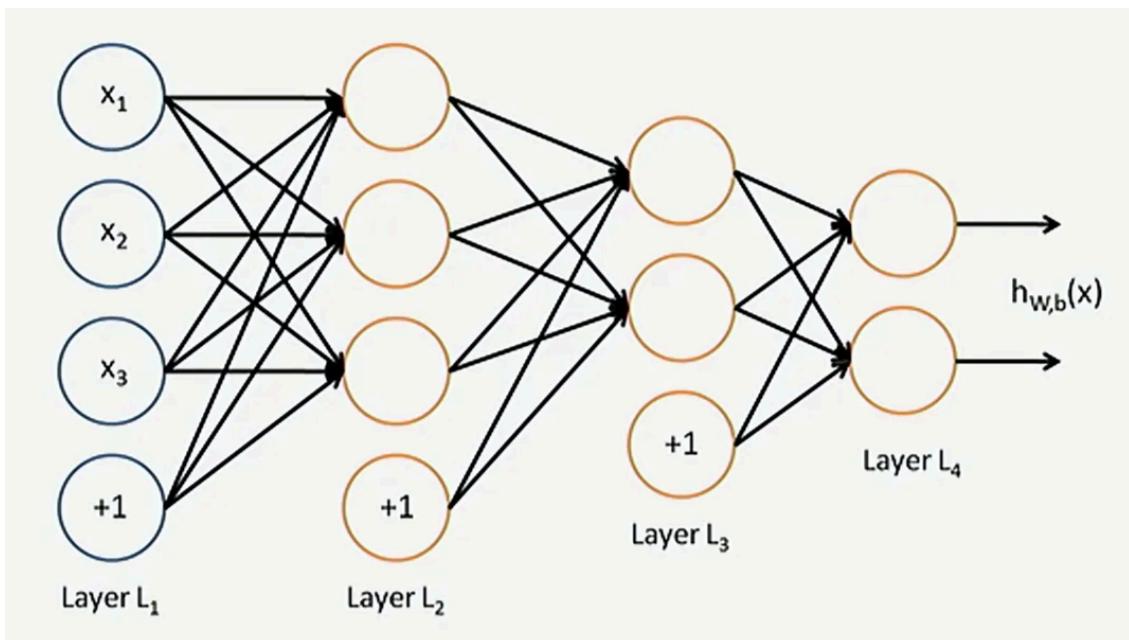


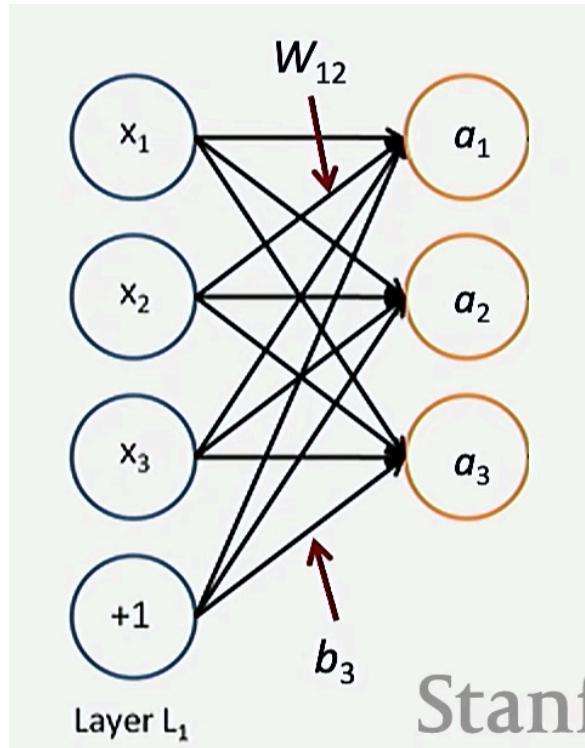
Lecture3

Spring 2024 : Backpropagation, Neural Network

Running several logistic regressions at the same time



non-linear



one output feeding with weights → next layer input (+bias term)

입력 없이도 출력을 조정할 수 있게 해주는 값

$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

etc.

