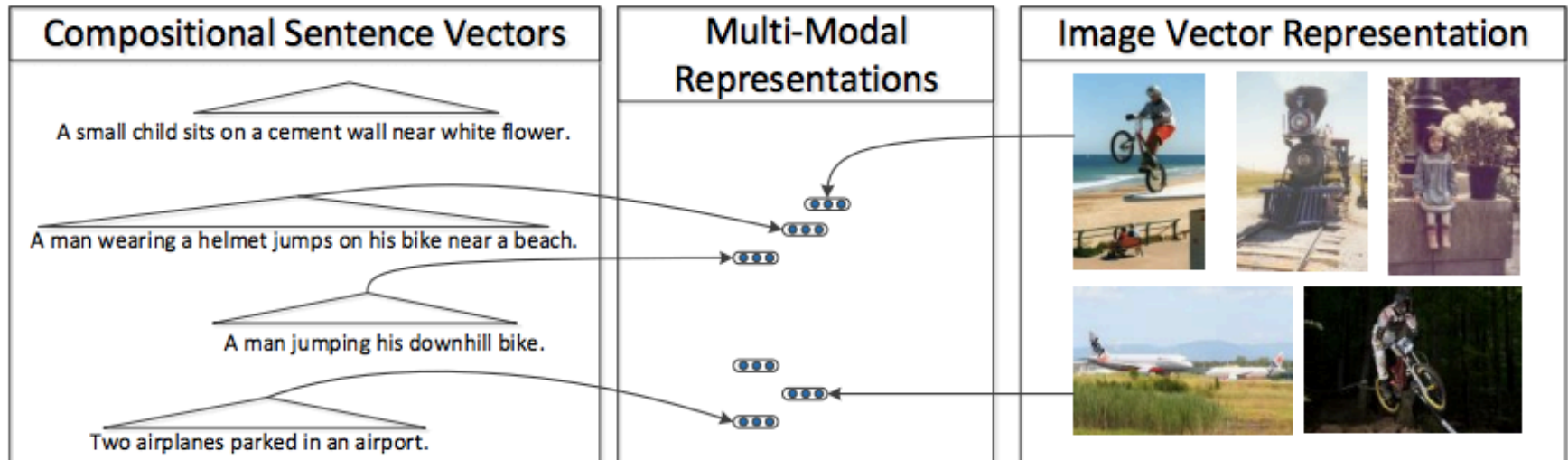
WIKIPEDIA
The Free Encyclopedia



WIKIPEDIA
Simple English

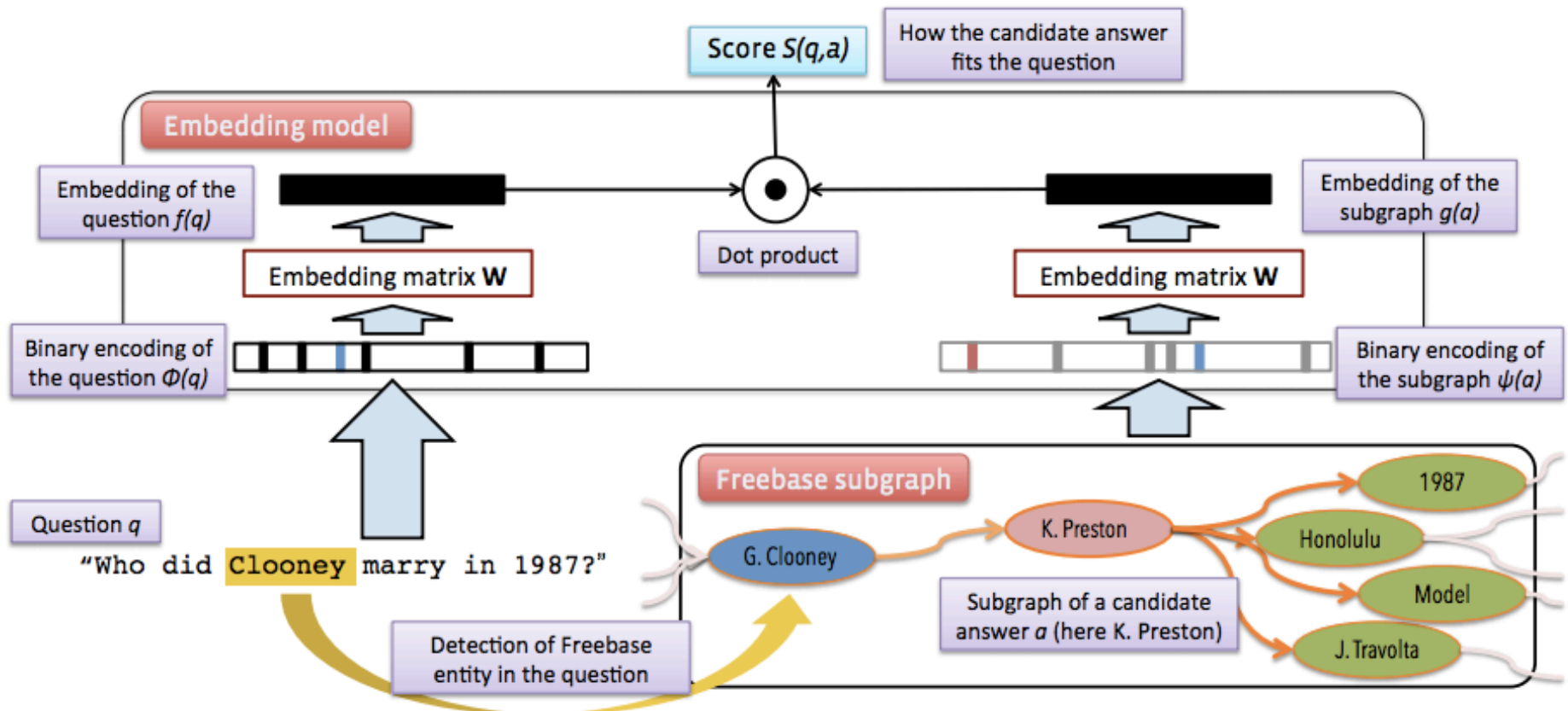
Captions and Images

- Richard Socher, Andrej Karpathy, Quoc V. Le, Christopher D. Manning, Andrew Y. Ng. "Grounded Compositional Semantics For Finding And Describing Images With Sentences," *TACL* 2014.



