

TTIC 31190: Natural Language Processing

Kevin Gimpel
Winter 2016

Lecture 10:
Neural Networks for NLP

Announcements

- Assignment 2 due Friday
- project proposal due Tuesday, Feb. 16
- midterm on Thursday, Feb. 18

Roadmap

- classification
- words
- lexical semantics
- language modeling
- sequence labeling
- neural network methods in NLP
- syntax and syntactic parsing
- semantic compositionality
- semantic parsing
- unsupervised learning
- machine translation and other applications

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

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

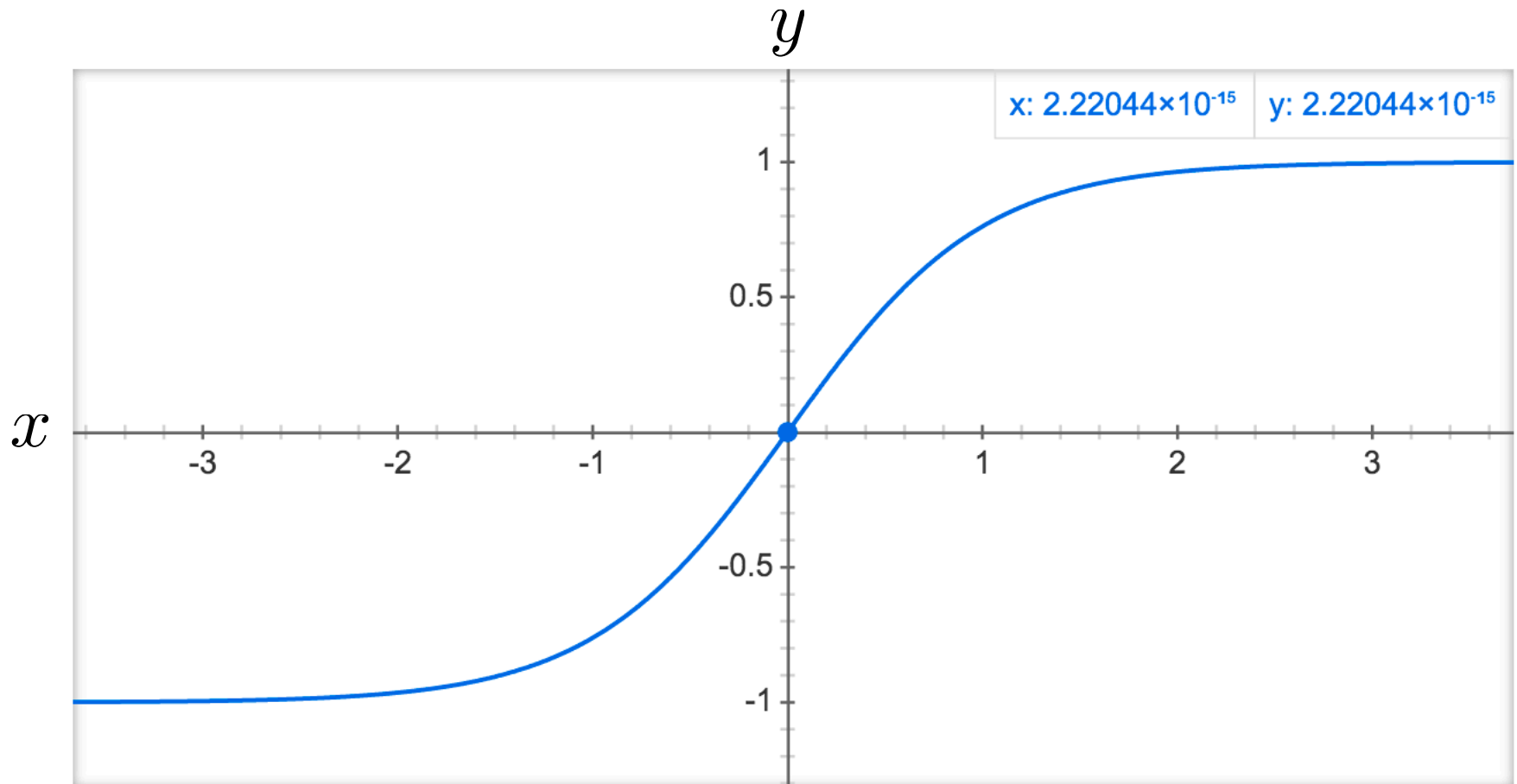
- this is a single “layer” of a neural network
- input vector is \mathbf{x}
- vector of “hidden units” is $\mathbf{z}^{(1)}$

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...