— — —

$$z = Wx + b$$

$$a = f(z)$$

— — —

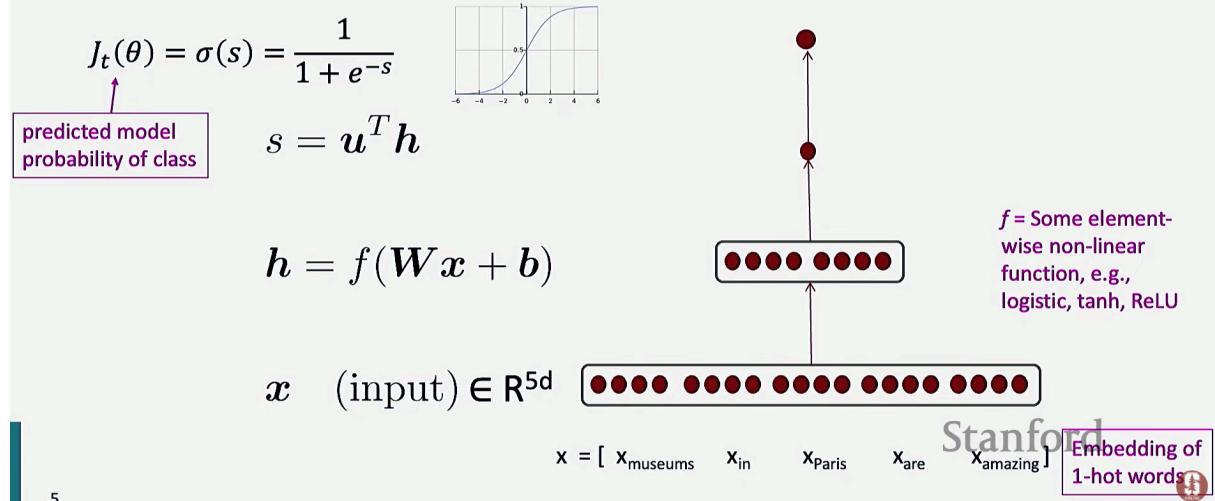
$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$

collapse to vector → metrics

scalar function : applying to each element of the vector

NER: Binary classification for center word being location

- We do supervised training and want high score if it's a location



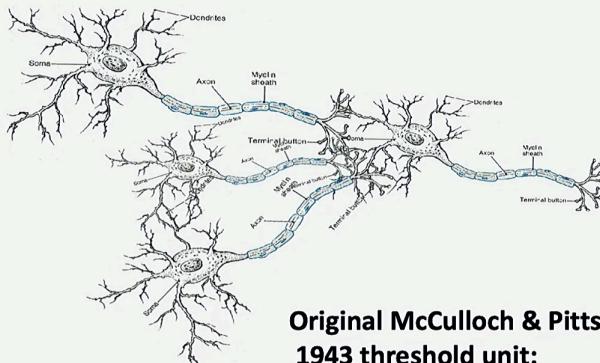
5

- Hidden Layer

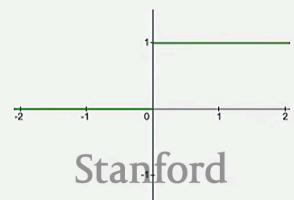
$$h = f(Wx + b)$$

function: non-linear ex) sigmoid

7. Neural computation



**Original McCulloch & Pitts
1943 threshold unit:**
 $\mathbf{1}(Wx > \theta)$
 $= \mathbf{1}(Wx - \theta > 0)$
 This function has no slope,
 so, no gradient-based learning



6

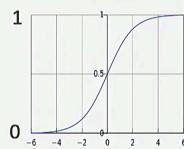
input과 weights의 linear combination $>\theta$ 에 따라 1 or 0 출력

→no gradient

Non-linearities, old and new

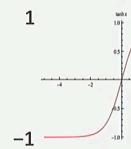
logistic ("sigmoid")

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



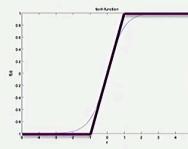
tanh

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



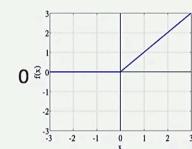
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$

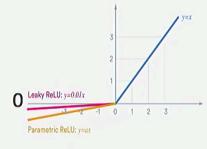


(Rectified Linear Unit)
ReLU

$$\text{ReLU}(z) = \max(z, 0)$$



Leaky ReLU /
Parametric ReLU



\tanh is just a rescaled and shifted sigmoid ($2 \times$ as steep, $[-1, 1]$):
 $\tanh(z) = 2\text{logistic}(2z) - 1$

Logistic and tanh are still used (e.g., logistic to get a probability)

However, now, for deep networks, the first thing to try is ReLU: it trains quickly and performs well due to good gradient backflow.