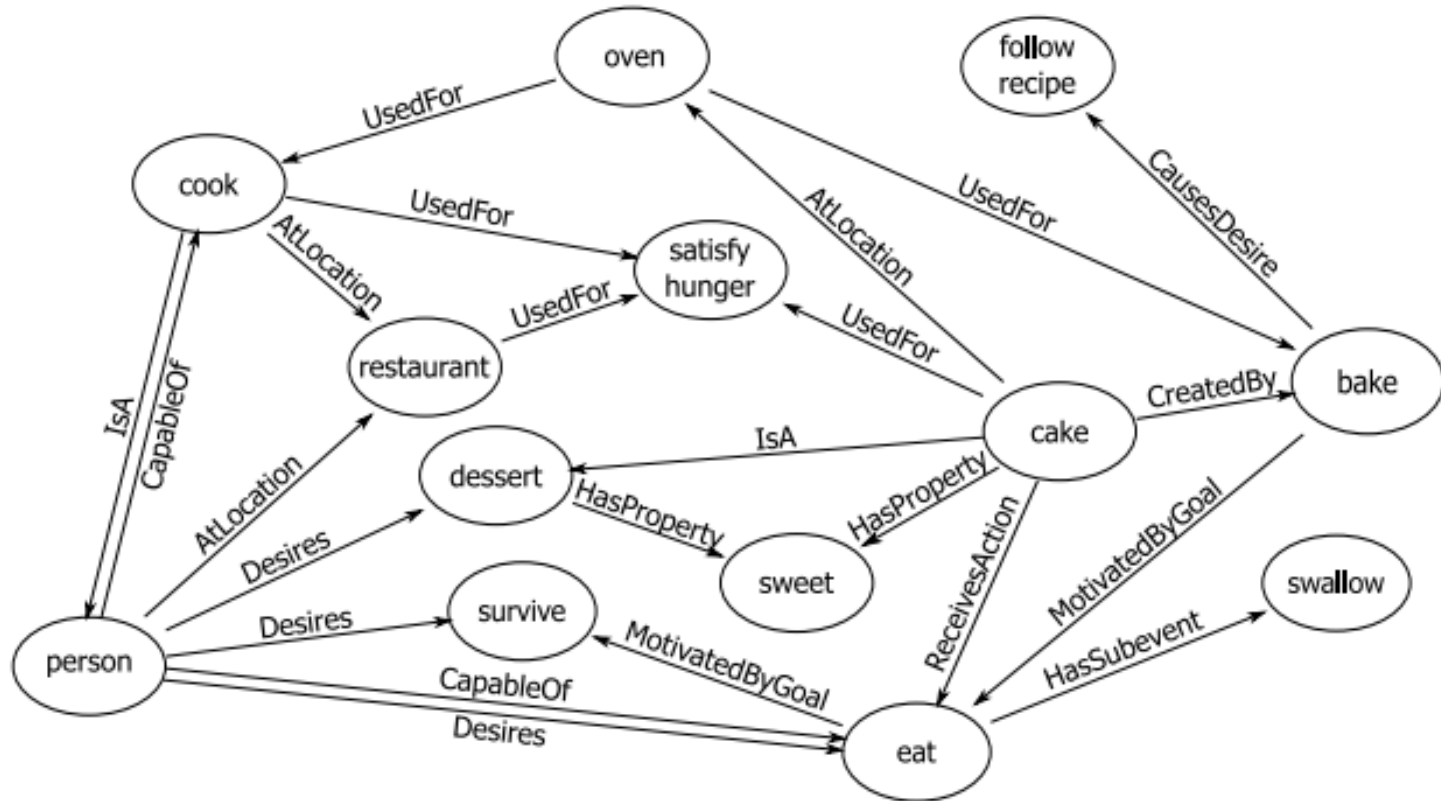
Questions and Answers

- Mohit Iyyer, Jordan Boyd-Graber, Leonardo Claudino, Richard Socher, and Hal Daumé III. “A Neural Network for Factoid Question Answering over Paragraphs,” *EMNLP* 2014
- Antoine Bordes, Sumit Chopra, and Jason Weston. “Question Answering with Subgraph Embeddings,” *EMNLP* 2014.



Commonsense Knowledge Triples

- Xiang Li, Aynaz Taheri, Lifu Tu, Kevin Gimpel. “Commonsense Knowledge Base Completion,” *ACL* 2016.



<“cake”, UsedFor, “satisfy hunger”>

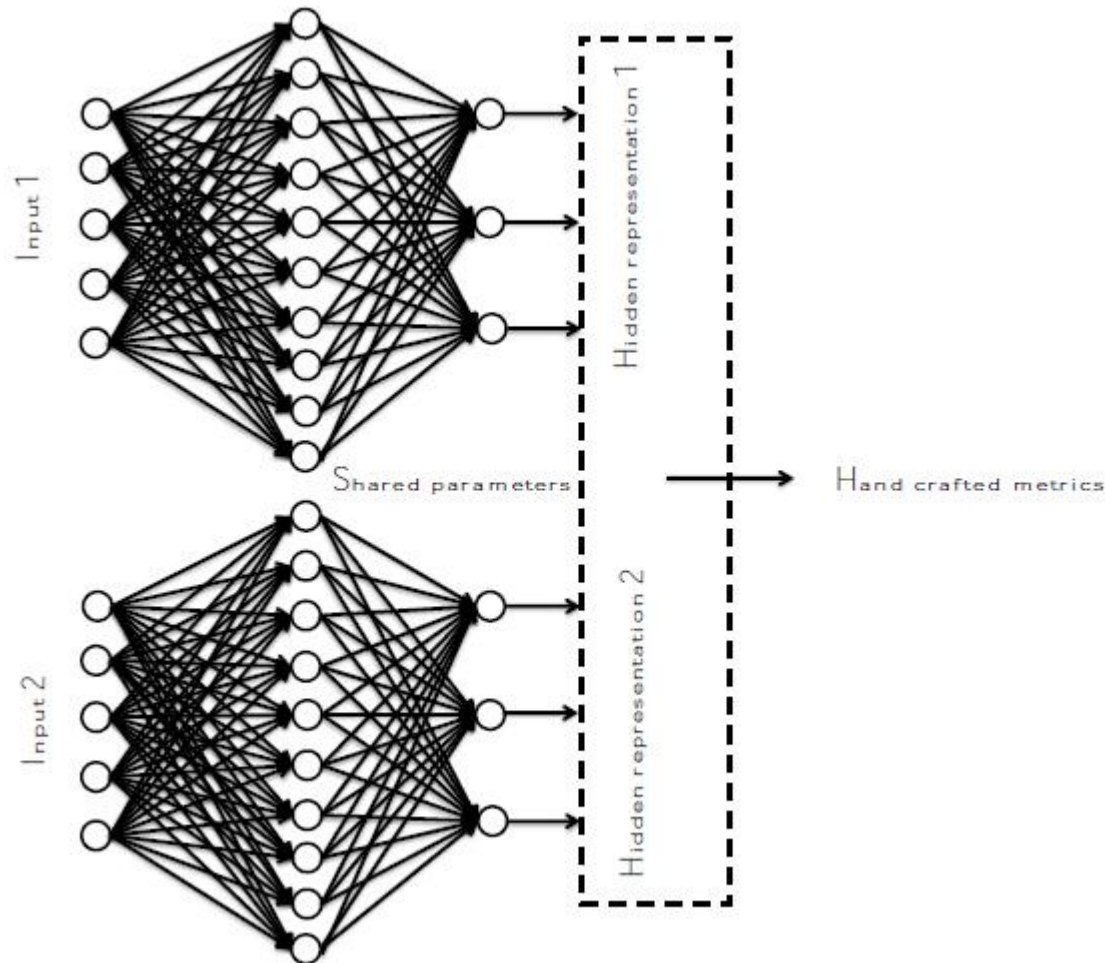
Stories and Endings

- story:
 - Jennifer has a big exam tomorrow. She got so stressed, she pulled an all-nighter. She went into class the next day, weary as can be. Her teacher stated that the test is postponed for next week.
- ending:
 - Jennifer felt bittersweet about it.
- from ROC Story Corpus (Mostafazadeh et al., 2016)

- Other examples/applications you can think of?
- Sometimes direction matters, sometimes not
- Sometimes there is a particular kind of relation being named for each pair, sometimes not (i.e., just one kind)

Neural Similarity Modeling

- “Siamese networks” (Bromley et al., 1993)
 - these typically share parameters, hence the name



Similarity Modeling

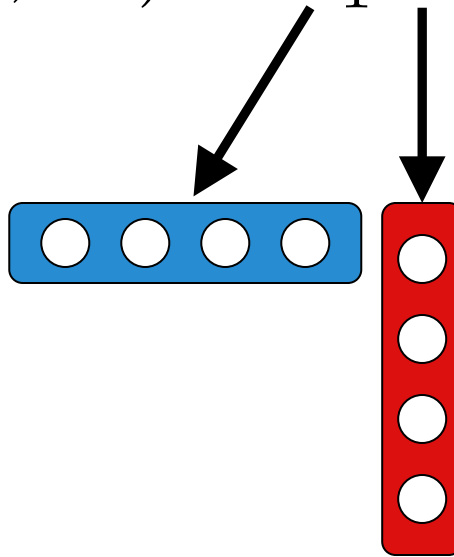
- Siamese networks typically share parameters across the two networks
- but it's also common to not share parameters when the items have different types, but we still want to relate them in some way
 - whether map to same space + compute sim or map each to some other space + compute sim

Similarity Functions

- many choices for similarity functions
- we'll go over some in the next few slides
- throughout, keep in mind:
 - output range
 - symmetric? $sim(\mathbf{x}_1, \mathbf{x}_2) = sim(\mathbf{x}_2, \mathbf{x}_1)$
 - introduces new parameters?
 - connections among similarity functions?
 - notes/tips on using these