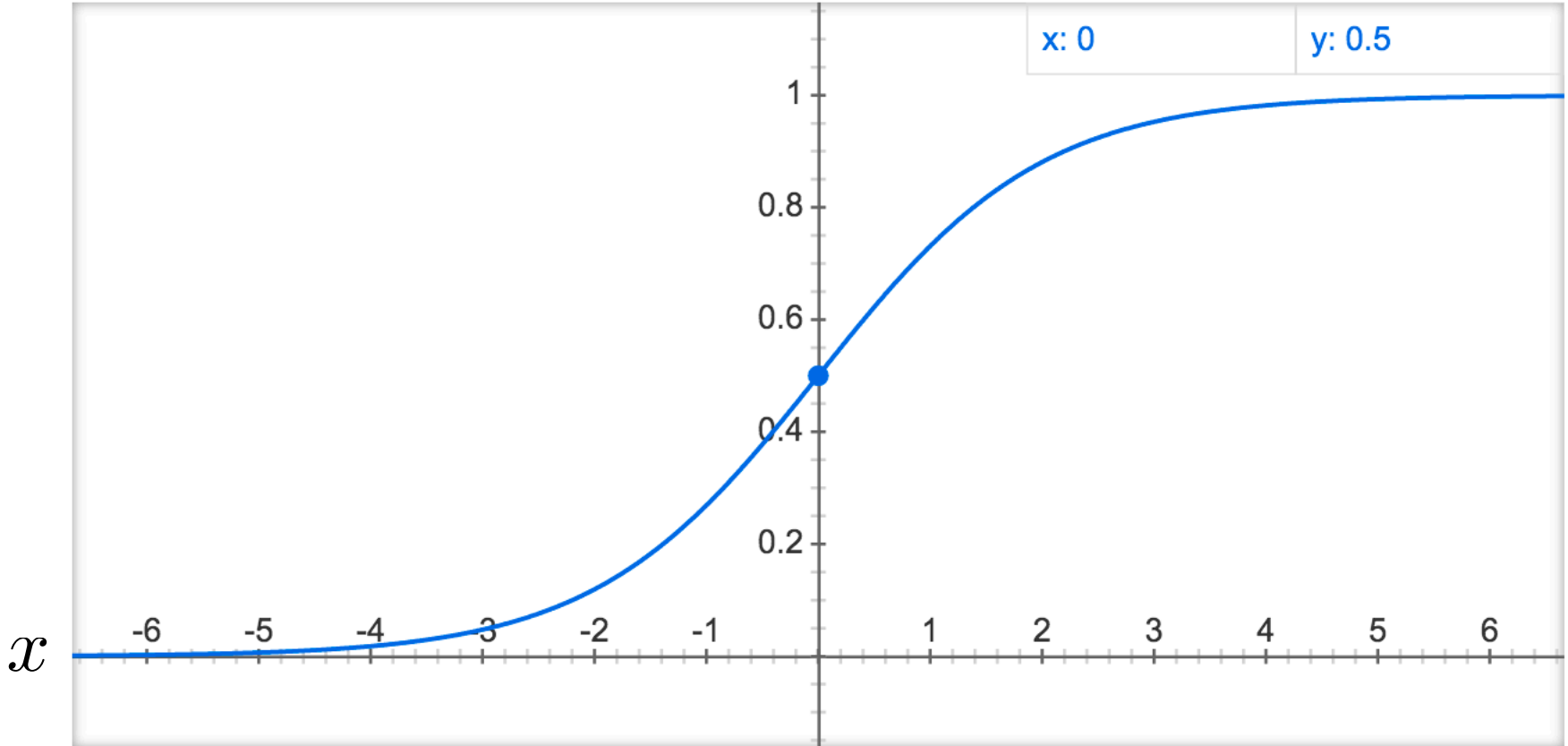
tanh: $y = \tanh(x)$



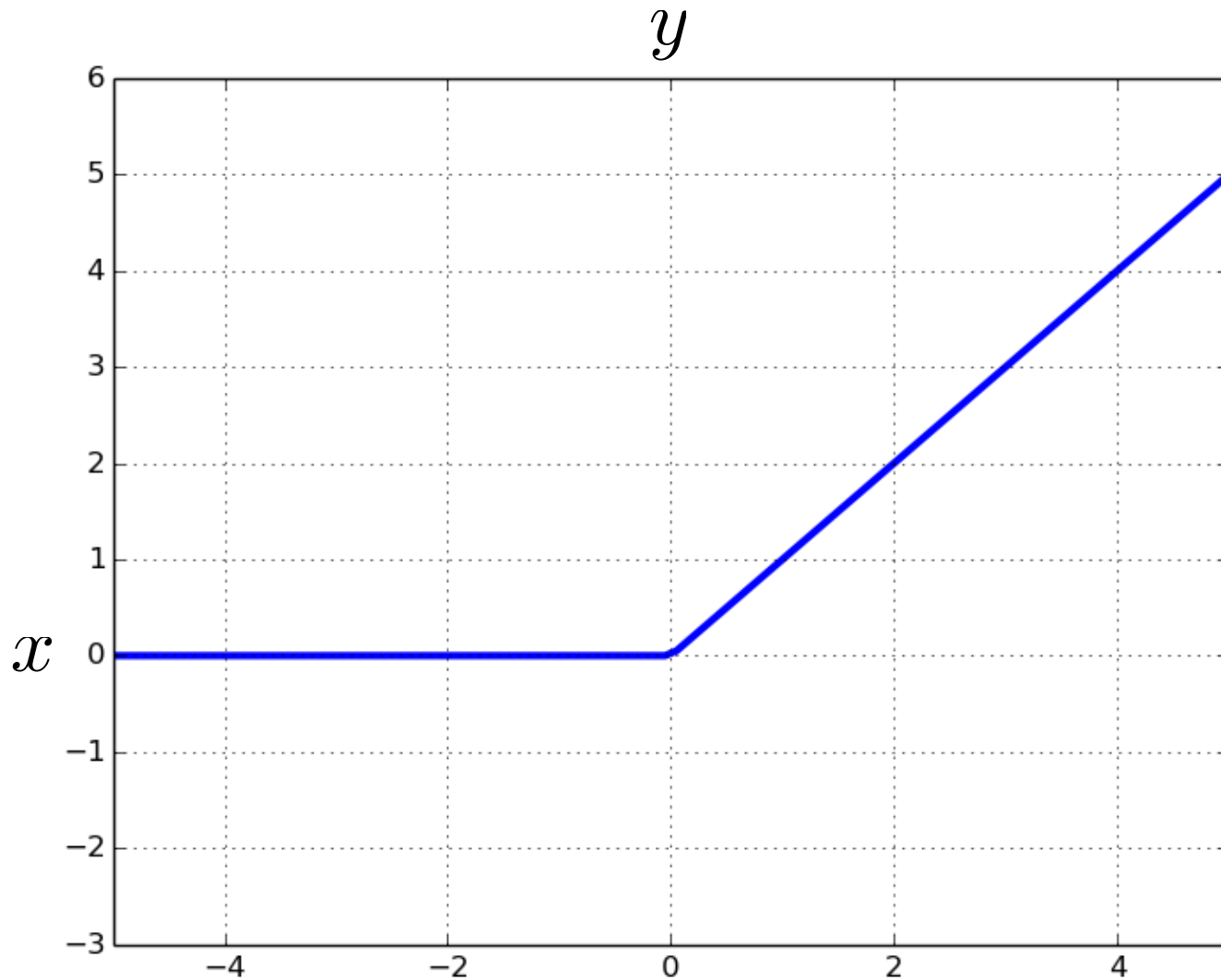
(logistic) sigmoid:

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

y



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



2-layer network

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

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

vector of label scores

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

2-layer neural network for sentiment classification

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

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



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

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)

Universal approximation theorem

From Wikipedia, the free encyclopedia

In the [mathematical](#) theory of [artificial neural networks](#), the **universal approximation theorem** states^[1] that a [feed-forward](#) network with a single hidden layer containing a finite number of [neurons](#) (i.e., a [multilayer perceptron](#)), can approximate [continuous functions](#) on [compact subsets](#) of \mathbf{R}^n , under mild assumptions on the activation function. The theorem thus states that simple neural networks can *represent* a wide variety of interesting functions when given appropriate parameters; however, it does not touch upon the algorithmic [learnability](#) of those parameters.

