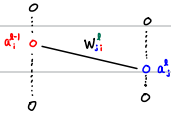


Layer $l-1$ l



$$a^l = \sigma(w^l a^{l-1} + b^l), \quad a^l, b^l: n_l \times 1 \text{ 행렬}, \quad w^l: n_{l-1} \times n_l \text{ 행렬}$$

$$a_j^l \text{의 번째 } b_j: a_j^l = \sigma\left(\sum_{i=1}^{n_{l-1}} w_{ji}^l a_i^{l-1} + b_j^l\right)$$

$\sigma(\cdot)$: 활성화 함수 (activation function)

C : 기준 함수 (criterion function) *) 손실 함수 Loss function 으로 생각해도 무방

[Input 2t Feedforward]

• Input : a^l

• $z^l = w^{l-1} \cdot a^{l-1} + b^l$ and $a^l = \sigma(z^l)$ where $l=2, 3, \dots, L$

[Output error 2t Backpropagate the error]

• Output error : $\delta^l = \frac{\partial C}{\partial z^l}$ 이라 하자.

l 이 $L-1, L-2, \dots, 2$ 에서 $\delta^l = (w^{l+1})^T \delta^{l+1}$ 을 만족한다.

$$\begin{aligned} p3) \quad \delta^{l-1} &= \frac{\partial C}{\partial z^{l-1}} = \left(\frac{\partial C}{\partial z_1^{l-1}}, \frac{\partial C}{\partial z_2^{l-1}}, \dots, \frac{\partial C}{\partial z_{n_{l-1}}^{l-1}} \right)^T \\ &= \left(\frac{\partial C}{\partial a_1^l} \frac{\partial a_1^l}{\partial z_1^{l-1}}, \frac{\partial C}{\partial a_2^l} \frac{\partial a_2^l}{\partial z_1^{l-1}}, \dots \right)^T \\ &= \left(\frac{\partial C}{\partial a_1^l}, \frac{\partial C}{\partial a_2^l}, \dots \right)^T \odot \sigma'(z^{l-1}) \\ &= (w^l)^T \delta^l \odot \sigma'(z^{l-1}) \quad \leftarrow \text{by 1)} \end{aligned}$$

[Weight update]

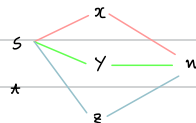
$$\begin{aligned} \frac{\partial C}{\partial w_{ji}^l} &= \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{ji}^l} = \delta_j^l a_i^{l-1} \Rightarrow \frac{\partial C}{\partial w_{ji}^l} = \delta_j^l (a_i^{l-1})^T \\ \frac{\partial C}{\partial b_j^l} &= \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = \delta_j^l \Rightarrow \frac{\partial C}{\partial b_j^l} = \delta_j^l \end{aligned}$$

① i 가 $1, 2, \dots, n_{l-1}$ 일 때

$$\begin{aligned} \frac{\partial C}{\partial a_i^{l-1}} &= \sum_{j=1}^{n_l} \frac{\partial C}{\partial a_j^l} \cdot \frac{\partial a_j^l}{\partial a_i^{l-1}} \quad \dots \text{ by 2)} \\ &= \sum_{j=1}^{n_l} \frac{\partial C}{\partial a_j^l} \cdot \frac{\partial a_j^l}{\partial z_j^l} \frac{\partial z_j^l}{\partial a_i^{l-1}} = \sum_{j=1}^{n_l} \delta_j^l \cdot w_{ji}^l = \left[(w^l)^T \delta^l \right]_i \end{aligned}$$

② Chain rule : 연쇄법칙

Case) $w = f(s, x, z)$, $x = p(s, x)$, $y = q(s, x)$, $z = r(s, x)$



$$\frac{dw}{ds} = \frac{dw}{dx} \cdot \frac{dx}{ds} + \frac{dw}{dy} \cdot \frac{dy}{ds} + \frac{dw}{dz} \cdot \frac{dz}{ds}$$

