

Week 2

CS224n NLP with Deeping Learning

Lecture 3 : Neural Networks

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

- take y^{th} row
- $W_y \cdot x = \sum_{i=1}^d W_{yi}x_i = f_y$

- objective: maximize the probability of the correct class y

\longleftrightarrow minimize negative log function

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

- Alternative: cross-entropy loss

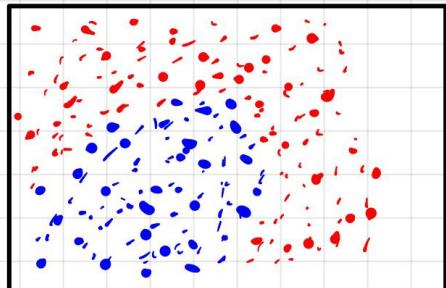
$$H(p, q) = - \sum_{c=1}^C p_{cc} \log q_{cc}$$

Assume the ground truth probability is 1 at the correct class,

the only term left is the negative log probability of the true class.

Neural Network Classifier

- notice basic classifier only gives linear boundary



\rightarrow a linear classifier won't work

- neural network can learn much more complex patterns and nonlinear decision boundary

- Common in NLP deep learning :

- learn both W and word vector x

- learn both parameters and representations

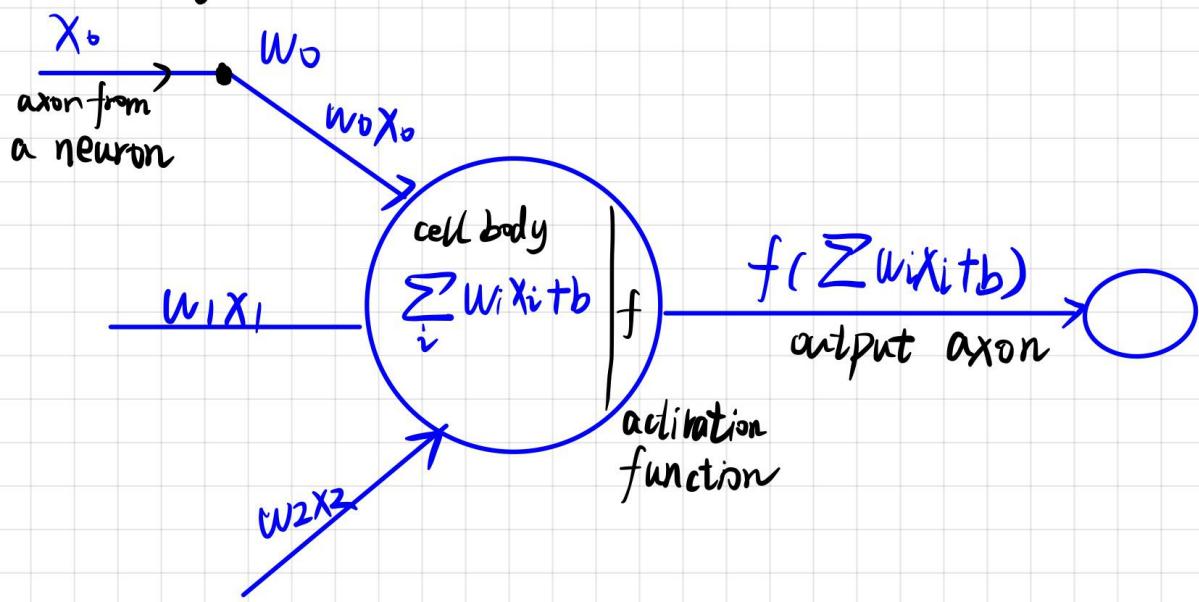
- The word vectors re-represent on-hot encoding — move them around in an intermediate layer vector space