Dot Product

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^\top \mathbf{x}_2$$



range? \mathbb{R}

symmetric or asymmetric? **symmetric**

introduces parameters? **no**

Cosine Similarity

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \cos(\mathbf{x}_1, \mathbf{x}_2) = \frac{\mathbf{x}_1^\top \mathbf{x}_2}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|}$$

$$\|\mathbf{x}\| = \sqrt{\sum_i x_i^2}$$

range? $[-1, 1]$

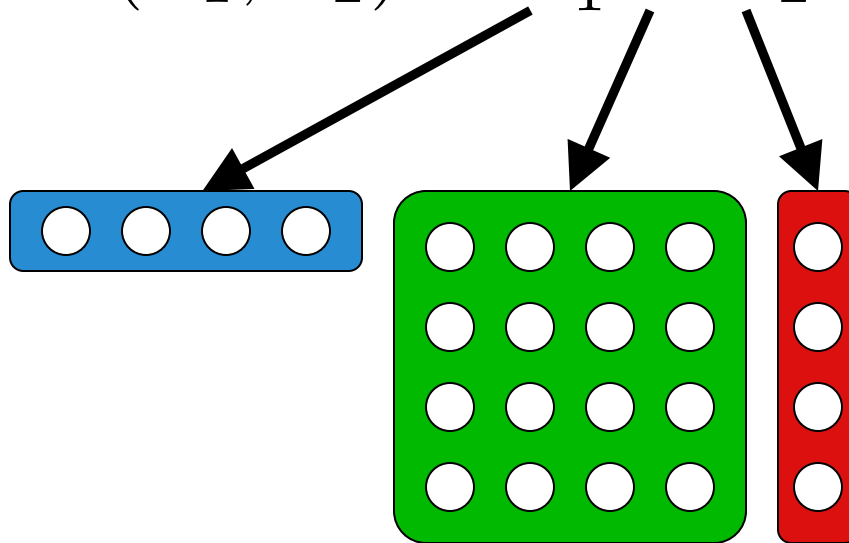
symmetric or asymmetric? **symmetric**

introduces parameters? **no**

generalizes anything? **dot product when vectors have norm 1**

Bilinear Function

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^\top W \mathbf{x}_2$$



range? \mathbb{R}

symmetric or asymmetric? **asymmetric in general**

introduces parameters? **yes**

generalizes anything? **dot product if W is identity**

Notes on Using Bilinear Functions

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^\top W \mathbf{x}_2$$

similarity can depend on relation by using different bilinear weight matrices for different relations:

$$\text{sim}(\mathbf{x}_1, r, \mathbf{x}_2) = \mathbf{x}_1^\top W_r \mathbf{x}_2$$

often W is initialized to the identity matrix (and regularized back to it)

potential issue: W might be very huge

L1/L2 Distances

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$$

range? $\mathbb{R}_{\geq 0}$

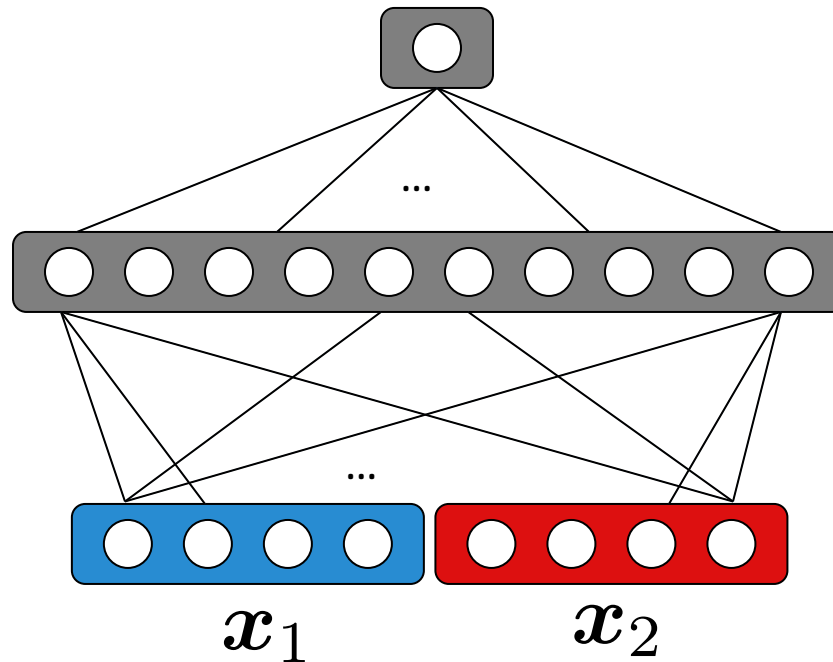
symmetric or asymmetric? **symmetric**

introduce parameters? **no**

Deep Neural Network

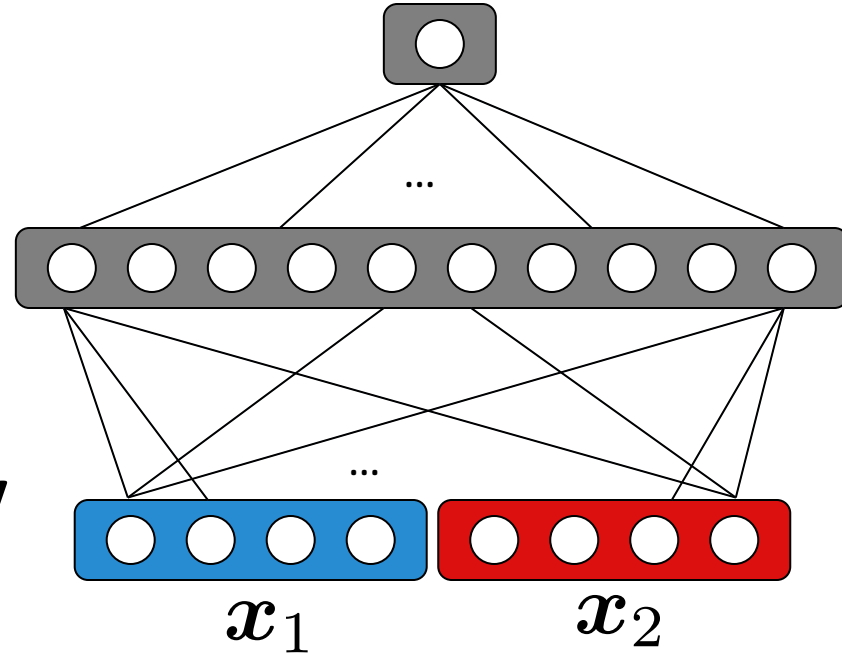
$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \text{DNN}(\text{cat}(\mathbf{x}_1, \mathbf{x}_2))$$

concatenate vectors, pass to DNN, use scalar for final output:



Deep Neural Network

$$\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = \text{DNN}(\text{cat}(\mathbf{x}_1, \mathbf{x}_2))$$



