

Back propagation

$$W^* = \arg \min_W \sum_{i=1}^n \underline{\underline{\text{Div}(F(x_i, w), \hat{y}_i)}}$$

Gradient descent

→ Computes gradient at w^0

$w^{t+1} \leftarrow w^t - \eta \nabla \text{Div}()$ — Move in -ve direction of gradient



$L=1$



$L=2$

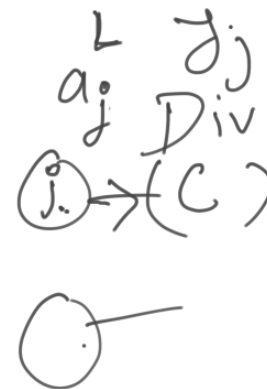
...



$L=l$



$L=L-1$



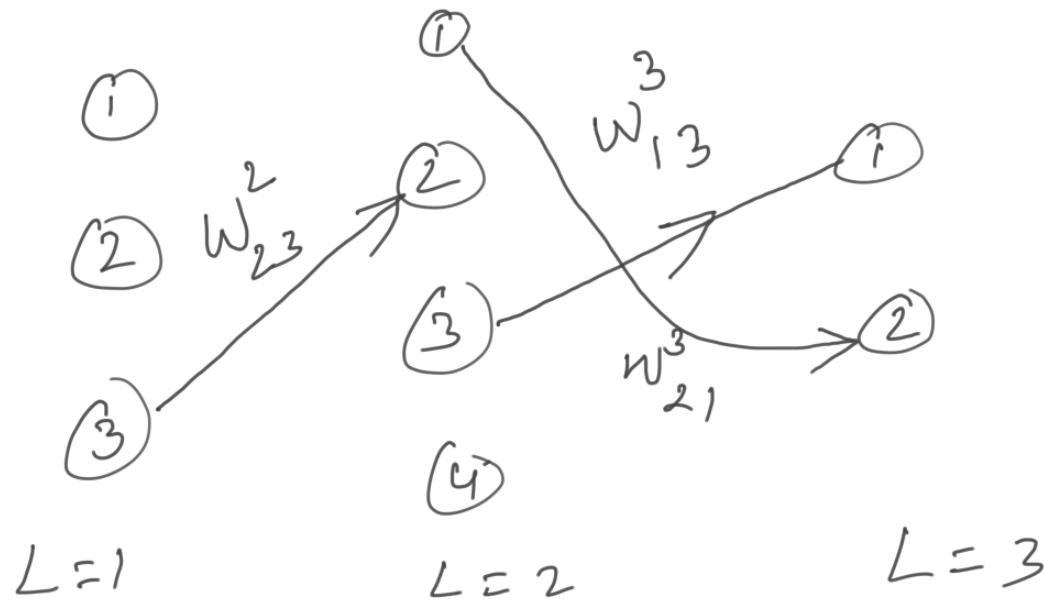
$$C = \frac{1}{2} \| y - a^L \|^2$$

$$= \frac{1}{2} \sum_{j=1}^u (y_j - a_j^L)^2$$

output of the L^{th}
layer = a^L
output has u neurons.

Goal: We want to compute gradient of C
w.r.t. weights.

w_{jk}^l
 $=$ → incoming weight to j^{th} neuron of l^{th} layer from k^{th} neuron of $(l-1)^{\text{th}}$ layer.



$a_j^l \rightarrow$ Activation (output) of j^{th} neuron in

$$a_j^l = \sigma \left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l \right)$$

$$= \sigma(z_j^l)$$

$$z_j^l = \sum_k w_{kj}^l a_k^{l-1} + b_j^l$$

$$\frac{\partial C}{\partial w_{jk}^l}, \frac{\partial C}{\partial b_j^l} \quad l=1, 2, \dots, \textcircled{L}$$

$$\begin{aligned} \frac{\partial C}{\partial w_{jk}^L} &= \frac{\partial}{\partial w_{jk}^L} \left(\frac{1}{2} (y_j - a_j^L)^2 \right) \\ &= \frac{1}{2} \frac{\partial}{\partial a_j^L} (y_j - a_j^L)^2 \cdot \frac{\partial a_j^L}{\partial w_{jk}^L} \end{aligned}$$

$$\left[\frac{\partial}{\partial x} t^2 = \frac{\partial}{\partial t} t^2 \cdot \frac{\partial t}{\partial x} \right]$$

$$= \frac{1}{2} \cancel{2} (y_j - a_j^L) \times (-1) \times \frac{\partial a_j^L}{\partial w_{jk}^L}$$

