

1. Read [Deep Learning: An Introduction for Applied Mathematicians](#). Consider a network as defined in (3.1) and (3.2). Assume that  $n_L = 1$ , find an algorithm to calculate  $\nabla a^{[L]}(x)$ .

$$\nabla a^{[L]}(x) = \begin{bmatrix} \frac{\partial a^{[L]}}{\partial w_{\lambda_j}^{[L]}} \\ \frac{\partial a^{[L]}}{\partial w_{\lambda_j}^{[L-1]}} \\ \vdots \\ \frac{\partial a^{[L]}}{\partial w_{\lambda_j}^{[2]}} \\ \frac{\partial a^{[L]}}{\partial b_j^{[L]}} \\ \vdots \\ \frac{\partial a^{[L]}}{\partial b_j^{[2]}} \end{bmatrix} \begin{matrix} \rightarrow 1 \leq \lambda \leq n_L, 1 \leq j \leq n_{L-1} \leftarrow \text{将 } n_L \times n_{L-1} \text{ 个数排成一行} \\ \rightarrow 1 \leq \lambda \leq n_{L-1}, 1 \leq j \leq n_{L-2} \leftarrow \text{将 } n_{L-1} \times n_{L-2} \text{ 个数排成一行} \\ \rightarrow 1 \leq \lambda \leq n_2, 1 \leq j \leq n_1 \leftarrow \text{将 } n_2 \times n_1 \text{ 个数排成一行} \\ \rightarrow 1 \leq j \leq n_L \\ \rightarrow 1 \leq j \leq n_2 \end{matrix}$$

Note:  $a^{[k]} = \sigma(z^{[k]})$  where  $z^{[k]} \in \mathbb{R}^{n_k}$   
 $z^{[k]} = W^{[k]} a^{[k-1]} + b^{[k]}$  where  $W^{[k]} \in M_{n_k \times n_{k-1}}$ ,  $a \in \mathbb{R}^{n_{k-1}}$ ,  $b^{[k]} \in \mathbb{R}^{n_k}$

$$\begin{aligned} a^{[L]} &= \sigma(z^{[L]}) = \sigma\left(W^{[L]} a^{[L-1]} + b^{[L]}\right) = \sigma\left(W^{[L]} \sigma(z^{[L-1]}) + b^{[L]}\right) \\ &= \sigma\left(W^{[L]} \sigma\left(W^{[L-1]} a^{[L-2]} + b^{[L-1]}\right) + b^{[L]}\right) \\ &= \sigma\left(W^{[L]} \sigma\left(W^{[L-1]} \sigma(z^{[L-2]}) + b^{[L-1]}\right) + b^{[L]}\right) \end{aligned}$$

$$= \sigma(W^{[L]} \sigma(W^{[L-1]} \sigma(W^{[L-2]} a^{[L-3]} + b^{[L-2]} + b^{[L-1]} + b^{[L]}))$$

$$\because n_L = 1$$

$$\because W^{[L]} \in M_{1 \times n_{L-1}} \quad \frac{\partial a^{[L]}}{\partial w_{ij}^{[L]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \cdot \frac{\partial z^{[L]}}{\partial w_{ij}^{[L]}}$$

$$\frac{\partial a^{[L]}}{\partial w_{ij}^{[L-1]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \frac{\partial a_j^{[L-1]}}{\partial z_{\lambda}^{[L-1]}} \frac{\partial z_{\lambda}^{[L-1]}}{\partial w_{ij}^{[L-1]}}$$

$$\frac{\partial a^{[L]}}{\partial w_{ij}^{[L-2]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \frac{\partial a_j^{[L-1]}}{\partial z_{\lambda}^{[L-1]}} \frac{\partial z_{\lambda}^{[L-1]}}{\partial a_j^{[L-2]}} \frac{\partial a_j^{[L-2]}}{\partial z_{\lambda}^{[L-2]}} \frac{\partial z_{\lambda}^{[L-2]}}{\partial w_{ij}^{[L-2]}}$$

$$\Rightarrow \frac{\partial a^{[L]}}{\partial w_{ij}^{[L-k]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \cdot \dots \cdot \frac{\partial a_j^{[L-k+1]}}{\partial z_{\lambda}^{[L-k+1]}} \frac{\partial z_{\lambda}^{[L-k]}}{\partial w_{ij}^{[L-k]}} \quad \text{for } k=0, \dots, L-2$$

1  
" "

$$W^{[L-k]} \in M_{n_{L-k} \times n_{L-k+1}}$$

$$\frac{\partial a^{[L]}}{\partial b^{[L]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \cdot \frac{\partial z^{[L]}}{\partial b^{[L]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}}$$

$$\frac{\partial a^{[L]}}{\partial b_j^{[L-1]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \frac{\partial a_j^{[L-1]}}{\partial z_j^{[L-1]}} \frac{\partial z_j^{[L-1]}}{\partial b_j^{[L-1]}}$$

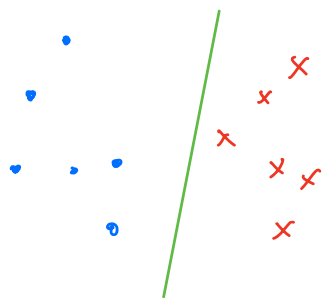
$$\frac{\partial a^{[L]}}{\partial b_j^{[L-2]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \frac{\partial a_j^{[L-1]}}{\partial z_j^{[L-1]}} \cdot \frac{\partial z_j^{[L-1]}}{\partial a_j^{[L-2]}} \frac{\partial a_j^{[L-2]}}{\partial z_j^{[L-2]}}$$

$$\Rightarrow \frac{\partial a^{[L]}}{\partial b_j^{[L-k]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial a_j^{[L-1]}} \cdot \dots \cdot \frac{\partial a_j^{[L-k]}}{\partial z^{[L-k]}} \quad \text{for } k=0, \dots, L-2, \quad b^{[L-k]} \in \mathbb{R}^{n_{L-k}}$$

2. There are unanswered questions during the lecture, and there are likely more questions we haven't covered. Take a moment to think about them and write them down here.

在 classification 問題中,

data :  $\{(\vec{x}_i, C_i)\}$ ,  $C_i \in \{0, 1\}$



方法1: Find a function  $H : \mathbb{R}^2 \rightarrow \mathbb{R}$  s.t.  $H(\vec{x}_i) = C_i$

方法2: (One-hot coding)

Let  $C_i \in \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$

Find a function  $\hat{H} : \mathbb{R}^2 \rightarrow \mathbb{R}$  s.t.  $\hat{H}(\vec{x}_i) = C_i$

Q: 在方法1中找的  $H$  會是 discontinuous function, 在分界會學不好  
為什麼方法2比較好? (可以克服方法1的問題?)