
CS224n Study

Back Propagation

DT 추진단 송석민

Why Again BackPropagation??

= Back Prop이 항상 잘 작동하지는 않아서

= 학습이 잘 안될 때, 모델을 개선시키려고

Remember: Stochastic Gradient Descent

- Update equation:

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$

$\alpha = \text{step size or learning rate}$

Back Propagation Equation

$Node \rightarrow Node$
오차를 전달

Computational Graph

$f(x)$ 미분값 계산하기

- This Lecture: How do we compute $\nabla_{\theta} J(\theta)$?
 - By hand
 - Algorithmically (the backpropagation algorithm)

Gradients

- Given a function with 1 output and n inputs

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$$

- Its gradient is a vector of partial derivatives

$$\frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

Jacobian Matrix: Generalization of the Gradient

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

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

- Its Jacobian is an **$m \times n$ matrix** of partial derivatives

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

$$\left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

Example Jacobian: Activation Function

$$\mathbf{h} = f(\mathbf{z}), \text{ what is } \frac{\partial \mathbf{h}}{\partial \mathbf{z}}? \quad \mathbf{h}, \mathbf{z} \in \mathbb{R}^n$$

$$h_i = f(z_i)$$

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

definition of Jacobian

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

regular 1-variable derivative

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

Other Jacobians

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

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

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

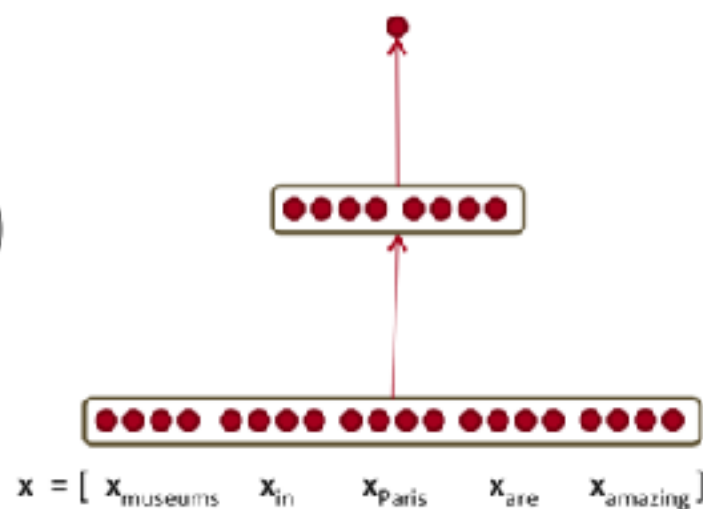
Back to Neural Nets!

- Let's find $\frac{\partial s}{\partial b}$
 - In practice we care about the gradient of the loss, but we will compute the gradient of the score for simplicity

$$s = u^T h$$

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

$$x \text{ (input)}$$



$$s = u^T h$$

$$h = f(z)$$

$$z = Wx + b$$

$$x \text{ (input)}$$

Useful Jacobians from previous slide

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

$$\frac{\partial}{\partial z}(f(z)) = \text{diag}(f'(z))$$

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

$$\begin{aligned} \frac{\partial s}{\partial b} &= \frac{\partial s}{\partial h} \quad \frac{\partial h}{\partial z} \quad \frac{\partial z}{\partial b} \\ &\quad \downarrow \quad \downarrow \quad \downarrow \\ &= u^T \text{diag}(f'(z)) I \\ &= u^T \circ f'(z) \end{aligned}$$

Derivative with respect to Matrix

- What does $\frac{\partial s}{\partial \mathbf{W}}$ look like? $\mathbf{W} \in \mathbb{R}^{n \times m}$
- 1 output, nm inputs: 1 by nm Jacobian?
 - Inconvenient to do $\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$

업데이트 식 모양이 깔끔하게

- Instead follow convention: shape of the gradient is shape of parameters

- So $\frac{\partial s}{\partial \mathbf{W}}$ is n by m :

$$\begin{bmatrix} \frac{\partial s}{\partial W_{11}} & \cdots & \frac{\partial s}{\partial W_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial s}{\partial W_{n1}} & \cdots & \frac{\partial s}{\partial W_{nm}} \end{bmatrix}$$