range? **depends on nonlinearity**

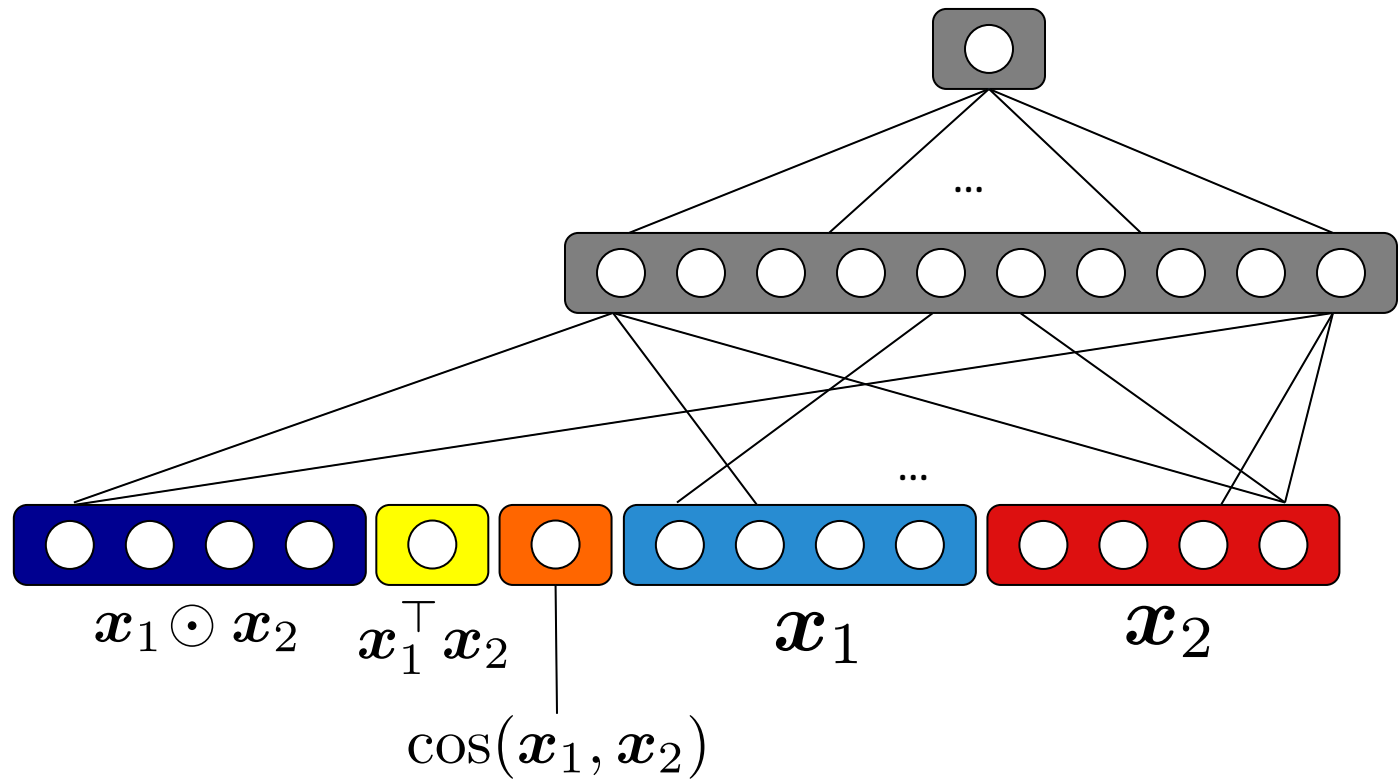
symmetric? **asymmetric**

introduces parameters? **yes**

generalizes anything? **yes, can represent any function!**

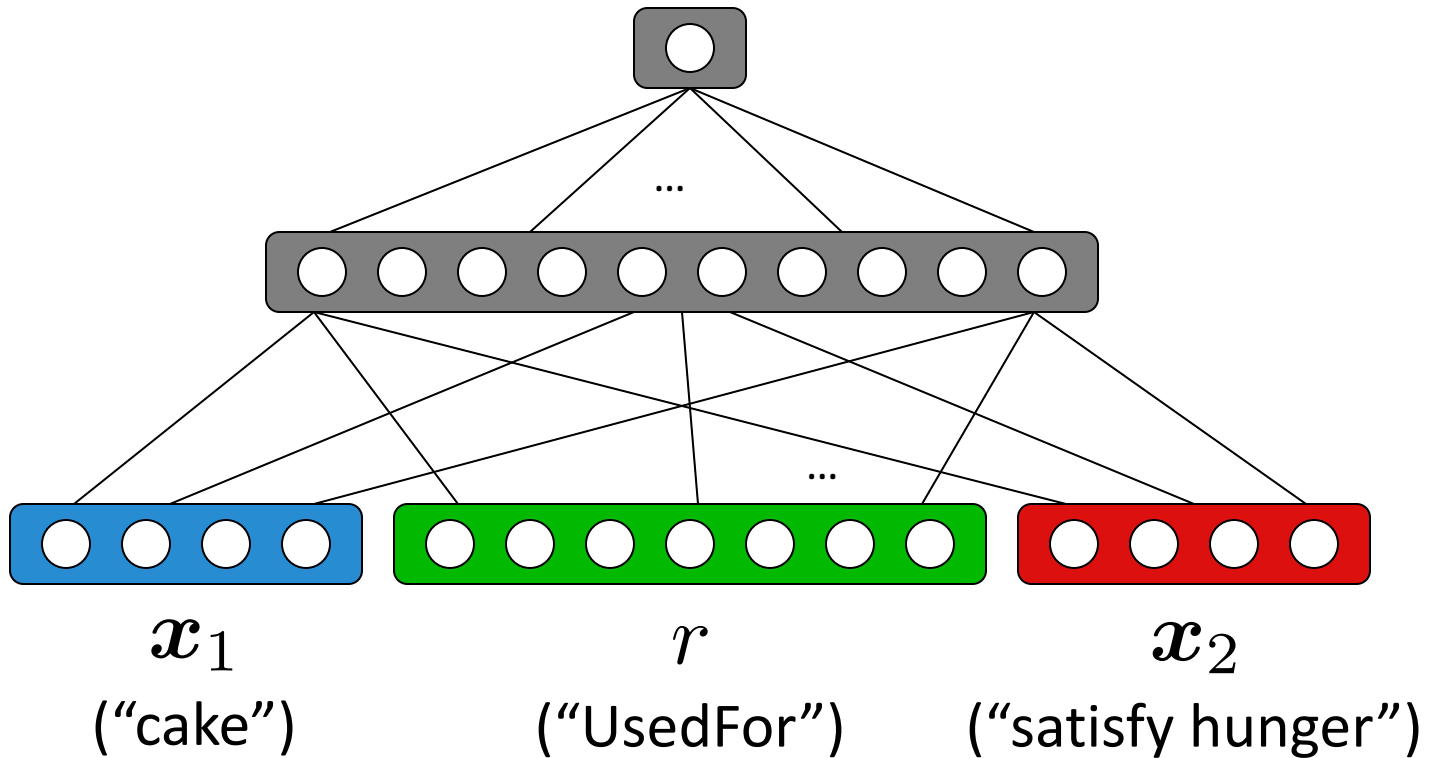
Notes on DNN Similarity Functions

since DNNs are so powerful, things could go horribly wrong.
in practice, often pass additional quantities:



Deep Neural Network

similarity can depend on relation:



Learning for Similarity

- We want to learn input representations as well as all parameters of $sim(\mathbf{x}_1, \mathbf{x}_2)$
- We'll just write all these parameters as θ
- How about this?

$$\min_{\theta} \sum_{\langle \mathbf{x}_1, \mathbf{x}_2 \rangle \in \text{train}} -sim(\mathbf{x}_1, \mathbf{x}_2)$$

- Any potential problems with this?

(Better) Learning for Similarity

- Contrastive hinge loss:

$$\min_{\theta} \sum_{\langle \mathbf{x}_1, \mathbf{x}_2 \rangle \in \text{train}} [-\text{sim}(\mathbf{x}_1, \mathbf{x}_2) + \text{sim}(\mathbf{x}_1, \mathbf{v})]_+$$

$$[a]_+ = \max(0, a)$$

- \mathbf{v} is a “negative” example
- Any potential problems with this?

Learning with Neural Networks

$$\text{classify}(\mathbf{x}, \boldsymbol{\theta}) = \underset{\mathbf{y} \in \mathcal{L}}{\text{argmax}} \text{score}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta})$$

$$\text{score}(\mathbf{x}, \text{positive}, \boldsymbol{\theta}) = s_1 = g \left(W_{1,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_1^{(1)} \right)$$

$$\text{score}(\mathbf{x}, \text{negative}, \boldsymbol{\theta}) = s_2 = g \left(W_{2,*}^{(1)} g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right) + \mathbf{b}_2^{(1)} \right)$$

- we can use any of our loss functions from before, as long as we can compute (sub)gradients
- algorithm for doing this efficiently:
backpropagation
- it's basically just the chain rule of derivatives

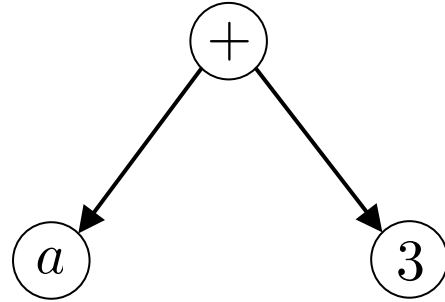
Computation Graphs

- 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 (Theano, (Py)Torch, TensorFlow, CNTK, DyNet, PENNE, etc.)