One of the first versions of the [theorem](#) was proved by [George Cybenko](#) in 1989 for [sigmoid](#) activation functions.^[2]

Kurt Hornik showed in 1991^[3] that it is not the specific choice of the activation function, but rather the multilayer feedforward architecture itself which gives neural networks the potential of being universal approximators. The output units are always assumed to be linear. For notational convenience, only the single output case will be shown. The general case can easily be deduced from the single output case.

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

Observation

- with linear models:
 - when using linear models for, say, sentiment classification, every feature included a label
 - **no** parameters were shared between labels
- with neural networks
 - we now have parameters shared across labels!
 - we still have some parameters that are devoted to particular labels
 - to define x , we design features that **only** look at the input (not at the labels)

Defining input features

- say we're doing sentiment classification and we want to use a neural network
- what should \mathbf{x} be?
 - it has to be independent of the label
 - it has to be a **fixed-length** vector

Empirical Risk Minimization with Surrogate Loss Functions

- 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
- we want to solve the following:

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

many possible loss
functions to consider
optimizing

Loss Functions

name	loss	where used
cost (“0-1”)	$\text{cost}(y, \text{classify}(\mathbf{x}, \boldsymbol{\theta}))$	intractable, but underlies “direct error minimization”
perceptron	$-\text{score}(\mathbf{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} \text{score}(\mathbf{x}, y', \boldsymbol{\theta})$	perceptron algorithm (Rosenblatt, 1958)
hinge	$-\text{score}(\mathbf{x}, y, \boldsymbol{\theta}) + \max_{y' \in \mathcal{L}} (\text{score}(\mathbf{x}, y', \boldsymbol{\theta}) + \text{cost}(y, y'))$	support vector machines, other large-margin algorithms
log	$-\log p_{\boldsymbol{\theta}}(y \mid \mathbf{x})$ $= \text{score}(\mathbf{x}, y, \boldsymbol{\theta}) + \log \sum_{y' \in \mathcal{L}} \exp\{\text{score}(\mathbf{x}, y', \boldsymbol{\theta})\}$	logistic regression, conditional random fields, maximum entropy models

(Sub)gradients of Losses for Linear Models

name	entry j of (sub)gradient of loss for linear model
cost ("0-1")	not subdifferentiable in general
perceptron	$-f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \hat{y})$, where $\hat{y} = \text{classify}(\mathbf{x}, \boldsymbol{\theta})$
hinge	$-f_j(\mathbf{x}, y) + f_j(\mathbf{x}, \tilde{y})$, where $\tilde{y} = \text{costClassify}(\mathbf{x}, y, \boldsymbol{\theta})$
log	$-f_j(\mathbf{x}, y) + \mathbb{E}_{p_{\boldsymbol{\theta}}(\cdot \mathbf{x})}[f_j(\mathbf{x}, \cdot)]$

Learning with Neural Networks

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

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

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

- 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

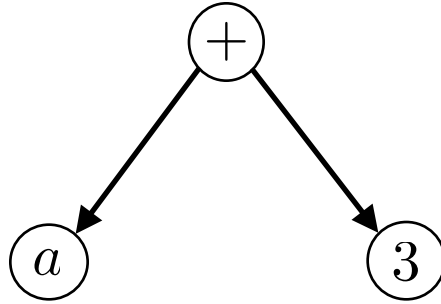
Computation Graphs

- 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, Torch, TensorFlow, CNTK, CNN, PENNE, etc.)