Derivative with respect to Matrix

- Remember $\frac{\partial s}{\partial \mathbf{W}} = \delta \frac{\partial \mathbf{z}}{\partial \mathbf{W}}$
- It turns out $\frac{\partial s}{\partial \mathbf{W}} = \delta^T \mathbf{x}^T$

$$\begin{aligned} \frac{\partial s}{\partial \mathbf{W}} &= \delta \frac{\partial \mathbf{z}}{\partial \mathbf{W}} \\ \frac{\partial s}{\partial \mathbf{b}} &= \delta \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \\ \delta &= \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \mathbf{u}^T \circ f'(\mathbf{z}) \end{aligned}$$

$$\frac{\partial s}{\partial \mathbf{W}} = \delta \frac{\partial \mathbf{z}}{\partial \mathbf{W}}$$

$$\frac{\partial s}{\partial \mathbf{b}} = \delta \frac{\partial \mathbf{z}}{\partial \mathbf{b}}$$

$$\delta = \frac{\partial s}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \mathbf{u}^T \circ f'(\mathbf{z})$$

How to Compute?

1. By Hand
2. Automatically → Computational Graph

Computational Graphs

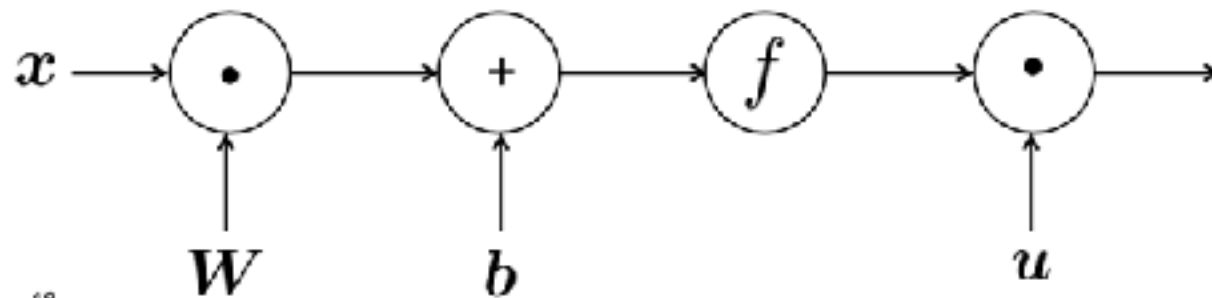
- Representing our neural net equations as a graph
 - Source nodes: inputs
 - Interior nodes: operations

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{z})$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

\mathbf{x} (input)