Backpropagation

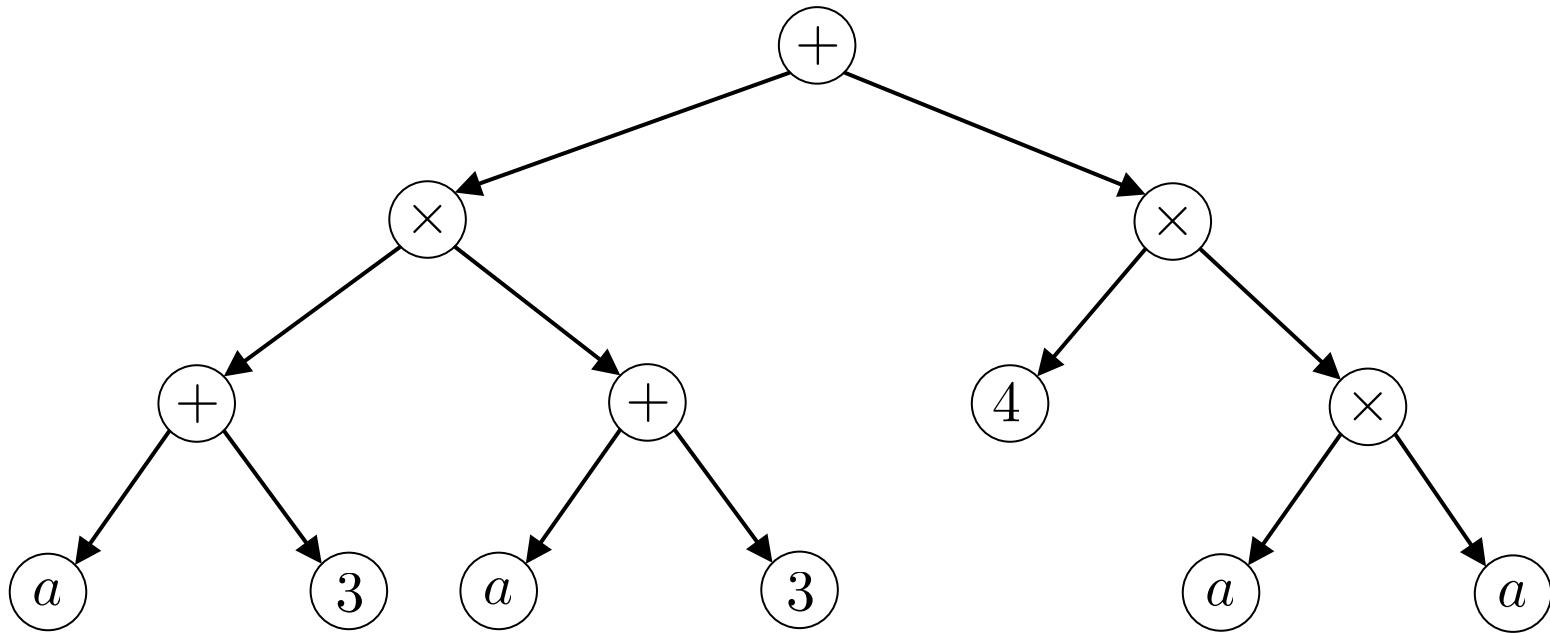
- backpropagation has become associated with neural networks, but it's much more general
- I also use backpropagation to compute gradients in **linear** models for structured prediction

A simple computation graph:



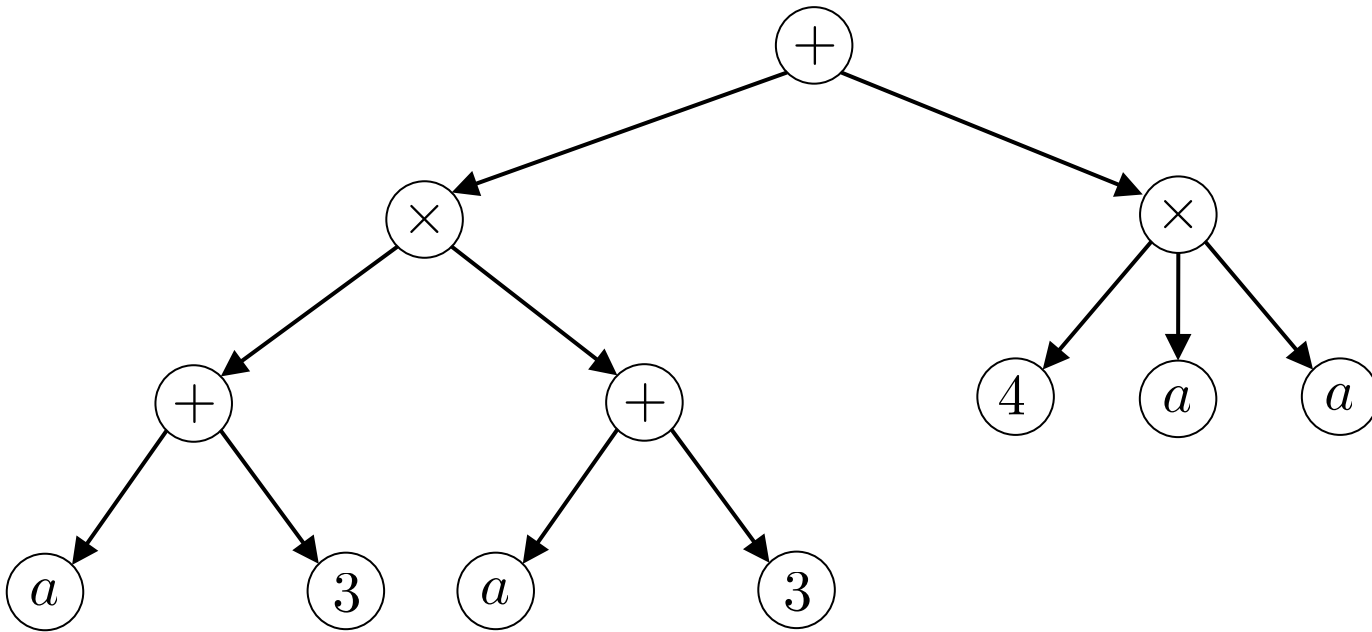
- represents expression “a + 3”

A slightly bigger computation graph:

