

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  $\nabla a^{[L]}(x)$ .

$$(3.1) \quad a^{[1]} = x \in \mathbb{R}^{n_1},$$

$$(3.2) \quad a^{[l]} = \sigma(W^{[l]} a^{[l-1]} + b^{[l]}) \in \mathbb{R}^{n_l} \quad \text{for } l = 2, 3, \dots, L.$$

Define : ① Input :  $a^{[1]} = x \in \mathbb{R}^{n_1}$

② Hidden layer :  $z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \in \mathbb{R}^{n_l}$ ,  $a^{[l]} = \sigma(z^{[l]}) \in \mathbb{R}^{n_l}$ .  $l = 2, \dots, L-1$

③ Output :  $z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]} \in \mathbb{R}^{n_L}$ ,  $a^{[L]} = \sigma(z^{[L]}) \in \mathbb{R}^{n_L}$

Use backpropagation,

Define loss function  $C = \frac{1}{2} \|y - a^{[L]}\|_2^2$ ,

$$\text{Let } \delta^{[L]} = \frac{\partial C}{\partial z^{[L]}} = \begin{bmatrix} \frac{\partial C}{\partial z_1^{[L]}} \\ \vdots \\ \frac{\partial C}{\partial z_{n_L}^{[L]}} \end{bmatrix}, \quad \frac{\partial a_j^{[L]}}{\partial z_j^{[L]}} = \sigma'(z_j^{[L]}) \quad \text{for } j = 1, 2, \dots, n_L$$

$$\text{Also, } \frac{\partial C}{\partial a_j^{[L]}} = \frac{\partial}{\partial a_j^{[L]}} \frac{1}{2} \sum_{k=1}^{n_L} \|y - a_k^{[L]}\|^2 = -(y - a_j^{[L]}) = a_j^{[L]} - y$$

$$\text{So, } \delta_j^{[L]} = \frac{\partial C}{\partial z_j^{[L]}} = \frac{\partial C}{\partial a_j^{[L]}} \frac{\partial a_j^{[L]}}{\partial z_j^{[L]}} = (a_j^{[L]} - y) \sigma'(z_j^{[L]}) \Rightarrow \underline{\delta^{[L]} = \sigma'(z^{[L]}) \cdot (a^{[L]} - y)}$$

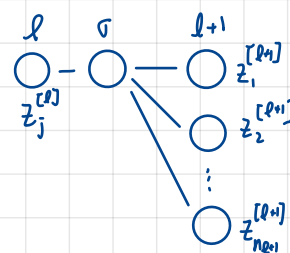
$$\delta_j^{[l]} = \frac{\partial C}{\partial z_j^{[l]}} = \sum_{k=1}^{n_{l+1}} \frac{\partial C}{\partial z_k^{[l+1]}} \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} = \sum_{k=1}^{n_{l+1}} \delta_k^{[l+1]} \frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}}$$

$$\text{Because } z_k^{[l+1]} = \sum_{s=1}^{n_l} w_{ks}^{[l+1]} a_s^{[l]} + b_k^{[l+1]},$$

$$\frac{\partial z_k^{[l+1]}}{\partial z_j^{[l]}} = \frac{\partial}{\partial z_j^{[l]}} \left( \sum_{s=1}^{n_l} w_{ks}^{[l+1]} a_s^{[l]} + b_k^{[l+1]} \right) = w_{kj}^{[l+1]} \sigma'(z_j^{[l]})$$

$$\text{Hence, } \delta_j^{[l]} = \sum_{k=1}^{n_{l+1}} \delta_k^{[l+1]} w_{kj}^{[l+1]} \sigma'(z_j^{[l]}) = \sigma'(z_j^{[l]}) \sum_{k=1}^{n_{l+1}} \delta_k^{[l+1]} w_{kj}^{[l+1]} = \sigma'(z_j^{[l]}) \left( (W^{[l+1]})^T \delta^{[l+1]} \right)_j$$

$$\Rightarrow \underline{\delta^{[l]} = \sigma'(z^{[l]}) \cdot \left( (W^{[l+1]})^T \delta^{[l+1]} \right)}$$



$$\text{Take } n_L = 1, \quad \nabla_x a^{[L]}(x) = \begin{bmatrix} \frac{\partial a^{[L]}}{\partial x_1} \\ \frac{\partial a^{[L]}}{\partial x_2} \\ \vdots \\ \frac{\partial a^{[L]}}{\partial x_{n_1}} \end{bmatrix}, \quad \text{where } \frac{\partial a^{[L]}}{\partial x_k} = \sum_{j=1}^{n_L} \frac{\partial a^{[L]}}{\partial z_j^{[L]}} \frac{\partial z_j^{[L]}}{\partial x_k} \quad \text{for } l = 2, \dots, L-1$$

$$\delta_j^{[2]} = \frac{\partial C}{\partial z_j^{[2]}} = \frac{\partial C}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z_j^{[2]}} = (a^{[2]} - y) \frac{\partial a^{[2]}}{\partial z_j^{[2]}} \Rightarrow \frac{\partial a^{[2]}}{\partial z_j^{[2]}} = \frac{\delta_j^{[2]}}{a^{[2]} - y}$$

$$\because z_j^{[2]} = \sum_{i=1}^{n_1} w_{ji}^{[2]} a_i^{[1]} + b_j^{[2]} = \sum_{i=1}^{n_1} w_{ji}^{[2]} x_i + b_j^{[2]} \Rightarrow \frac{\partial z_j^{[2]}}{\partial x_k} = w_{jk}^{[2]}$$

$$\text{Therefore, } \nabla_x a^{[2]}(x) = (W^{[2]})^T \frac{\delta^{[2]}}{a^{[2]} - y} = \frac{1}{a^{[2]} - y} (W^{[2]})^T \delta^{[2]}$$

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.

In our class, when deriving the relationship between MSE and MLE, we assumed that the error terms follow an i.i.d. Gaussian distribution. However, in real-world data, these assumptions often do not hold.

If the error distribution is skewed (non-symmetric) or heavy-tailed, how would this affect the asymptotic properties of the regression parameter estimators (such as consistency and asymptotic normality)?

In such cases, are there alternative methods (such as robust regression) that can provide better estimates?