

$$1. \nabla a^{[L]}(x) = \left( \frac{\partial a^{[L]}}{\partial x_1}, \dots, \frac{\partial a^{[L]}}{\partial x_{n_1}} \right), \quad z^{[L]} = (x_1, \dots, x_{n_1})$$

$$\text{let } \delta^{[L]} = \frac{\partial a^{[L]}}{\partial z^{[L]}} = \sigma'(z^{[L]})$$

$$\text{then } \delta^{[L]} = \frac{\partial a^{[L]}}{\partial z^{[L]}} = \frac{\partial a^{[L]}}{\partial z^{[L+1]}} \cdot \frac{\partial z^{[L+1]}}{\partial z^{[L]}} = \sigma'(z^{[L]}) \cdot [(W^{[L+1]})^T \delta^{[L+1]}], \quad z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

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

prove it:

when  $l = L-1$ :

$$\delta^{[L-1]} = \frac{\partial a^{[L]}}{\partial z^{[L-1]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \cdot \frac{\partial z^{[L]}}{\partial z^{[L-1]}} = \delta^{[L]} \cdot \left( \frac{\partial z_1^{[L]}}{\partial z_1^{[L-1]}} \dots \frac{\partial z_{n_{L-1}}^{[L]}}{\partial z_{n_{L-1}}^{[L-1]}} \right) = \delta^{[L]} \cdot (w_1^{[L]} \sigma'(z_1^{[L-1]}), \dots, w_{n_{L-1}}^{[L]} \sigma'(z_{n_{L-1}}^{[L-1]}))$$

$$= (\sigma'(z_1^{[L-1]}) [(W^{[L]})^T \delta^{[L]}]_1, \dots, \sigma'(z_{n_{L-1}}^{[L-1]}) [(W^{[L]})^T \delta^{[L]}]_{n_{L-1}})$$

$$= \sigma'(z^{[L-1]}) \cdot [(W^{[L]})^T \delta^{[L]}]$$

Suppose when  $l = k$ :

$$\delta^{[k]} = \sigma'(z^{[k]}) \cdot [(W^{[k+1]})^T \delta^{[k+1]}] \text{ holds}$$

when  $l = k-1$ :

$$\delta^{[k-1]} = \frac{\partial a^{[L]}}{\partial z^{[k-1]}} = \frac{\partial a^{[L]}}{\partial z^{[k]}} \cdot \frac{\partial z^{[k]}}{\partial z^{[k-1]}} = (\delta_1^{[k]}, \dots, \delta_{n_k}^{[k]}) \cdot \begin{bmatrix} \frac{\partial z_1^{[k]}}{\partial z_1^{[k-1]}} & \dots & \frac{\partial z_1^{[k]}}{\partial z_{n_{k-1}}^{[k-1]}} \\ \vdots & & \vdots \\ \frac{\partial z_{n_k}^{[k]}}{\partial z_1^{[k-1]}} & \dots & \frac{\partial z_{n_k}^{[k]}}{\partial z_{n_{k-1}}^{[k-1]}} \end{bmatrix}$$

$$= \left( \sum_{t=1}^{n_k} \delta_t^{[k]} w_{t1}^{[k]} \sigma'(z_1^{[k-1]}), \dots, \sum_{t=1}^{n_k} \delta_t^{[k]} w_{tn_{k-1}}^{[k]} \sigma'(z_{n_{k-1}}^{[k-1]}) \right)$$

$$= (\sigma'(z_1^{[k-1]}) [(W^{[k]})^T \delta^{[k]}]_1, \dots, \sigma'(z_{n_{k-1}}^{[k-1]}) [(W^{[k]})^T \delta^{[k]}]_{n_{k-1}})$$

$$= \sigma'(z^{[k-1]}) \cdot (W^{[k]})^T \delta^{[k]} \text{ also holds.}$$

Thus, for  $l = 1, \dots, L-1$ ,  $\delta^{[l]} = \sigma'(z^{[l]}) \cdot (W^{[l+1]})^T \delta^{[l+1]}$

Algorithm:

① Forward pass:

- $a^{[1]} = x$

- For  $l = 2, \dots, L$ :

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = \sigma(z^{[l]}).$$

save each  $z^{[l]}$  and  $a^{[l]}$  for the backward pass.

② Backward pass:

- Compute  $\delta^{[L]} = \frac{\partial a^{[L]}}{\partial z^{[L]}} = \sigma'(z^{[L]})$

- For  $l = L, L-1, \dots, 3, 2$ :

$$\delta^{[l-1]} = (W^{[l]})^T \delta^{[l]} \cdot \sigma'(z^{[l-1]})$$

- When we reach layer 1 (the input layer), note  $a^{[1]} = x$  has no activation derivative. So

$$\nabla_x a^{[1]}(x) = \delta^{[1]} = (W^{[2]})^T \delta^{[2]}.$$

This produces  $\nabla_x a^{[1]}(x) \in \mathbb{R}^{n_1}$ .