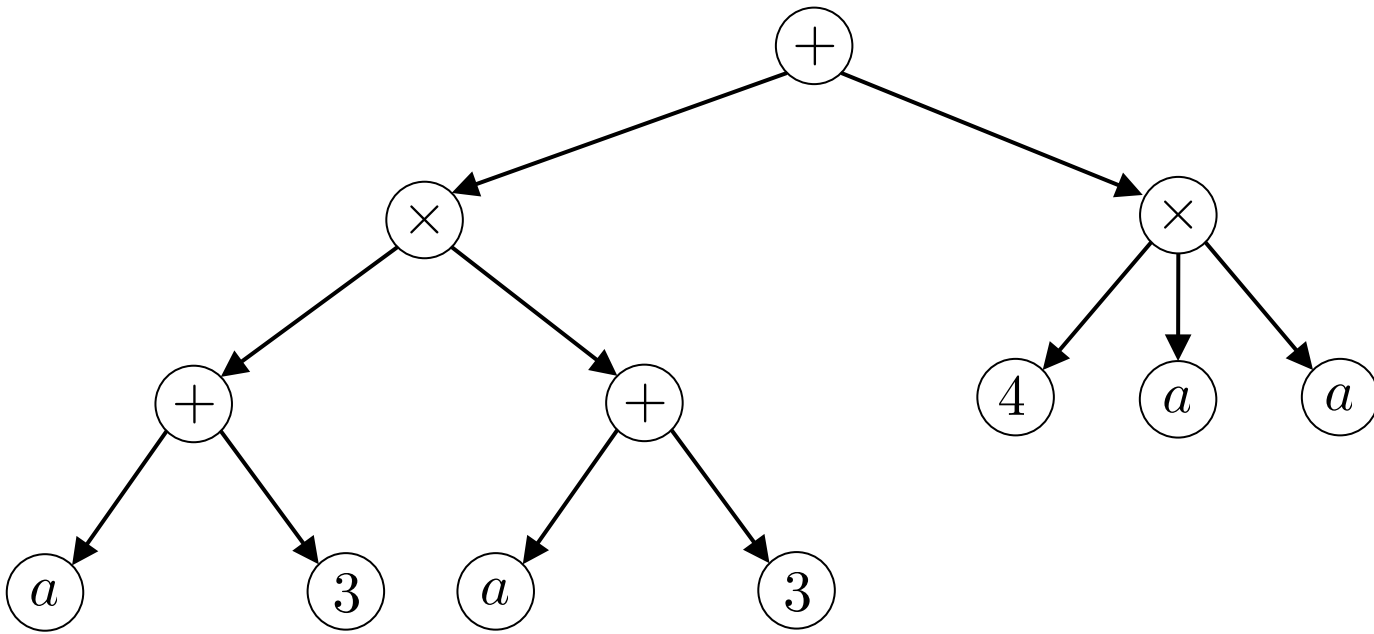


- represents expression $(a + 3)^2 + 4a^2$

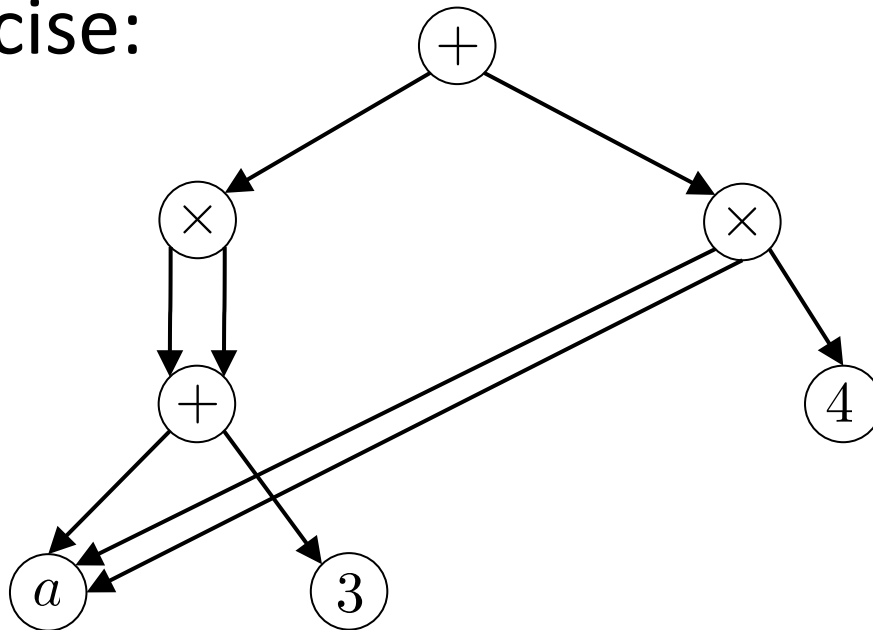
Operators can have more than 2 operands:



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



- more concise:



Overfitting & Regularization

- when we can fit any function, **overfitting** becomes a big concern
- overfitting: learning a model that does well on the training set but doesn't generalize to new data
- there are many strategies to reduce overfitting (we'll use the general term **regularization** for any such strategy)
- you used **early stopping** in Assignment 1, which is one kind of regularization

Regularization Terms

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

- most common: penalize large parameter values
- intuition: large parameters might be instances of overfitting
- examples:

L_2 regularization: $R_{L_2}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2 = \sum_i \theta_i^2$

(also called Tikhonov regularization
or ridge regression)

L_1 regularization: $R_{L_1}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_i |\theta_i|$

(also called basis pursuit or LASSO)

Regularization Terms

L_2 regularization: $R_{L_2}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2 = \sum_i \theta_i^2$

differentiable, widely-used

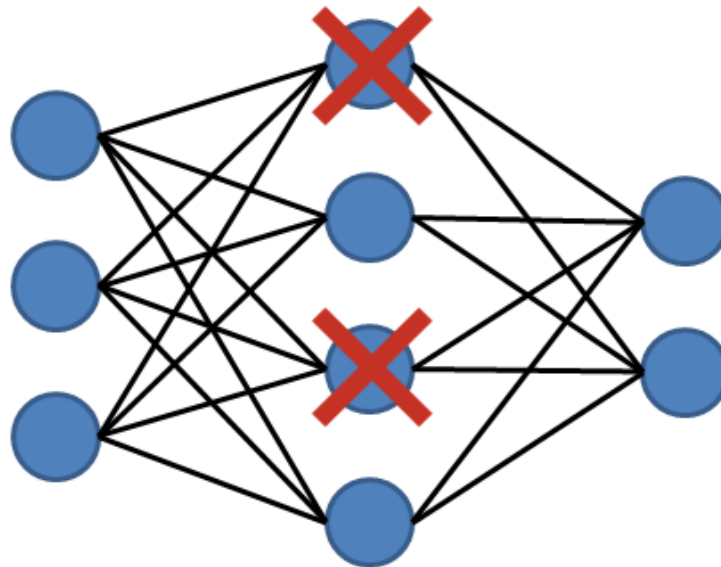
L_1 regularization: $R_{L_1}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_i |\theta_i|$

not differentiable (but is subdifferentiable)

leads to sparse solutions (many parameters become zero!)

Dropout

- popular regularization method for neural networks
- randomly “drop out” (set to zero) some of the vector entries in the layers



Optimization Algorithms

- many choices:
 - SGD
 - AdaGrad
 - AdaDelta
 - RMSProp
 - Adam
 - SGD with momentum
- we don't have time to go through these in class, but you should try using them! (most toolkits have implementations of these and others)

2-transformation (1-layer) network

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

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



vector of label scores

- we'll call this a “2-transformation” neural network, or a “1-layer” neural network
- input vector is \mathbf{x}
- score vector is \mathbf{s}
- one hidden vector $\mathbf{z}^{(1)}$ (“hidden layer”)

1-layer neural network for sentiment classification

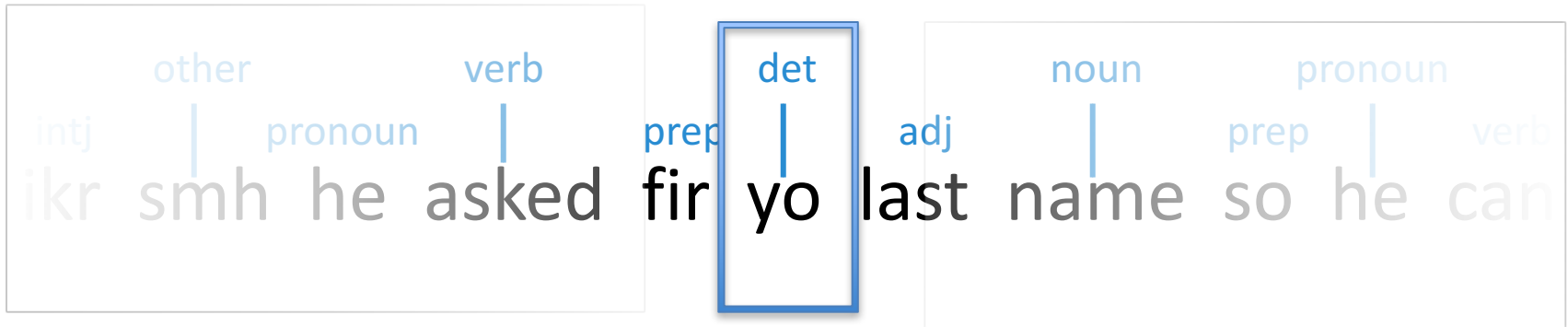
$$\mathbf{z}^{(1)} = g \left(W^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

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



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

Neural Networks for Twitter Part-of-Speech Tagging



- let's use the center word + two words to the right:

$$\mathbf{x} = [0.4 \quad \dots \quad 0.9 \quad 0.2 \quad \dots \quad 0.7 \quad 0.3 \quad \dots \quad 0.6]^\top$$

vector for *yo* vector for *last* vector for *name*

- if *name* is to the right of *yo*, then *yo* is probably a form of *your*
- but our \mathbf{x} above uses separate dimensions for each position!
 - i.e., *name* is two words to the right
 - what if *name* is one word to the right?

Convolution

\mathbf{c} = “feature map”, has an entry for each word position in context window / sentence

$$\mathbf{x} = [0.4 \ \dots \ 0.9 \ 0.2 \ \dots \ 0.7 \ 0.3 \ \dots \ 0.6]^\top$$

vector for *yo* vector for *last* vector for *name*

$$c_1 = \mathbf{w} \cdot \mathbf{x}_{1:d}$$

$$c_2 = \mathbf{w} \cdot \mathbf{x}_{d+1:2d}$$

$$c_3 = \mathbf{w} \cdot \mathbf{x}_{2d+1:3d}$$

Pooling

\mathbf{c} = “feature map”, has an entry for each word position in context window / sentence

how do we convert this into a fixed-length vector?

use **pooling**:

max-pooling: returns maximum value in \mathbf{c}

average pooling: returns average of values in \mathbf{c}

vector for *yo* vector for *last* vector for *name*

$$c_1 = \mathbf{w} \cdot \mathbf{x}_{1:d}$$

$$c_2 = \mathbf{w} \cdot \mathbf{x}_{d+1:2d}$$

$$c_3 = \mathbf{w} \cdot \mathbf{x}_{2d+1:3d}$$

Pooling

c = “feature map”, has an entry for each word position in context window / sentence

how do we convert this into a fixed-length vector?

use **pooling**:

max-pooling: returns maximum value in c

average pooling: returns average of values in c

vector for *yo* vector for *last* vector for *name*

$$c_1 = w \cdot x_{1:d}$$

then, this single filter w produces a single feature value (the output of some kind of pooling).