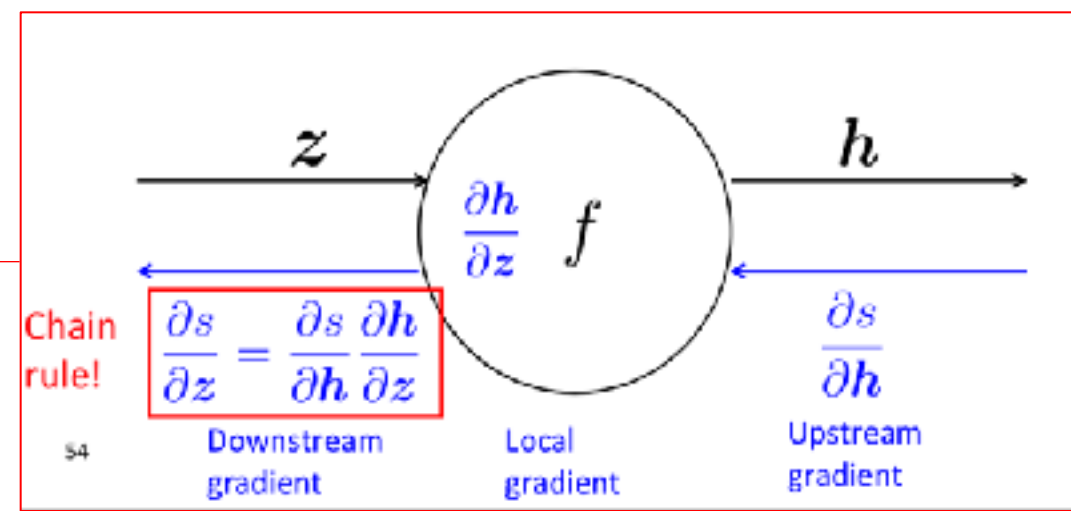
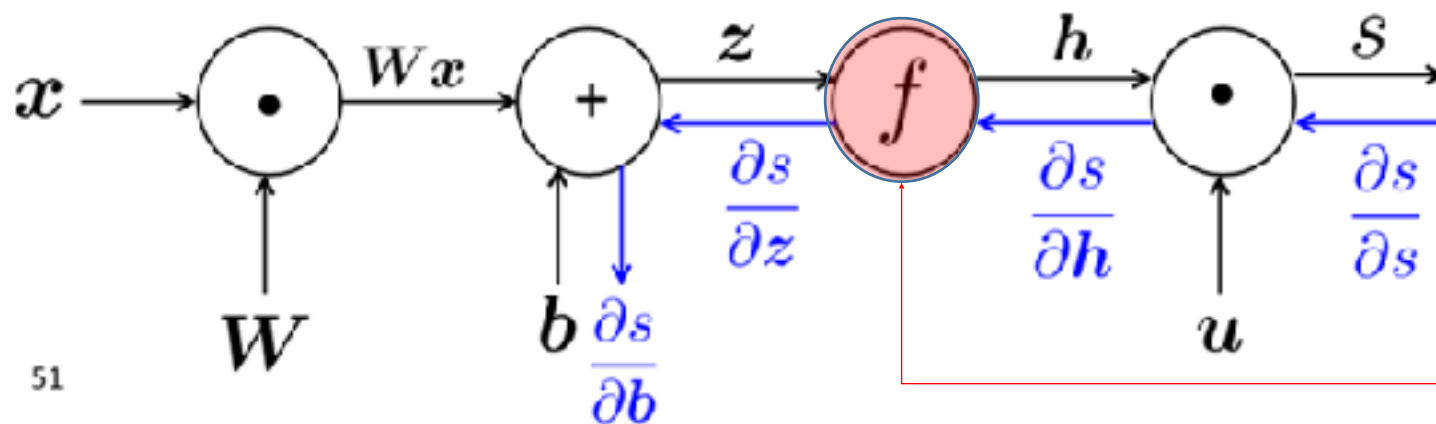
- Go backwards along edges
 - Pass along **gradients**