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_1} \\ \vdots \\ \nabla_{W_d} \\ \nabla_x \text{ abandon} \\ \nabla_x \text{ zebra} \end{bmatrix}$$

$$\in \mathbb{R}^{Cd + \boxed{Vd}}$$

very large number of parameters

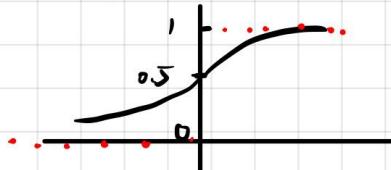
An artificial neuron



a neuron can be a binary logistic regression unit

$$h_{w,b}(x) = f(w^T x + b)$$

$$f(z) = \frac{1}{1 + e^{-z}}$$



a neural network = running several logistic regressions at the same time
... which we can feed into another logistic regression

Matrix notation :
 $z = Wx + b$
 $a = f(z)$

activation f is applied element-wise : $f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$

★ non-linearity is most important for deep neural network
(axiom: affine transformation's composition is still affine)

4. Named Entity Recognition (NER)

Task : find and classify names in Texts

Why difficult :

- boundary of entity
- if something is an entity
- hard to know class of novel entity
- entity depends on context

• Binary word window classification

Idea: classify a word in its context window

Δ: unnormalized scores

use neural activation "a" simple to give unnormalized score
 $\text{Score}(x) = \sum a \in \mathbb{R}$

add extra layer for non-linearity information.

How to compute $\nabla_{\theta} J(\theta)$? By hand or by back propagation

Gradients: 1. $f(x) = x^3$, $\frac{df}{dx} = 3x^2$

$$2. f(\mathbf{x}) = f(x_1, x_2, \dots, x_n), \quad \frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$$

$$3. f(\mathbf{x}) = [f_1(x_1, \dots, x_n), \dots, f_m(x_1, \dots, x_n)]$$

$$\text{Jacobian: } \frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

4. Chain Rule:

for multiple variables at once: multiply Jacobians

$$h = f(z)$$

$$z = Wx + b$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial x} = \dots$$

5. $h = f(z)$, element-wise activation function f

$$h_i = f(z_i)$$

What is $\frac{\partial h}{\partial z}$?

$$\left(\frac{\partial h}{\partial z} \right)_{ij} = \frac{\partial h_i}{\partial z_j} = \frac{\partial}{\partial z_j} f(z_i)$$

$$= \begin{cases} f'(z_i) & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Therefore: } \frac{\partial h}{\partial z} = \begin{pmatrix} f'(z_1) & & & \\ & \ddots & & \\ & & 0 & \\ 0 & & & f'(z_n) \end{pmatrix} = \text{diag}(f'(z))$$

other Jacobian

$$\frac{\partial}{\partial x} (Wx + b) = W$$

$$\frac{\partial}{\partial b} (Wx + b) = I \text{ (identity matrix)}$$

$$\frac{\partial}{\partial u} (u^T h) = h^T$$

Back to neural Network:

$$\begin{aligned} s &= u^T h \\ h &= f(Wx + b) \\ x &\text{ (input)} \end{aligned}$$

let's find $\frac{\partial s}{\partial b}$

$$\begin{aligned} &\Rightarrow \begin{cases} s = u^T h \\ h = f(z) \\ z = Wx + b \\ x \text{ (input)} \end{cases} \\ &\Rightarrow \frac{\partial s}{\partial b} = \frac{\partial s}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial b} \\ &= u^T \cdot f'(z) \cdot I \\ &= u^T \text{diag}(f'(z)) \cdot I \\ &= u^T \circ \text{diag}(f'(z)) \end{aligned}$$

may use "local error signal". δ

Instead of following convention: shape of the gradient is the shape of the parameters