$$= (a_j^L - y_j) \times \frac{\partial a_j^L}{\partial w_{jk}^L}$$

$$= (a_j^L - y_j) \times \frac{\partial}{\partial w_{jk}^L} \sigma(z_j^L)$$

$$= \underline{(a_j^L - y_j) \sigma'(z_j^L)} \cdot \frac{\partial z_j^L}{\partial w_{jk}^L}$$

$$\frac{\partial \mathcal{L}}{\partial w_{jk}^L} = (a_j^L - y_j) \cdot \sigma'(z_j^L) \frac{\partial}{\partial w_{jk}^L} \left(\sum_k w_{jk}^L a_k^{L-1} + b_j^0 \right)$$

$$\left[z_j^L = \sum_k w_{jk}^L a_k^{L-1} + b_j^0 \right]$$

$$= \boxed{(a_j^L - y_j) \sigma'(z_j^L)} a_k^{L-1}$$

$$= \delta_j^L \times a_k^{L-1}$$

$$\frac{\partial C}{\partial w_{jk}^{L-1}} = \frac{\partial}{\partial w_{jk}^{L-1}} \frac{1}{2} (y_j - a_j^L)^2$$

$$= \frac{\frac{\partial}{\partial a_j^L} \frac{1}{2} (y_j - a_j^L)^2}{\frac{\partial a_j^L}{\partial w_{jk}^{L-1}}}$$

$$= (a_j^L - y_j) \frac{\partial a_j^L}{\partial w_{jk}^{L-1}}$$

$$= (a_j^L - y_j) \frac{\partial}{\partial w_{jk}^{L-1}} \left(\sum_k w_{jk}^{L-1} a_k^{L-1} + b_j^L \right)$$

$$= \frac{(a_j - y_j) \sigma'(z_j^L)}{\frac{\partial}{\partial w_{jk}^{L-1}} \left(\sum_k w_{jk}^L a_k^{L-1} + b_j^L \right)}$$

$$= \delta_j^L \times \frac{\partial}{\partial a_k^{L-1}} \left(\sum_k w_{jk}^L a_k^{L-1} + b_j^L \right) \cdot \frac{\partial a_k^{L-1}}{\partial w_{jk}^{L-1}}$$

$$= \delta_j^L \times w_{jk}^L \times \frac{\partial}{\partial w_{jk}^{L-1}} \sigma \left(\sum_k w_{jk}^{L-1} a_k^{L-2} + b_j^{L-1} \right)$$

$$\frac{\partial C}{\partial w_{jk}^{L-1}} = \delta_j^L \times w_{jk}^L \times \sigma'(z_j^{L-1}) \times a_k^{L-2}$$

$$\frac{\partial C}{\partial w_{jk}^L} = \delta_j^L \times a_k^{L-1}$$

$$\frac{\partial C}{\partial w_{jk}^{L-1}} = \delta_j^{L-1} \times w_{jk}^L \times \sigma'(z_j^{L-1}) \times a_k^{L-2}$$

$$\frac{\partial C}{\partial b_j^L} = \delta_j^L$$

$$\delta_j^l = \delta_j^{l+1} \cdot w_{jk}^{l+1}$$

(l ≠ L)

$$\frac{\partial C}{\partial w_{jk}^{L-1}} = (\sigma_j^L) \times w_{jk}^L \times \sigma'(z_j^{L-1}) \times a_k^{L-2}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \boxed{\sigma_j^{l+1} \times w_{jk}^{l+1}} \times \sigma'(z_j^l) \times a_k^{l-1}$$

$$(l \neq L) \quad \parallel \quad D_j^l \times \sigma'(z_j^l) \times a_k^{l-1}$$

where $\boxed{D_j^l} = \sigma_j^{l+1} \times w_{jk}^{l+1}$

$$\frac{\partial C}{\partial w_{jk}^l} = \delta_j^l \times a_k^{l-1}$$

where $\delta_j^l = \delta_j^{l+1} \times w_{jk}^{l+1}$ } non final layer

$$\frac{\partial C}{\partial w_{jk}^L} = (a_j^L - y_j) \sigma'(z_j^L) \}$$

final layer

$$a_j^l = \sigma \left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l \right)$$

$$= \sigma \left(\underbrace{\begin{bmatrix} \quad \end{bmatrix}_{m \times n}^T}_{\substack{\swarrow \\ \text{}}}} \begin{bmatrix} \quad \end{bmatrix}_{n \times 1} + \begin{bmatrix} \quad \end{bmatrix}_{n \times 1} \right)$$

$l-1$ layer has m neuron

l layer has n neuron

Forward propagation



getting output

Backward pass

$$\frac{\partial C}{\partial w} \quad \frac{\partial C}{\partial b}$$

Computes
gradient

y_1
 y_2 } ground truth

- ① Forward pass — output
- ② Backward pass — gradient
- ③ gradient descent — updates the weight