Backpropagation

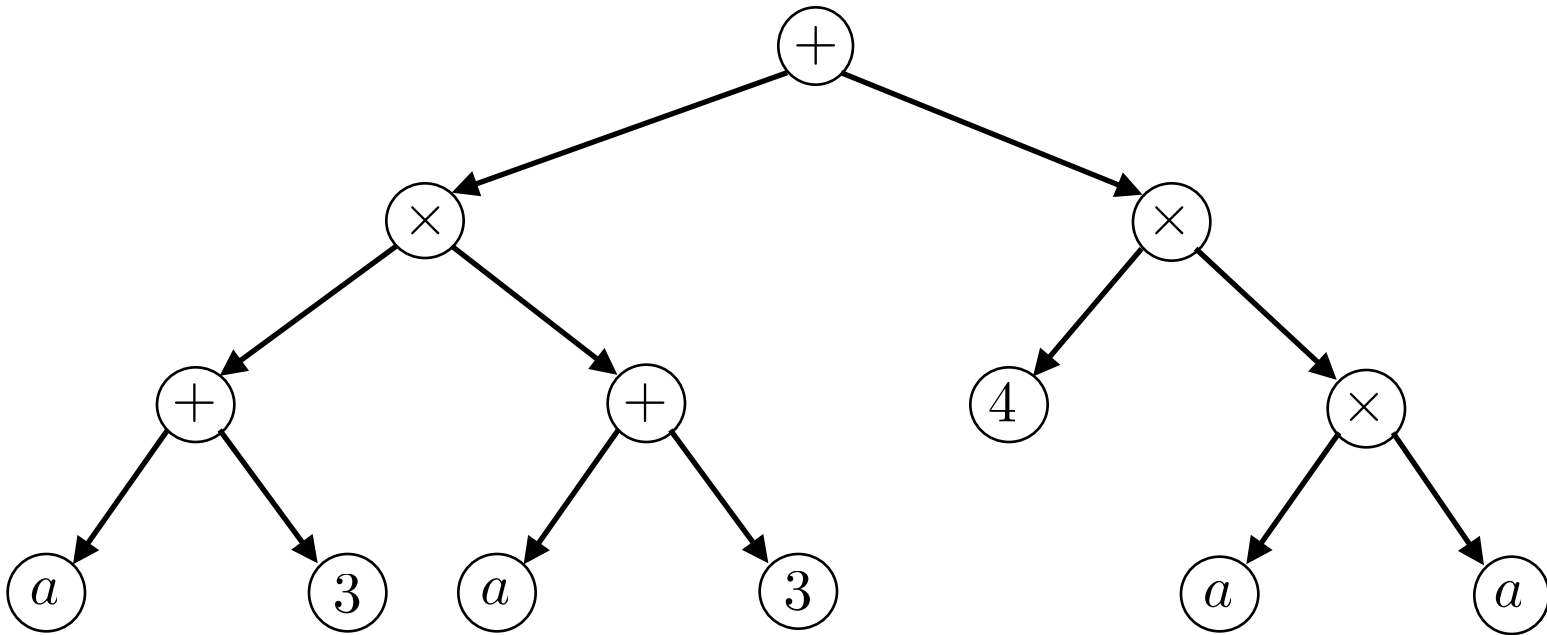
- backpropagation has become associated with neural networks, but it's much more general
- I also use backpropagation to compute gradients in **linear** models for structured prediction

A simple computation graph:



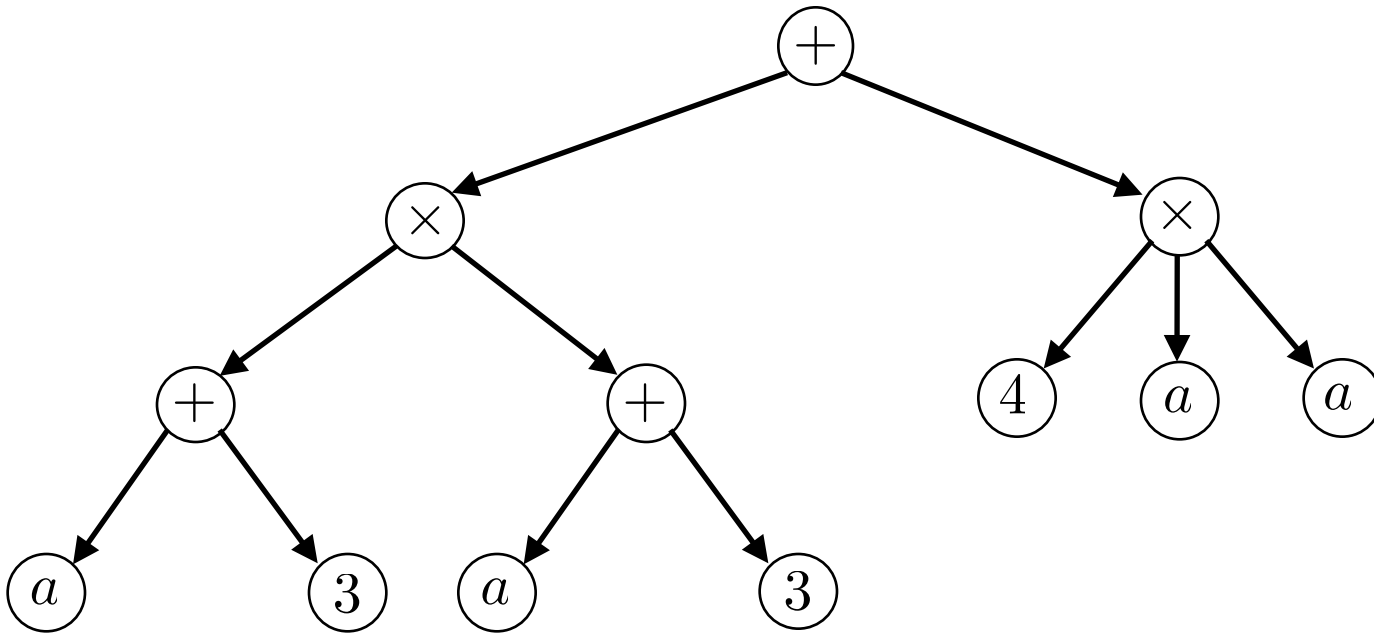
- represents expression “a + 3”

A slightly bigger computation graph:

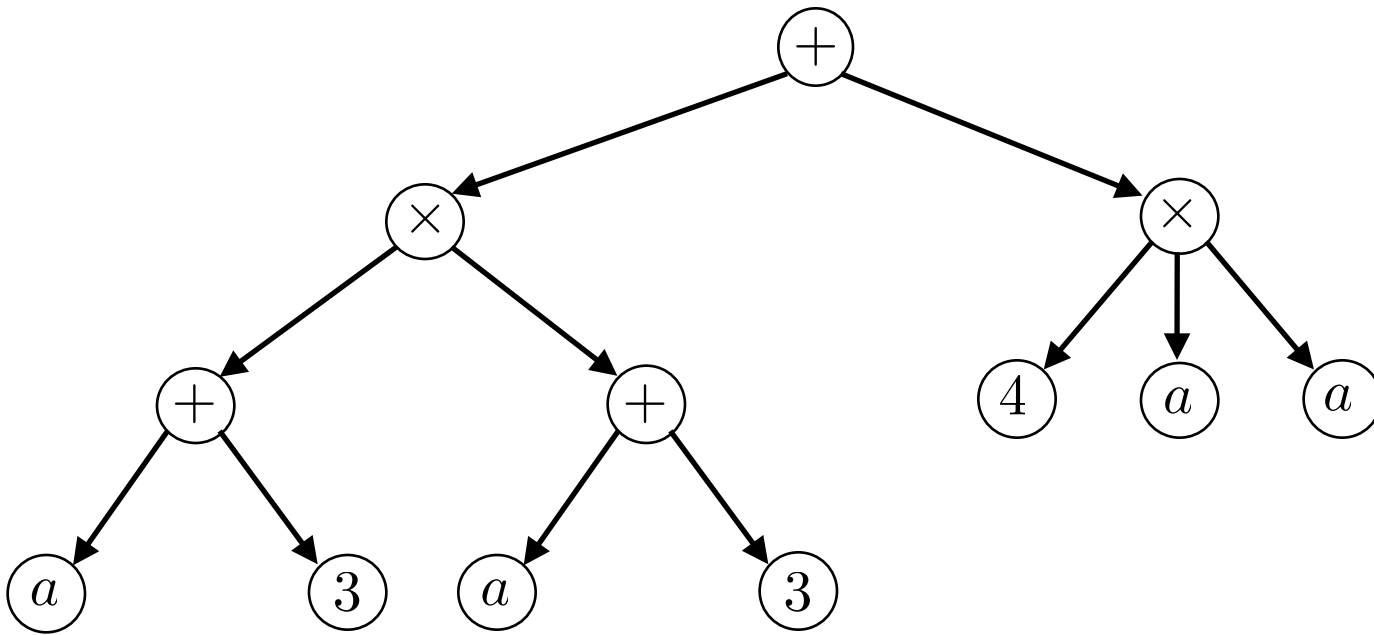


- represents expression $(a + 3)^2 + 4a^2$

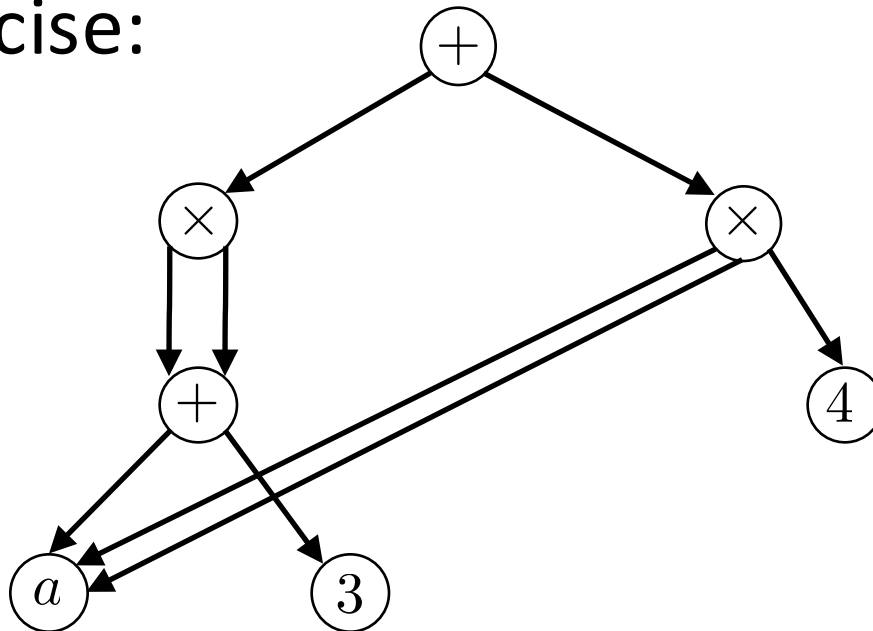
Operators can have more than 2 operands:



- 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

Empirical 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
- we want to solve the following:

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

Regularized Empirical 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

- we want to solve the following:

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

regularization
strength



regularization
term

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

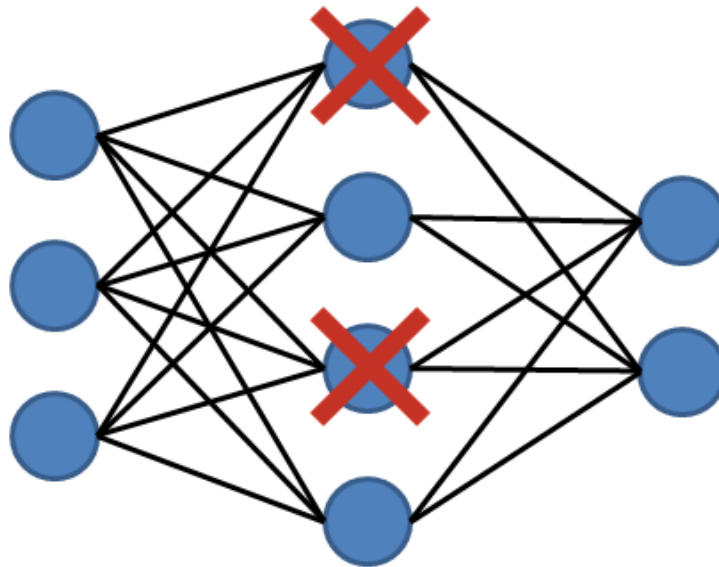
L_1 regularization: $R_{L_1}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = \sum_i |\theta_i|$

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

- you used stochastic gradient descent (SGD) in Assignment 1
- but there are many other choices:
 - AdaGrad
 - AdaDelta
 - 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)