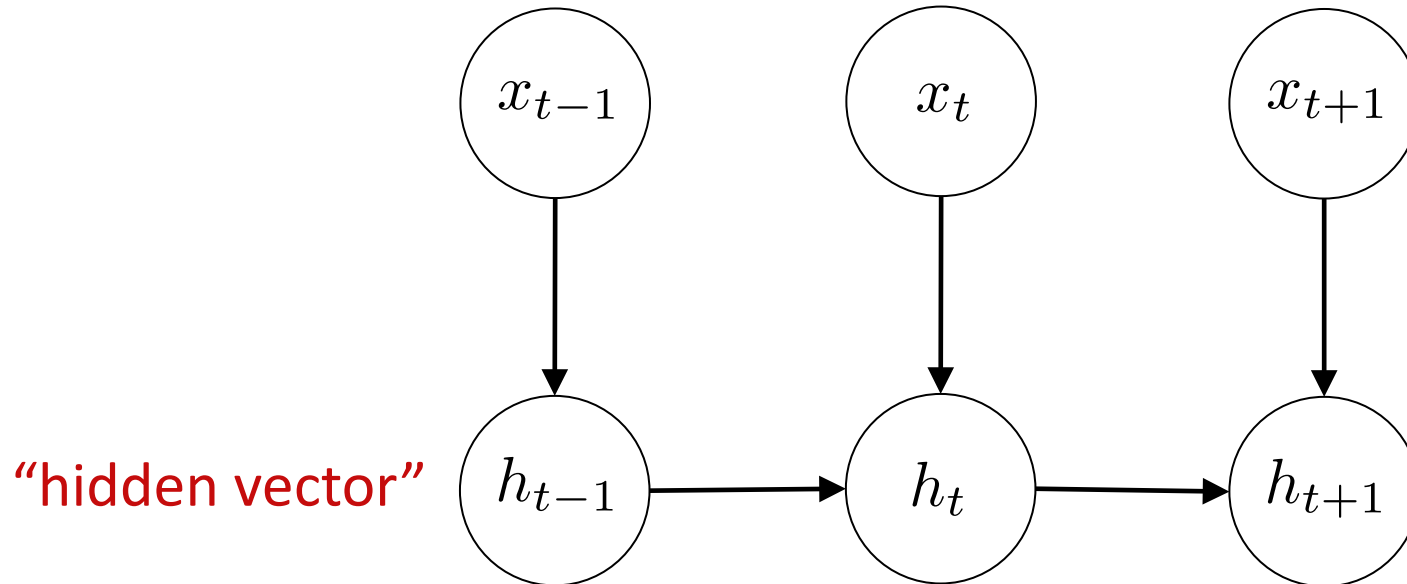
in practice, we use many filters of many different lengths (e.g., n -grams rather than words).

Convolutional Neural Networks

- convolutional neural networks (**convnets** or **CNNs**) use filters that are “convolved with” (matched against all positions of) the input
- think of convolution as “perform the same operation everywhere on the input in some systematic order”
- “convolutional layer” = set of filters that are convolved with the input vector (whether \mathbf{x} or hidden vector)
- could be followed by more convolutional layers, or by a type of pooling
- often used in NLP to convert a sentence into a feature vector

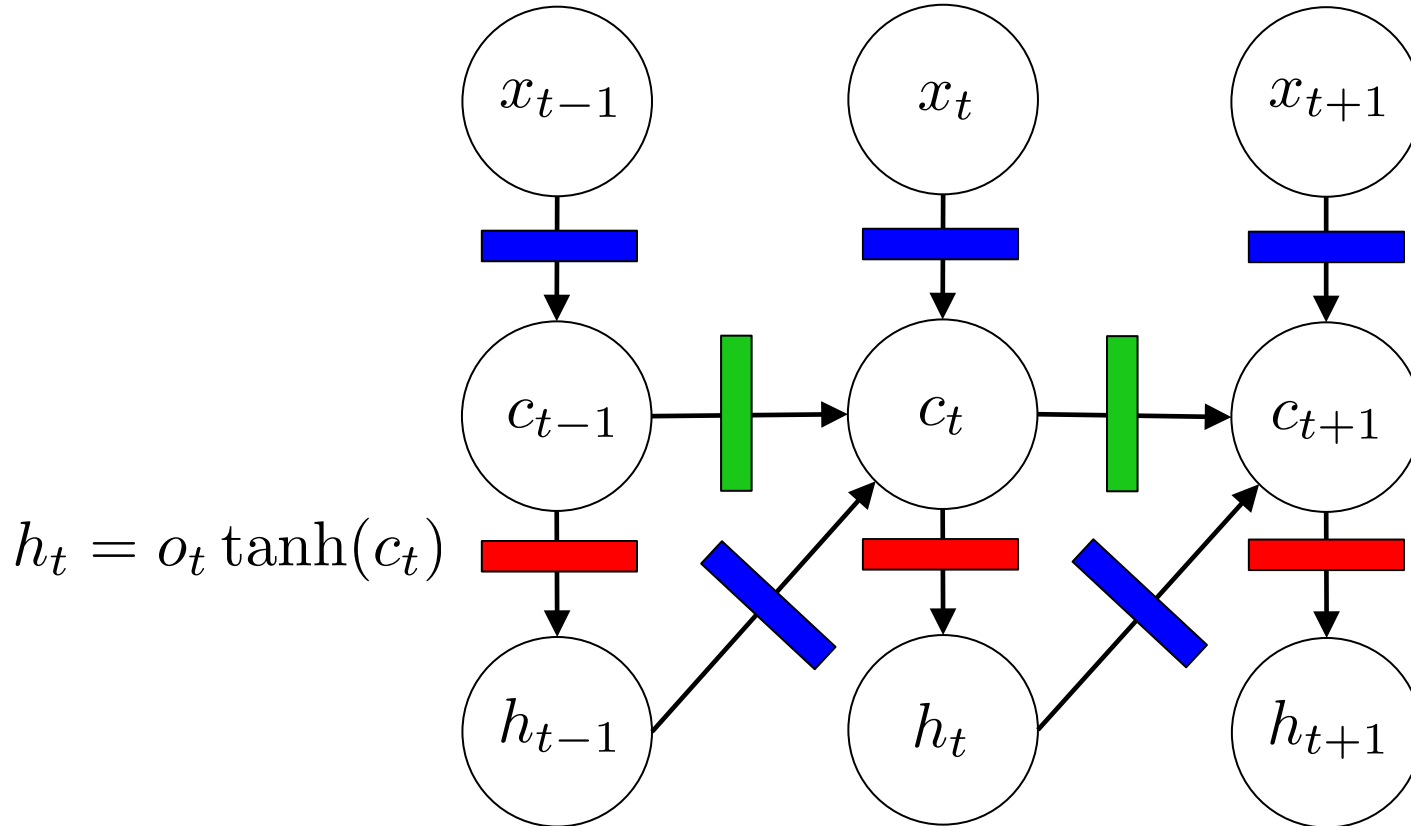
Recurrent Neural Networks

$$h_t = \tanh \left(W^{(xh)} x_t + W^{(hh)} h_{t-1} + b^{(h)} \right)$$



Long Short-Term Memory (LSTM) Recurrent Neural Networks

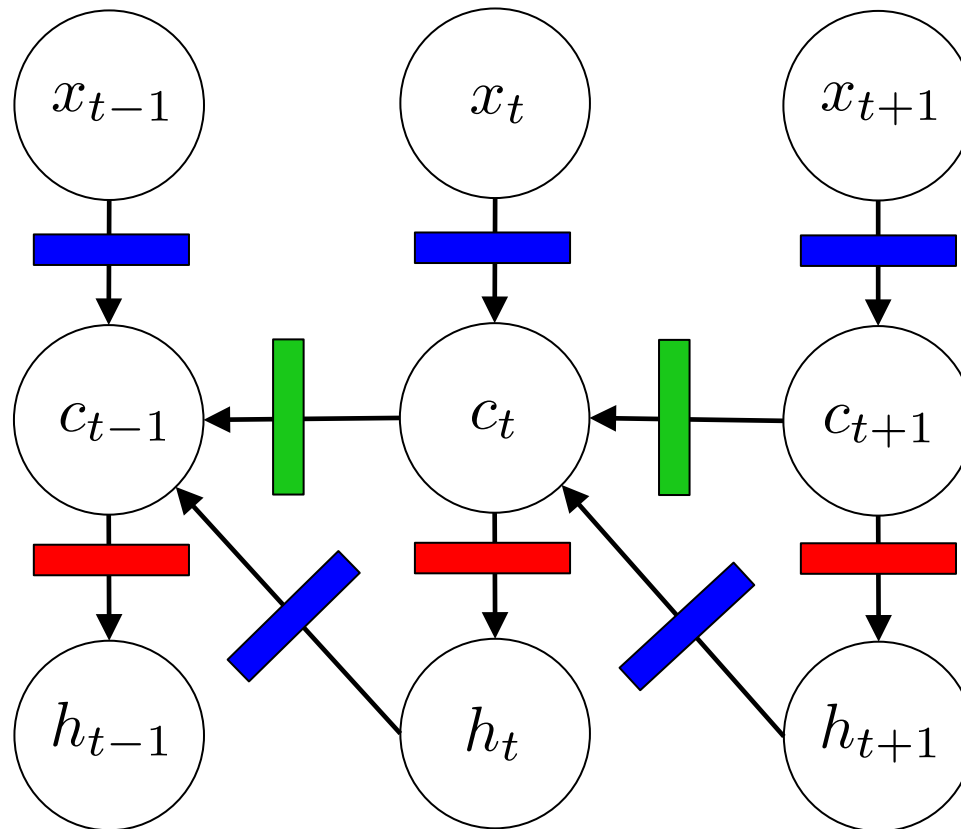
$$c_t = f_t c_{t-1} + i_t \tanh \left(W^{(xc)} x_t + W^{(hc)} h_{t-1} + b^{(c)} \right)$$