[문제 1]

input : a^0 와 target : t 가 $a^1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $t = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$ 이고

$$W^1 = \begin{pmatrix} 1 & -2 \\ 2 & 4 \\ -3 & 1 \end{pmatrix}, \quad b^1 = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$$

$$W^3 = \begin{pmatrix} 1 & 2 & -3 \\ 2 & -1 & 3 \end{pmatrix}, \quad b^3 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad \text{activation : } \sigma(x) = x$$

$$\begin{pmatrix} a_1^1 \\ a_2^1 \end{pmatrix} = \sigma \left(W^3 \cdot \sigma \left(W^1 \cdot a^0 + b^1 \right) + b^3 \right) \text{ 일때}$$

$$\text{Loss} = \frac{1}{2} \left[(a_1^1 - t_1)^2 + (a_2^1 - t_2)^2 \right] \text{에 대하여 } lr = 0.01$$

Weight update를 한 결과를 구해보자.

[Forward]

$$\text{Layer: } l \quad \text{forward} \xrightarrow{W^l a + b} z^l \xrightarrow{\sigma(z)} a^l$$

$$1 \quad \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$2 \quad W^1 a^0 + b^1 \rightarrow \begin{pmatrix} 0 \\ 8 \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ 8 \\ 1 \end{pmatrix}$$

$$3 \quad W^3 a^1 + b^3 \rightarrow \begin{pmatrix} 15 \\ -4 \end{pmatrix} \rightarrow \begin{pmatrix} 15 \\ -4 \end{pmatrix}$$

[Backward]

$$\text{Layer: } l \quad \text{backward} \longrightarrow \delta^l$$

$$3 \quad \frac{\partial L}{\partial z^3} \longrightarrow \begin{pmatrix} 13 \\ -7 \end{pmatrix}$$

$$2 \quad ((W^3)^T \cdot \delta^3) \odot \sigma'(z^2) \longrightarrow \begin{pmatrix} -1 \\ 33 \\ -60 \end{pmatrix}$$

[Weight update]

$$\frac{\partial L}{\partial W^3} = \begin{pmatrix} 13 \\ -7 \end{pmatrix} (0 \ 8 \ 1) = \begin{pmatrix} 0 & 104 & 13 \\ 0 & -56 & -7 \end{pmatrix}, \quad \frac{\partial L}{\partial b^3} = \begin{pmatrix} 13 \\ -7 \end{pmatrix}$$

$$W^3 - 0.01 \frac{\partial L}{\partial W^3} = \begin{pmatrix} 1 & 0.96 & -3.13 \\ 2 & -0.94 & 3.07 \end{pmatrix}$$

$$b^3 - 0.01 \frac{\partial L}{\partial b^3} = \begin{pmatrix} 2.99 \\ 1.01 \end{pmatrix}$$

$$\frac{\partial L}{\partial W^1} = \begin{pmatrix} -1 \\ 33 \\ -60 \end{pmatrix} (1 \ 1) = \begin{pmatrix} -1 & -1 \\ 33 & 33 \\ -60 & -60 \end{pmatrix}, \quad \frac{\partial L}{\partial b^1} = \begin{pmatrix} -1 \\ 33 \\ -60 \end{pmatrix}$$

$$W^1 - 0.01 \frac{\partial L}{\partial W^1} = \begin{pmatrix} 1.01 & -1.99 \\ 1.67 & 3.67 \\ -2.4 & 1.6 \end{pmatrix}$$

$$b^1 - 0.01 \frac{\partial L}{\partial b^1} = \begin{pmatrix} 1.99 \\ 1.67 \\ 3.6 \end{pmatrix}$$

[Step 2]

input : a^1 and target : t 가 $a^1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, $t = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$ 이고

$$W^1 = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}, \quad b^1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$W^2 = \begin{pmatrix} 1 & -1 \\ 2 & -1 \end{pmatrix}, \quad b^2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \text{activation : } \sigma(x) = x^2$$

$$\begin{pmatrix} a_1^2 \\ a_2^2 \end{pmatrix} = \sigma \left(W^2 \cdot \sigma \left(W^1 \cdot a^1 + b^1 \right) + b^2 \right) \text{ 일때}$$

$$\text{Loss} = \frac{1}{2} \left[(a_1^2 - t_1)^2 + (a_2^2 - t_2)^2 \right] \text{에 대하여 } |r| = 0.012$$

Weight update 를 한 결과를 구해보자.