ReLU has a negative "dead zone" that recent proposals mitigate
GELU/Swish often used with Transformers (BERT, RoBERTa, etc.)

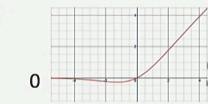
Swish [arXiv:1710.05941](https://arxiv.org/abs/1710.05941)

$$\text{swish}(x) = x \cdot \text{logistic}(x)$$

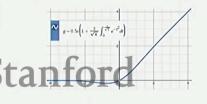
GELU [arXiv:1606.08415](https://arxiv.org/abs/1606.08415)

$$\text{GELU}(x) = x \cdot P(X \leq x), X \sim N(0, 1)$$

$$\approx x \cdot \text{logistic}(1.702x)$$



Stanford



③

logistic("sigmoid") : 음의 값 x

tanh: logistic * scaling +shift

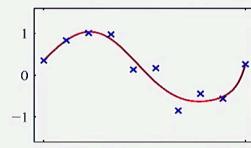
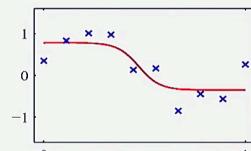
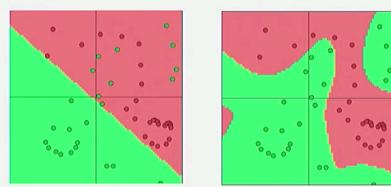
hard tanh: gradient :1 효과적

ReLU: 음수 입력에선 gradient 0 → "dead zone"

Leaky ReLU / Parametric ReLU: 음수도 기울기

Non-linearities (i.e., " f " on previous slide): Why they're needed

- Neural networks do function approximation, e.g., regression or classification
 - Without non-linearities, deep neural networks can't do anything more than a linear transform
 - Extra layers could just be compiled down into a single linear transform: $W_1 W_2 x = Wx$
 - But, with more layers that include non-linearities, they can approximate any complex function!



Stanford

왜 비선형적이어야하는가?

to make function approximation

층을 여러 개 쌓아도 모두 선형이면 결국 **하나의 선형 변환**과 같음:

$$W_1 W_2 x = Wx$$

Jacobian Matrix: Generalization of the Gradient

- Given a function with **m outputs** and **n inputs**

$$\mathbf{f}(\mathbf{x}) = [f_1(x_1, x_2, \dots, x_n), \dots, f_m(x_1, x_2, \dots, x_n)]$$

- It's Jacobian is an **$m \times n$ matrix** of partial derivatives

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \boxed{\left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)_{ij} = \frac{\partial f_i}{\partial x_j}}$$

Jacobian Matrix: 각 f_i 에 대한 gradient vectors

$$J = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \in \mathbb{R}^{m \times n}$$

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

$i=j$ 일 때만 미분 값 존재

$i \neq j$ 이면 서로 독립 \rightarrow 미분값 0

\rightarrow Diagonal Matrix 꼴

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

$$\frac{\partial}{\partial \mathbf{b}}(\mathbf{W}\mathbf{x} + \mathbf{b}) = \mathbf{I} \text{ (Identity matrix)}$$

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

1.

- \mathbf{W} 는 고정된 행렬, x 는 벡터
- $\mathbf{W}\mathbf{x}$ 는 선형변환
- x 에 대해 편미분하면 변화율이 \mathbf{W} 로 일정

2.

- \mathbf{b} 는 벡터
- $\mathbf{W}\mathbf{x}$ 는 \mathbf{b} 와 무관하므로 미분하면 0
- $\frac{\partial \mathbf{b}}{\partial \mathbf{b}} = \mathbf{I}$

3.

- $\mathbf{u}^T \mathbf{h}$ 는 벡터 내적 \rightarrow 결과는 스칼라 값
- \mathbf{u} 에 대해 미분하면 내적에서는 **상대 벡터(\mathbf{h})**가 남게 됨

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}, \quad \mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \Rightarrow \mathbf{u}^\top \mathbf{h} = u_1 h_1 + u_2 h_2 + u_3 h_3$$

$$\begin{aligned} \frac{\partial}{\partial u_1}(\mathbf{u}^\top \mathbf{h}) &= h_1 \\ \frac{\partial}{\partial u_2}(\mathbf{u}^\top \mathbf{h}) &= h_2 \\ \frac{\partial}{\partial u_3}(\mathbf{u}^\top \mathbf{h}) &= h_3 \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{u}}(\mathbf{u}^\top \mathbf{h}) = [h_1 \quad h_2 \quad h_3] = \mathbf{h}^\top$$