7 Matrix Calculus Identity

- (1) Matrix times column vector with respect to the column vector
 $Z = Wx$, then $\frac{\partial Z}{\partial x} = W$
- (2) Row vector times matrix with respect to the row vector
 $Z = xW$, then $\frac{\partial Z}{\partial x} = W^T$
- (3) A vector with itself
 $Z = x$, then $\frac{\partial Z}{\partial x} = I$
- (4) An element-wise function applied to a vector
 $Z_i = f(x_i)$ $Z = f(x)$, then $\frac{\partial Z}{\partial x} = \text{diag}(f'(x))$
- (5) Matrix times column vector with respect to the matrix
 $Z = Wx$, $\delta = \frac{\partial J}{\partial Z}$, what is $\frac{\partial J}{\partial W} = \frac{\partial J}{\partial Z} \frac{\partial Z}{\partial W} = \delta \frac{\partial Z}{\partial W}$
for keep the shape, $\frac{\partial J}{\partial W} = \delta^T x^T$
- (6) Row vector times matrix with respect to the matrix
 $Z = xW$, $\delta = \frac{\partial J}{\partial Z}$, what is $\frac{\partial J}{\partial W} = \delta \frac{\partial Z}{\partial W}$
for keep the shape, $\frac{\partial J}{\partial W} = x^T \delta$
- (7) Cross-entropy loss with respect to logits
 $\hat{y} = \text{softmax}(\theta)$, $J = CE(\hat{y}, y)$
 $\frac{\partial J}{\partial \theta} = \hat{y} - y$

Lecture 4: Backpropagation

- The gradient that arrives at and updates the word vectors can be split up for each word vector

with $X_{\text{window}} = [X_{\text{museum}} \ X_{\text{in}} \ X_{\text{paris}} \ X_{\text{are}} \ X_{\text{amazing}}]$

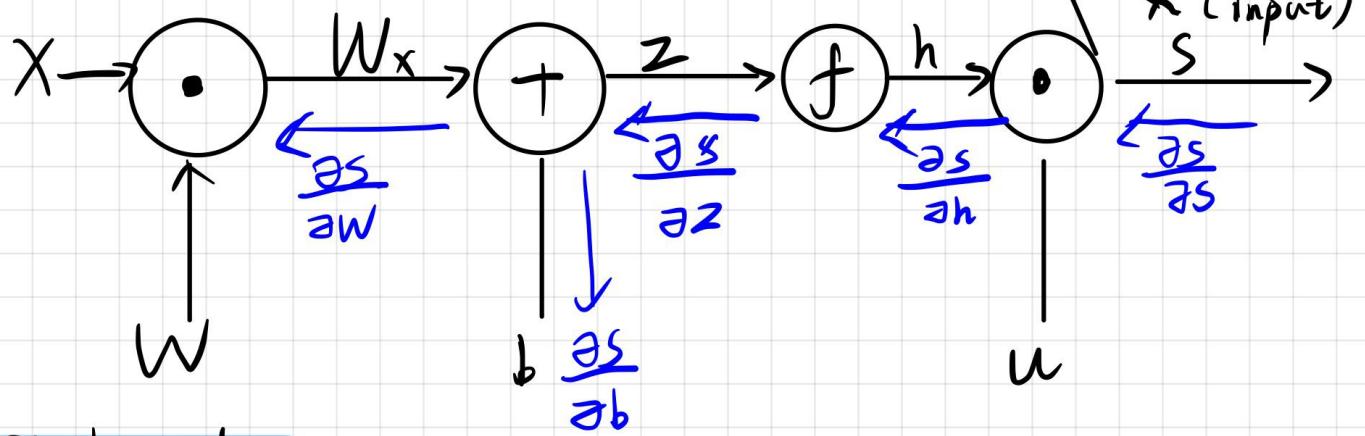
we have $S_{\text{window}} = \begin{bmatrix} \nabla X_{\text{museum}} \\ \nabla X_{\text{in}} \\ \nabla X_{\text{paris}} \\ \nabla X_{\text{are}} \\ \nabla X_{\text{amazing}} \end{bmatrix} \in \mathbb{R}^{5d}$

\Rightarrow This pushes word vectors around, might be helpful in determining named entity recognition.