[Forward]

$$\text{Layer: } l \quad \text{forward} \xrightarrow{W^l a + b} z^l \xrightarrow{\sigma(z)} a^l$$

$$1 \quad \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$2 \quad W^1 a^1 + b^1 \rightarrow \begin{pmatrix} -1 \\ 2 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 4 \end{pmatrix}$$

$$3 \quad W^2 a^2 + b^2 \rightarrow \begin{pmatrix} -3 \\ -2 \end{pmatrix} \rightarrow \begin{pmatrix} 9 \\ 4 \end{pmatrix}$$

[Backward]

$$\text{Layer: } l \quad \text{backward} \longrightarrow \delta^l$$

$$3 \quad \frac{\partial L}{\partial z^3} \longrightarrow \begin{pmatrix} -42 \\ -4 \end{pmatrix}$$

$$2 \quad ((W^3)^T \cdot \delta^3) \odot \sigma'(z^2) \longrightarrow \begin{pmatrix} 100 \\ 184 \end{pmatrix}$$

$$\frac{\partial L}{\partial z^3} = \frac{\partial L}{\partial a^3} \cdot \frac{\partial a^3}{\partial z^3} = \begin{pmatrix} a_1^3 - t_1 \\ a_2^3 - t_2 \end{pmatrix} \odot \begin{pmatrix} 2 \cdot z_1^3 \\ 2 \cdot z_2^3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \odot \begin{pmatrix} -6 \\ -4 \end{pmatrix} = \begin{pmatrix} -6 \\ -4 \end{pmatrix}$$

$$((W^3)^T \cdot \delta^3) \odot \sigma'(z^2) = \begin{pmatrix} 1 & 2 \\ -1 & -1 \end{pmatrix} \begin{pmatrix} -6 \\ -4 \end{pmatrix} \odot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 100 \\ 184 \end{pmatrix}$$

[Weight update]

$$\frac{\partial L}{\partial W^3} = \begin{pmatrix} -42 \\ -4 \end{pmatrix} \begin{pmatrix} 1 & 4 \end{pmatrix} = \begin{pmatrix} -42 & -168 \\ -4 & -16 \end{pmatrix}, \quad \frac{\partial L}{\partial b^3} = \begin{pmatrix} -42 \\ -4 \end{pmatrix}$$

$$W^3 - 0.01 \frac{\partial L}{\partial W^3} = \begin{pmatrix} 1.42 & 0.68 \\ 2.04 & -0.84 \end{pmatrix}$$

$$b^3 - 0.01 \frac{\partial L}{\partial b^3} = \begin{pmatrix} 0.42 \\ 0.04 \end{pmatrix}$$

$$\frac{\partial L}{\partial W^2} = \begin{pmatrix} 100 \\ 184 \end{pmatrix} \begin{pmatrix} 1 & 1 \end{pmatrix} = \begin{pmatrix} 100 & 100 \\ 184 & 184 \end{pmatrix}, \quad \frac{\partial L}{\partial b^2} = \begin{pmatrix} 100 \\ 184 \end{pmatrix}$$

$$W^2 - 0.01 \frac{\partial L}{\partial W^2} = \begin{pmatrix} 0 & -3 \\ -0.84 & -0.84 \end{pmatrix}$$

$$b^2 - 0.01 \frac{\partial L}{\partial b^2} = \begin{pmatrix} -1 \\ -1.84 \end{pmatrix}$$