An Example

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

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

Forward prop steps

$$a = x + y$$

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

$$f = ab$$

$$\frac{\partial f}{\partial x} = 2$$

$$\frac{\partial f}{\partial y} = 3 + 2 = 5$$

$$\frac{\partial f}{\partial z} = 0$$

Local gradients

$$\frac{\partial a}{\partial x} = 1 \quad \frac{\partial a}{\partial y} = 1$$

$$\frac{\partial b}{\partial y} = \mathbf{1}(y > z) = 1 \quad \frac{\partial b}{\partial z} = \mathbf{1}(z > y) = 0$$

$$\frac{\partial f}{\partial a} = b = 2 \quad \frac{\partial f}{\partial b} = a = 3$$

$$x \xrightarrow{1}$$

$$y \xrightarrow{2}$$

$$z \xrightarrow{0}$$

$$a \xrightarrow{2}$$

$$b \xrightarrow{3}$$

$$f \xrightarrow{6}$$

$$x \xrightarrow{2}$$

$$y \xrightarrow{3}$$

$$z \xrightarrow{0}$$

$$+ \xrightarrow{2}$$

$$\max \xrightarrow{3}$$

$$* \xrightarrow{1}$$

$$a = x + y = 1 + 2 = 3$$

$$b = \max(y, z) = \max(2, 0) = 2$$

$$f = ab = 3 \cdot 2 = 6$$

$$\frac{\partial f}{\partial a} = b = 2$$

$$\frac{\partial f}{\partial b} = a = 3$$

$$\frac{\partial a}{\partial x} = 1, \quad \frac{\partial a}{\partial y} = 1$$

$$(\text{Since } y > z) \Rightarrow \frac{\partial b}{\partial y} = 1, \quad \frac{\partial b}{\partial z} = 0$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \cdot \frac{\partial a}{\partial x} = 2 \cdot 1 = \boxed{2}$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial a} \cdot \frac{\partial a}{\partial y} + \frac{\partial f}{\partial b} \cdot \frac{\partial b}{\partial y} = 2 \cdot 1 + 3 \cdot 1 = \boxed{5}$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial b} \cdot \frac{\partial b}{\partial z} = 3 \cdot 0 = \boxed{0}$$

□

Deriving local input gradient in backprop

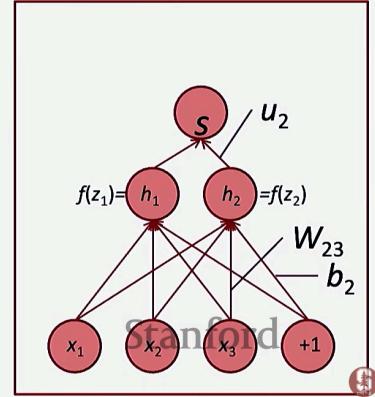
- For $\frac{\partial \mathbf{z}}{\partial W}$ in our equation:

$$\frac{\partial s}{\partial W} = \delta \frac{\partial \mathbf{z}}{\partial W} = \delta \frac{\partial}{\partial W} (Wx + b)$$

- Let's consider the derivative of a single weight W_{ij}
- W_{ij} only contributes to z_i
 - For example: W_{23} is only used to compute z_2 not z_1

$$\begin{aligned}\frac{\partial z_i}{\partial W_{ij}} &= \frac{\partial}{\partial W_{ij}} W_i \cdot x + b_i \\ &= \frac{\partial}{\partial W_{ij}} \sum_{k=1}^d W_{ik} x_k = x_j\end{aligned}$$

45



What shape should derivatives be?

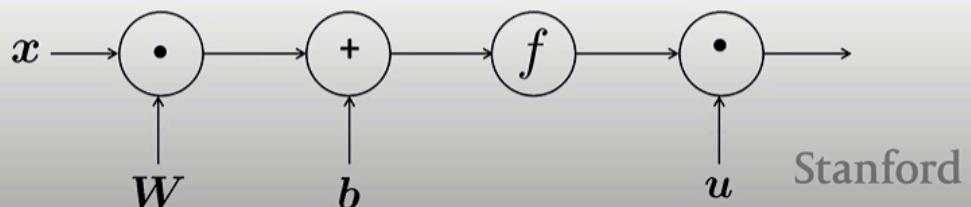
Two options for working through specific problems:

- Use Jacobian form as much as possible, reshape to follow the shape convention at the end:
 - What we just did. But at the end transpose $\frac{\partial s}{\partial b}$ to make the derivative a column vector, resulting in δ^T
- Always follow the shape convention
 - Look at dimensions to figure out when to transpose and/or reorder terms
 - The error message δ that arrives at a hidden layer has the same dimensionality as that hidden layer

Stanford

Computation Graphs and Backpropagation

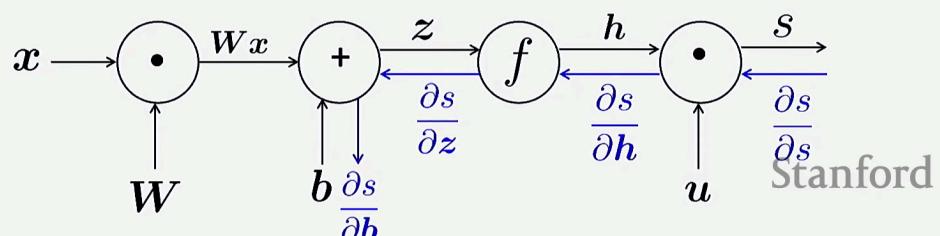
- Software represents our neural net equations as a graph
 - Source nodes: inputs
 - Interior nodes: operations
- $$s = \mathbf{u}^T \mathbf{h}$$
- $$\mathbf{h} = f(\mathbf{z})$$
- $$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$
- $$\mathbf{x} \quad (\text{input})$$



Forward Propagation

Backpropagation

- Then go backwards along edges
 - Pass along **gradients**
- $$s = \mathbf{u}^T \mathbf{h}$$
- $$\mathbf{h} = f(\mathbf{z})$$
- $$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$
- $$\mathbf{x} \quad (\text{input})$$

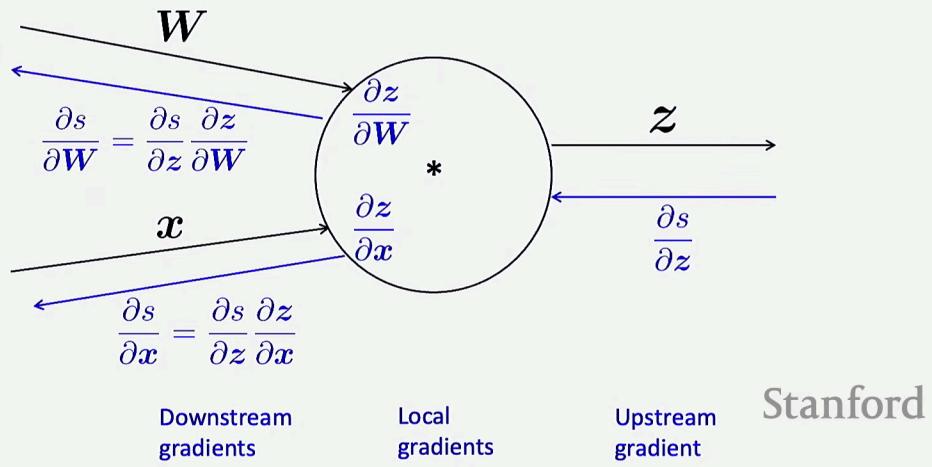


52

Backpropagation: Single Node

- Multiple inputs → multiple local gradients

$$z = Wx$$



58

Downstream gradients = local gradients * Upstream gradients

Same Big O

1. 정방향 계산 (Forward)

$$z = Wx + b, \quad h = f(z)$$

시간 복잡도:

- $W \in \mathbb{R}^{n \times m}, x \in \mathbb{R}^m$
- $\rightarrow Wx$ 의 계산 비용은 $O(nm)$

2. 역전파 계산 (Backward)

우리는 다음을 계산해야 함:

- $\frac{\partial L}{\partial W}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial x}$
- 이 역시 행렬-벡터 곱, 원소곱, 전치 곱 등으로 구성되어 있음
- \Rightarrow 계산량도 역시 $O(nm)$

Summary

We've mastered the core technology of neural nets! 

- **Backpropagation:** recursively (and hence efficiently) apply the chain rule along computation graph
 - $[\text{downstream gradient}] = [\text{upstream gradient}] \times [\text{local gradient}]$
- **Forward pass:** compute results of operations and save intermediate values
- **Backward pass:** apply chain rule to compute gradients