Backward & Bidirectional LSTMs

bidirectional:

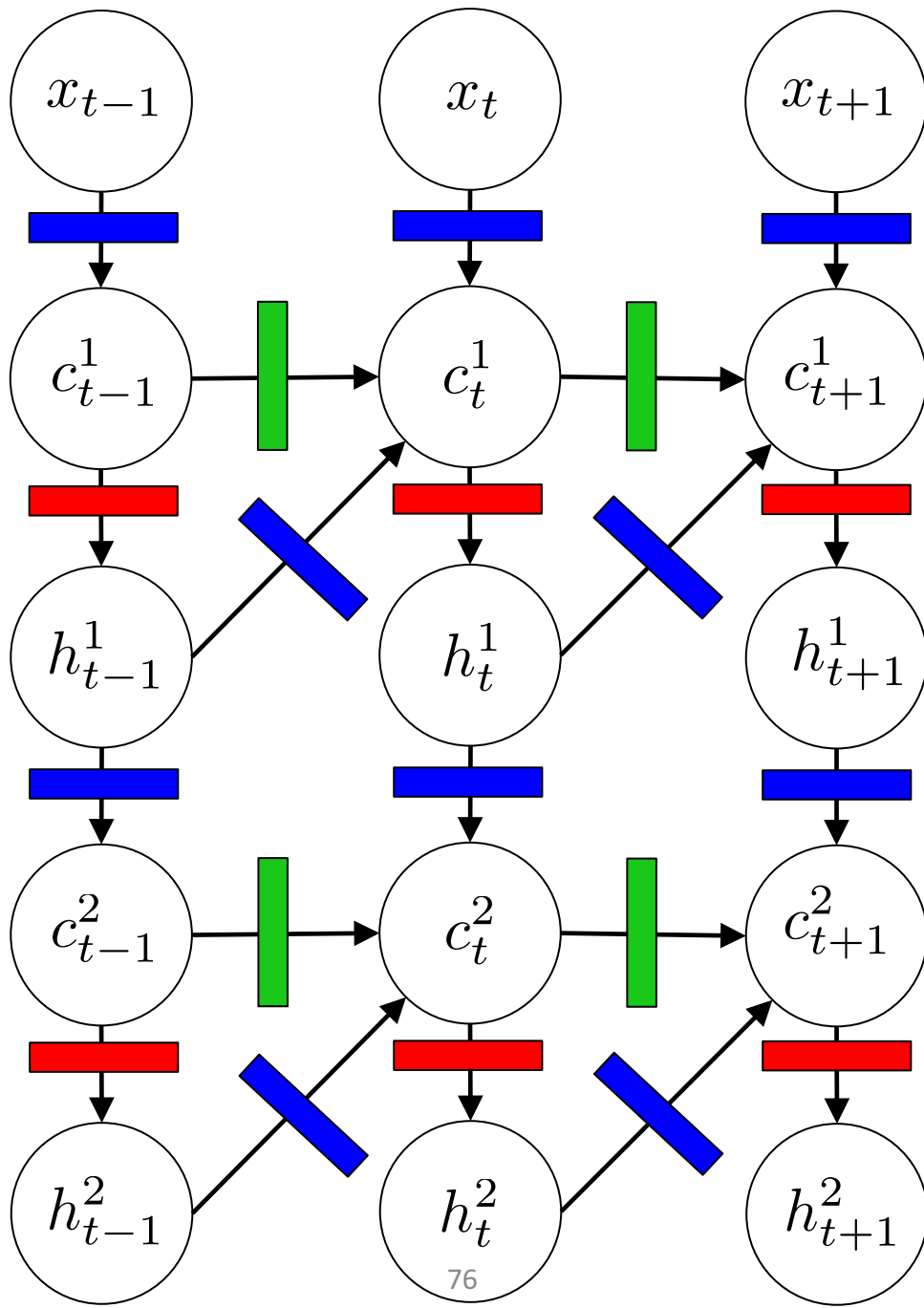
if shallow, just use forward and backward LSTMs in parallel, concatenate final two hidden vectors, feed to softmax



Deep LSTM (2-layer)

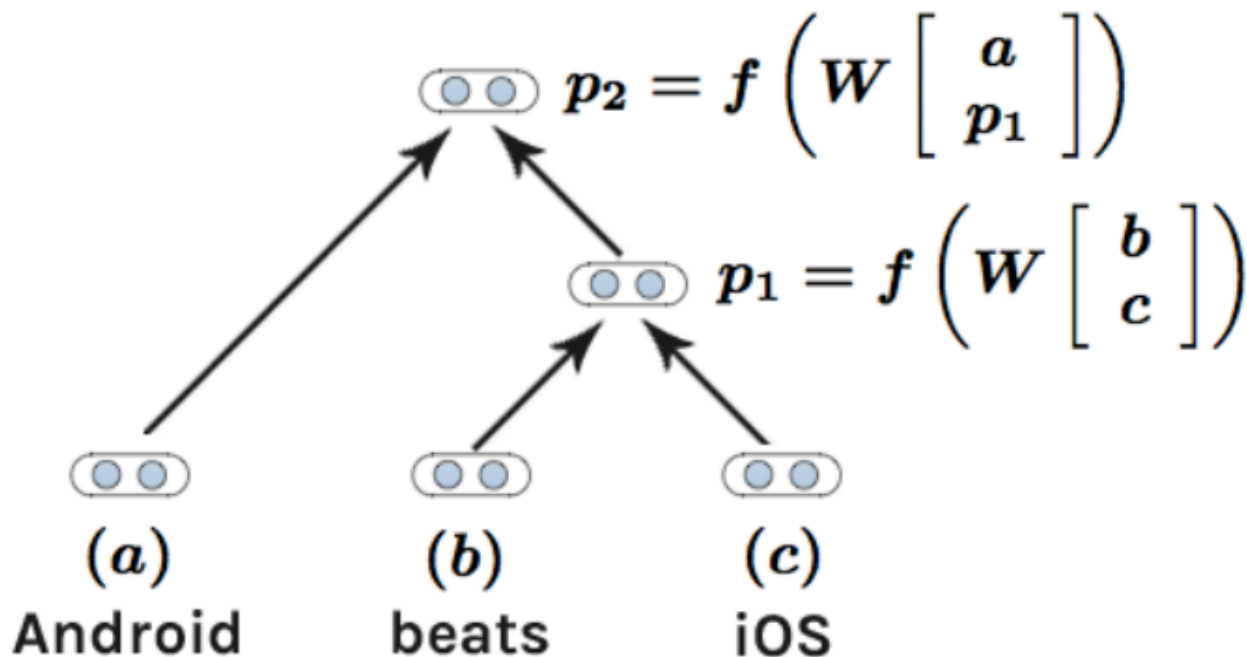
layer 1

layer 2



Recursive Neural Networks for NLP

- first, run a constituent parser on the sentence
- convert the constituent tree to a binary tree (each rewrite has exactly two children)
- construct vector for sentence recursively at each rewrite (“split point”):



Cost Functions

- **cost function**: scores output against a gold standard

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

- should reflect the evaluation metric for your task
- usual conventions: $\text{cost}(y, y) = 0$ $\text{cost}(y, y') = \text{cost}(y', y)$
- for classification, what cost should we use?

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