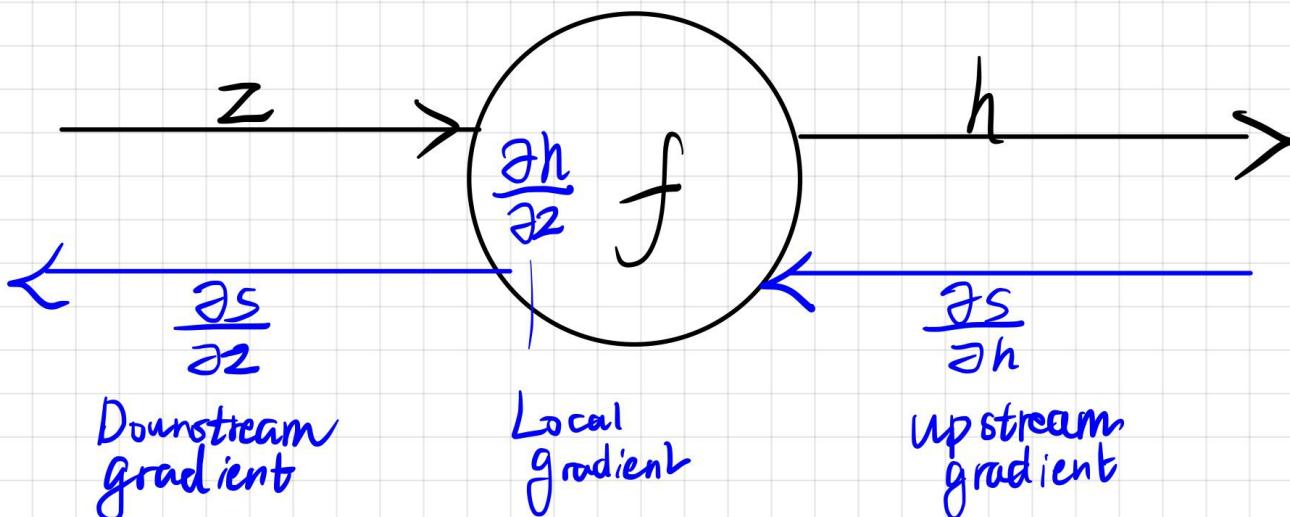
\Rightarrow pitfall: similar words may not be moved together if not appear in training data.

- Back propagation: we re-use derivative.

- Computational Graph:



Single node



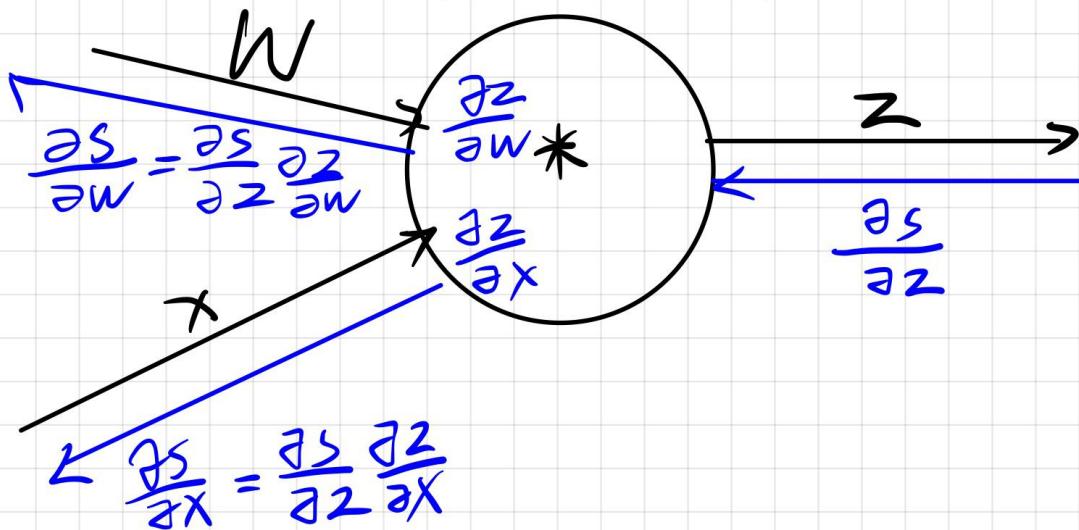
$$\boxed{\frac{\partial S}{\partial z} = \frac{\partial S}{\partial h} \frac{\partial h}{\partial z}}$$

