

Backpropagation

Note: $z_j^l = \left(\sum_k w_{jk}^l a_k^{l-1} \right) + b_j^l$

$$a_k^{l-1} = \sigma(z_k^{l-1})$$

To compute ∇L , we must compute:

$$\frac{\partial L}{\partial w_{jk}^l} = \frac{\partial L}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial w_{jk}^l} = d_j^l a_j^{l-1}$$

$$\frac{\partial L}{\partial b_j^l} = \frac{\partial L}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial b_j^l} = d_j^l$$

$$d_j^l = \frac{\partial L}{\partial z_j^l}$$

sometimes
called "error"

Final Layer ($l = F$)

Consider L as a fxn of the nodes in the final layer a_1^F, \dots, a_n^F .

$$a_j^F = \sigma(z_j^F)$$

$$d_j^F = \frac{\partial L}{\partial z_j^F} = \sum_k \frac{\partial L}{\partial a_k^F} \frac{\partial a_k^F}{\partial z_j^F}$$

$$= \frac{\partial L}{\partial a_j^F} \frac{\partial a_j^F}{\partial z_j^F}$$

$$= \frac{\partial L}{\partial a_j^F} \sigma'(z_j^F)$$

$\frac{\partial L}{\partial a_j^F}$ depends on choice of loss fcn
 σ' depends on choice of activation

$$\vec{d}^F = \begin{bmatrix} d_1^F \\ d_2^F \\ \vdots \\ d_n^F \end{bmatrix}$$

n nodes in the final layer

$$\boxed{\vec{d}^F = \nabla_{a^F} L \odot \sigma'(\vec{z}^F)}$$

Intermediate Layers ($l=2, 3, \dots, F-1$)

$$d_j^l = \frac{\partial L}{\partial z_j^l}$$

Question: Where do z_j^l 's appear?

→ In a_j^l 's, which appear in z_j^{l+1} 's

$$d_j^l = \frac{\partial L}{\partial z_j^l} = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

Sum over nodes in $l+1$ layer

$$= \sum_k d_k^{l+1} \cdot \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$z_k^{l+1} = \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1}$$

Sum over nodes in l^{th} layer

$$= \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1}$$

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l)$$

$$\Rightarrow d_j^l = \sum_k w_{kj}^{l+1} d_k^{l+1} \sigma'(z_j^l)$$

$$\Rightarrow \boxed{\vec{d}^l = \left((W^{l+1})^T \vec{d}^{l+1} \right) \odot \sigma'(\vec{z}^l)}$$

#8 from WS 5

$$L_T = \frac{1}{n} \sum_{i=1}^n L_i = \frac{1}{n} \sum_{i=1}^n L(\vec{y}_i, \tilde{y}_i)$$

\vec{y}_i = one training sample

$\tilde{y}_i = \vec{a}_i^F$ is model prediction

$n = \#$ samples in training data

Backprop: Gives us an algorithm for computing

$$\frac{\partial L_i}{\partial w_{jk}^l} \quad \& \quad \frac{\partial L_i}{\partial b_j^l}$$

* For Grad. Desc., we need

$$\frac{\partial L_T}{\partial w_{jk}^l} \quad \& \quad \frac{\partial L_T}{\partial b_j^l}$$

$$L_T = \frac{1}{n} \sum_{i=1}^n L_i$$

$$\frac{\partial L_T}{\partial w_{jk}^l} = \sum_{i=1}^n \frac{\partial L_T}{\partial L_i} \cdot \frac{\partial L_i}{\partial w_{jk}^l}$$

Know from back prop

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial w_{jk}^l}$$

$$\frac{\partial L_T}{\partial b_j^l} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i}{\partial b_j^l}$$

#10 WS 5

#4 JN3

$$L = \frac{1}{2} \|\vec{y} - \vec{\tilde{y}}\|^2$$

$$\nabla_{\vec{a}^F} L = \vec{a}^F - \vec{y}$$

$$\nabla_{\vec{a}^F} L = \begin{bmatrix} \frac{\partial L}{\partial a_1^F} \\ \frac{\partial L}{\partial a_2^F} \\ \vdots \\ \frac{\partial L}{\partial a_k^F} \end{bmatrix} \stackrel{\star}{=} \begin{bmatrix} a_1^F - y_1 \\ a_2^F - y_2 \\ \vdots \\ a_k^F - y_k \end{bmatrix} = \vec{a}^F - \vec{y}$$

$$L(\vec{y}, \vec{\tilde{y}}) = \frac{1}{2} \|\vec{y} - \vec{\tilde{y}}\|^2$$

$$= \frac{1}{2} \|\vec{y} - \vec{a}^F\|^2$$

$$= \frac{1}{2} \left(\sqrt{(y_1 - a_1^F)^2 + (y_2 - a_2^F)^2 + \dots + (y_k - a_k^F)^2} \right)^2$$

$$= \frac{1}{2} \left[(y_1 - a_1^F)^2 + (y_2 - a_2^F)^2 + \dots + (y_k - a_k^F)^2 \right]$$

$$\frac{\partial L}{\partial a_k^F} = \frac{1}{2} \left[(-2) \cdot (y_k - a_k^F) \right]$$

$$\stackrel{\star}{=} \boxed{a_k^F - y_k}$$