$$s = u^T h$$

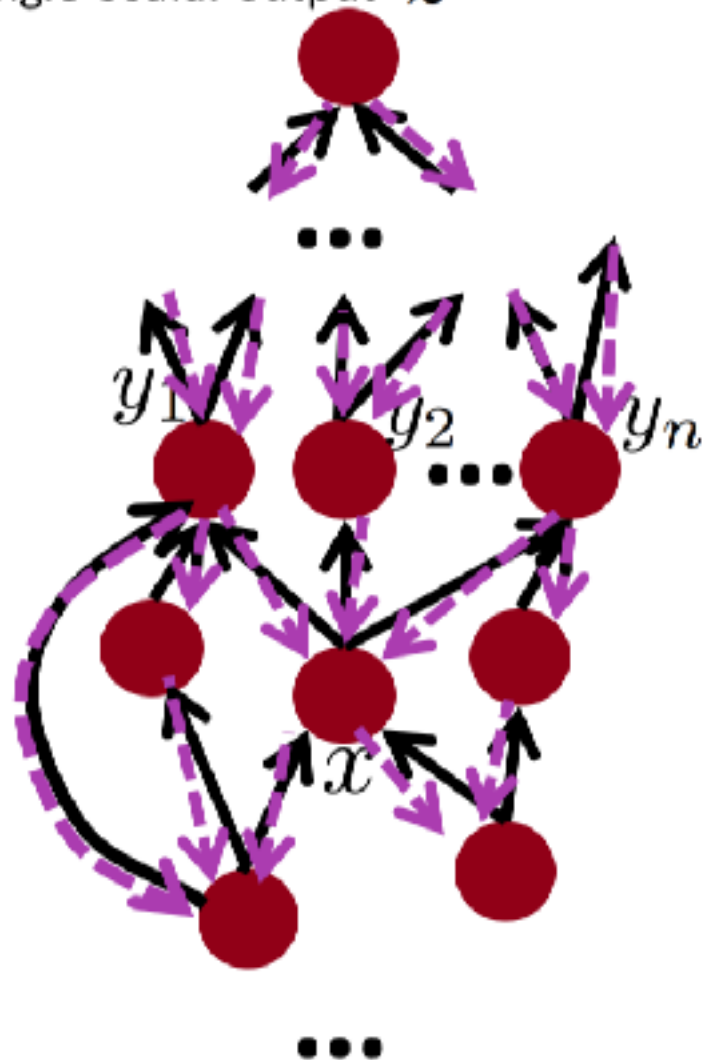
$$h = f(z)$$

$$z = Wx + b$$

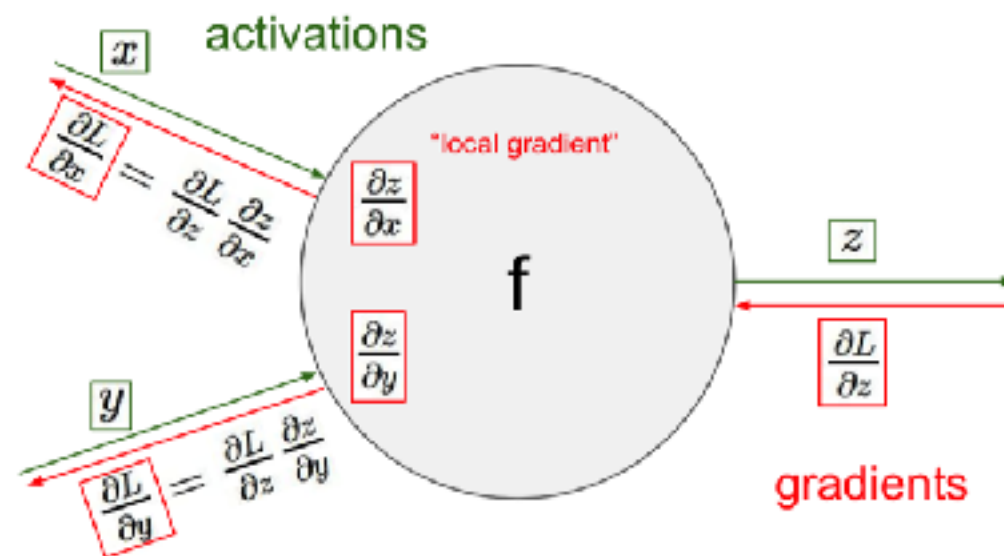
$$x \text{ (input)}$$



Single scalar output z



Recursively apply chain rule through each node



A diagram showing a node f with input x and outputs y_1 and y_2 . To the right of the diagram is the equation for the partial derivative of z with respect to x :

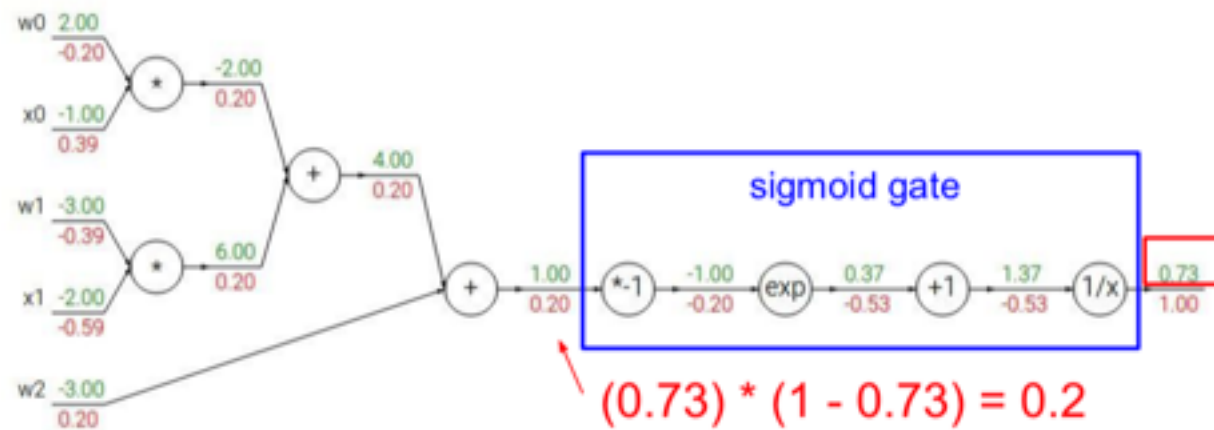
$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2 x_2)}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x)) \sigma(x)$$



Presentation

Thanks
for Watching