Chain rule

[downstream gradient] = [upstream gradient] \times [Local gradient]

multiple inputs

$$z = \sqrt{x}$$



example :

Forward prop steps

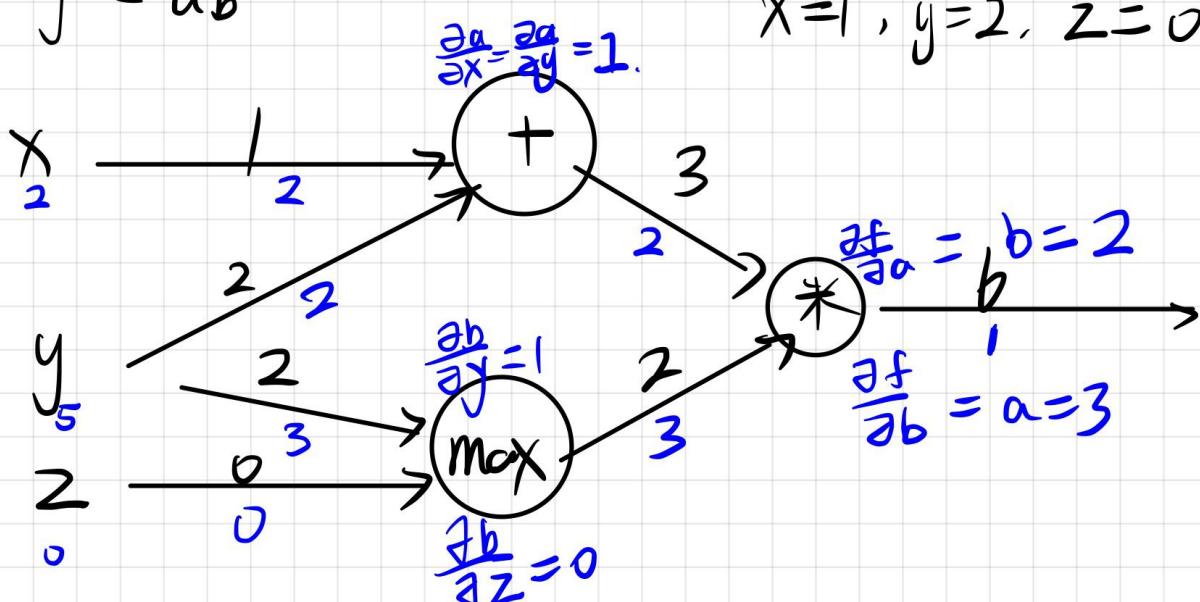
$$a = x + y$$

$$b = \max(y, z)$$

$$f = ab$$

$$f(x, y, z) = (x+y)\max(y, z)$$

$$x=1, y=2, z=0$$



*: Gradients sum at outward branches (加法, 非乘法)

"+" distribute

"*" switch — the upstream gradient

"max" route

correct way: avoid repeated computation.
compute all gradient (local) at once and save it for later usage.

Gradient checking: $f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$, for $h \approx 1e-4$

