



$$a^{[l]} = \sigma(z^{[l]})$$

$$z^{[l]} = w^{[l]} x + b^{[l]}$$

\uparrow \uparrow \uparrow \uparrow
 \mathbb{R}^{n_l} $\mathbb{R}^{n_{l+1}}$ $\mathbb{R}^{n_{l+1}}$ $\mathbb{R}^{n_{l+1}}$

$$\Rightarrow z^{[l]} = w^{[l-1]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = \sigma(z^{[l]}) \quad (*)$$

$$H(x) = a^{[L]}, \quad H: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_L}$$

Define $\delta_j^{[l]} = \frac{\partial H}{\partial z_j^{[l]}}$, $l=2, \dots, L$, $\delta_j^{[1]} = \frac{\partial H}{\partial z_j^{[1]}} = \sigma'(z_j^{[1]})$

$$H(x) = \sigma(z^{[L]}), \quad z^{[L]} \in \mathbb{R}$$

$$\delta_j^{[l-1]} = \frac{\partial H}{\partial z_j^{[l-1]}} \cdot \frac{\partial z_j^{[l-1]}}{\partial z_j^{[l-1]}}$$

* HW 3/1

$n_L=1$

$$H: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^1$$

Calculate $\nabla H = \begin{bmatrix} \frac{\partial H}{\partial x_1} \\ \vdots \\ \frac{\partial H}{\partial x_{n_1}} \end{bmatrix}$

<ans> $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \quad (*)$, $a^{[l]} = \sigma(z^{[l]})$, $H(x) = a^{[L]}$, calculate $\frac{\partial H(x)}{\partial x}$

Let $(MSE) C = \frac{1}{2} \|y - H(x)\|_2^2$, $\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}}$, $1 \leq j \leq n_l$, $2 \leq l \leq L$

$$\frac{\partial a_j^{[l]}}{\partial z_j^{[l]}} = \sigma'(z_j^{[l]}), \quad \frac{\partial C}{\partial a_j^{[l]}} = \frac{\partial}{\partial a_j^{[l]}} \left(\frac{1}{2} \sum_{k=1}^{n_L} (y_k - a_k^{[L]})^2 \right) = -(y_j - a_j^{[L]})$$

using chain rule, $\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \frac{\partial C}{\partial a_j^{[l]}} \frac{\partial a_j^{[l]}}{\partial z_j^{[l]}} = (a_j^{[L]} - y_j) \sigma'(z_j^{[l]}) = \sigma'(z_j^{[l]}) \circ (a_j^{[L]} - y_j)$

$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \sum_{k=1}^{n_L} \frac{\partial C}{\partial a_k^{[L]}} \frac{\partial a_k^{[L]}}{\partial z_j^{[l]}} = \sum_{k=1}^{n_L} \delta_k^{[L]} \frac{\partial a_k^{[L]}}{\partial z_j^{[l]}} = \sum_{k=1}^{n_L} \delta_k^{[L]} w_{kj}^{[L+1]} \sigma'(z_j^{[l]}) = \sigma'(z_j^{[l]}) (w^{[L+1]})^T \delta^{[L+1]}_j = (w^{[L+1]})^T \delta^{[L+1]} \circ \sigma'(z_j^{[l]})$$

By $(*)$, $a_k^{[l+1]} = \sum_{s=1}^{n_l} w_{ks}^{[l+1]} \sigma(z_s^{[l]}) + b_k^{[l+1]}$. Hence, $\frac{\partial a_k^{[l+1]}}{\partial z_j^{[l]}} = w_{kj}^{[l+1]} \sigma'(z_j^{[l]})$

By $(*)$ and $(*)$, $z_j^{[l]} = (w^{[l]} \sigma(z^{[l-1]}))_j + b_j^{[l]} \Rightarrow \frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = 1 \xrightarrow{\text{chain rule}} \frac{\partial C}{\partial b_j^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} \frac{\partial z_j^{[l]}}{\partial b_j^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} = \delta_j^{[l]}$

$z_j^{[l]} = \sum_{k=1}^{n_{l-1}} w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]}$ and $\frac{\partial z_j^{[l]}}{\partial a_k^{[l-1]}} = w_{jk}^{[l]}$, indep. of j , and $\frac{\partial z_j^{[l]}}{\partial a_k^{[l-1]}} = 0$ for $s \neq j$

$$\Rightarrow \frac{\partial C}{\partial w_{jk}^{[l]}} = \sum_{s=1}^{n_L} \frac{\partial C}{\partial a_s^{[L]}} \frac{\partial a_s^{[L]}}{\partial w_{jk}^{[l]}} = \frac{\partial C}{\partial a_j^{[L]}} \frac{\partial a_j^{[L]}}{\partial w_{jk}^{[l]}} = \frac{\partial C}{\partial z_j^{[l]}} a_k^{[l-1]} = \delta_j^{[l]} a_k^{[l-1]}$$

$$\frac{\partial z^{[l]}}{\partial a^{[l-1]}} = w^{[l]}, \quad a^{[l]} = \sigma(z^{[l]}) \Rightarrow \frac{\partial a^{[l]}}{\partial z^{[l]}} = \text{diag}(\sigma'(z^{[l]})) = D^{[l]}$$

layer 1 $\frac{\partial a^{[1]}}{\partial x} = \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial x} = D^{[1]} w^{[1]}$

layer 2 $\frac{\partial a^{[2]}}{\partial x} = \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial x} = D^{[2]} w^{[2]} (D^{[1]} w^{[1]})$

layer l. $\frac{\partial a^{[l]}}{\partial x} = \frac{\partial a^{[l]}}{\partial z^{[l]}} \frac{\partial z^{[l]}}{\partial a^{[l-1]}} \frac{\partial a^{[l-1]}}{\partial x} = D^{[l]} w^{[l]} \left(\frac{\partial a^{[l-1]}}{\partial x} \right)$

$$\frac{\partial H}{\partial x} = \frac{\partial a^{[L]}}{\partial x} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a^{[L-1]}} \dots \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial x}$$

$$\Rightarrow \frac{\partial H}{\partial x} = (D^{[L]} w^{[L]}) (D^{[L-1]} w^{[L-1]}) \dots (D^{[2]} w^{[2]}) (D^{[1]} w^{[1]})$$

The derivation is based on the paper Higham 2019 - DeepLearningIntroAppliedMathematic.