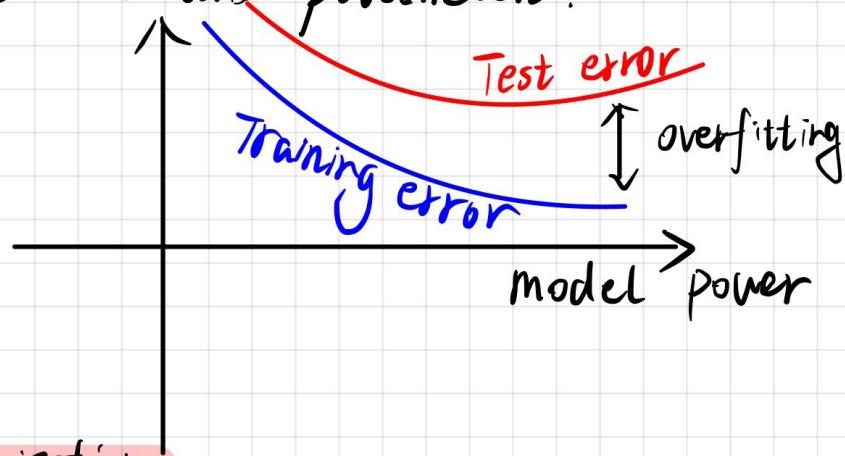
- easy to implement correctly

- But approximate and very slow

— have to recompute f for every parameter

Regularization

- In practice loss function usually includes regularization over all parameters Θ , e.g. L2-regularization
- Regularization (largely) prevents overfitting when have a lot of features and parameters.

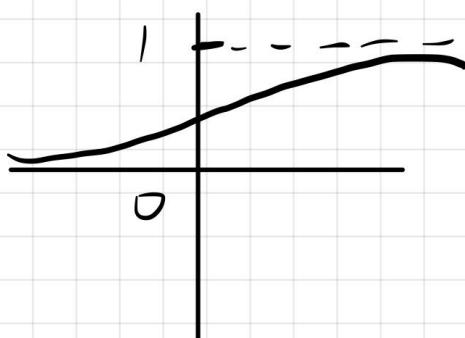


Vectorization

Non-linearity

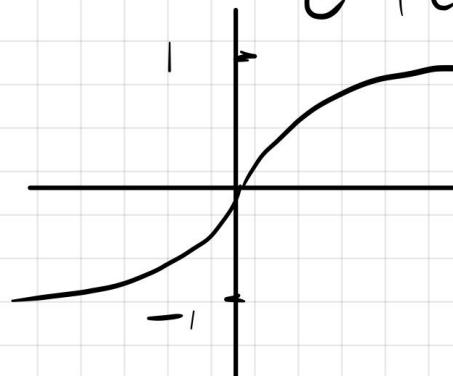
logistic ("sigmoid")

$$f(z) = \frac{1}{1 + \exp(-z)}$$

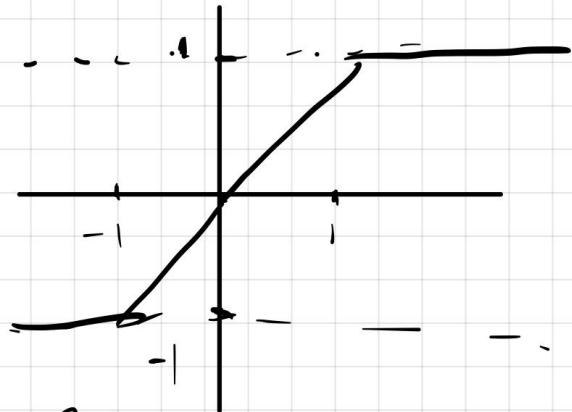


tanh

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



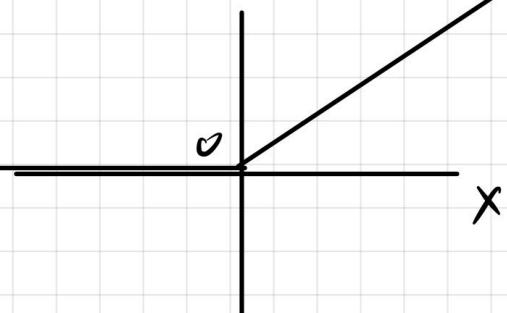
hard tanh



$$\text{hard tanh} = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } [-1, 1] \\ 1 & \text{if } x > 1 \end{cases}$$

if $x < -1$
if $[-1, 1]$
if $x > 1$

ReLU = $\max(x, 0)$



Parameter Initialization

- must initialize weights to small values
 - to avoid symmetries that prevent learning
- Xavier

$$\text{Var}(W_i) = \frac{2}{n_{\text{in}} + n_{\text{out}}}$$