

$b_j^l$ : bias of the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

$a_j^l$ : activation of the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

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

↑  
sum over all nodes in  $l-1$

$$a^l = \sigma \left( [a_1^{l-1}, a_2^{l-1}, \dots, a_k^{l-1}] \cdot \begin{bmatrix} w_{11}^l & w_{12}^l & \dots & w_{1k}^l \\ w_{21}^l & w_{22}^l & \dots & w_{2k}^l \\ \vdots & \vdots & \ddots & \vdots \\ w_{j1}^l & w_{j2}^l & \dots & w_{jk}^l \end{bmatrix} + [b_1^l, b_2^l, \dots, b_j^l] \right)$$

dot product

$$a^l = \sigma \left( [a_1^{l-1}, a_2^{l-1}, \dots, a_k^{l-1}, 1] \cdot \begin{bmatrix} w_{11}^l & w_{12}^l & \dots & w_{1k}^l \\ w_{21}^l & w_{22}^l & \dots & w_{2k}^l \\ \vdots & \vdots & \ddots & \vdots \\ w_{j1}^l & w_{j2}^l & \dots & w_{jk}^l \\ b_1^l & b_2^l & \dots & b_j^l \end{bmatrix} \right)$$

↓  
element-wise activation function

$$z_j^l = a^{l-1} \cdot w_j^l$$

$$a^l = \sigma \left( \underbrace{a^{l-1} \cdot w^l}_{z^l} \right) = \sigma(z^l)$$

$z^l$ : weighted input to the neurons in layer  $l$ .

Goal of backpropagation:  $\frac{\partial C}{\partial W}$  and  $\frac{\partial C}{\partial b}$

→ partial derivative of the cost function  $C$  w.r.t <sup>any</sup>  $W$  and  $b$  in the network

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$$

$n$ : total # of training examples.

$\sum_x$ : sum over all training data points.  $x$ .

$y(x)$ : Corresponding output.

$L$ : # of layers in the network

$a^L(x)$ : vector output of the network for input  $x$

Assume: calculating the cost for a single data point  $x$  for now.  
⇒  $x$  becomes a constant, so does  $y(x)$ .

Assume: Cost ( $C$ ) is a function of  $a^L$ . i.e.,  $C = C(a^L)$

i.e., 
$$C = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$

Assume: element-wise multiplication.  $\begin{bmatrix} 1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \end{bmatrix}$

Backpropagation: understand how changing weights & biases in a network changes the cost function.  
i.e.,  $\frac{\partial C}{\partial W_{jk}^l}$ ,  $\frac{\partial C}{\partial b_j^l}$

$\delta_j^l$ : error in the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer.  
 $\delta_j^l \rightarrow$  relate  $\delta_j^l$  to  $W_{jk}^l$  &  $b_j^l$ ,  
 $\frac{\partial C}{\partial W_{jk}^l}$  &  $\frac{\partial C}{\partial b_j^l}$

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}, \quad \delta^l = \frac{\partial C}{\partial z^l} \quad (\text{vectorized})$$

Error in the output layer,  $\delta^L$

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

↑  
component form

$z_j^L$ : compute from  $xW + b$

$\sigma'(z_j^L)$ : compute by taking the derivative of  $\sigma(z_j^L)$

$$\frac{\partial C}{\partial a_j^L} = \text{Take derivative of } C = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$

$$\hookrightarrow \frac{\partial C}{\partial a_j^L} = (a_j^L - y_j)$$

Matrix form:

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \quad (1)$$

~~Error in terms of the~~  
 Error  $\delta^l$  in terms of the error in the next layer  $\delta^{l+1}$  =

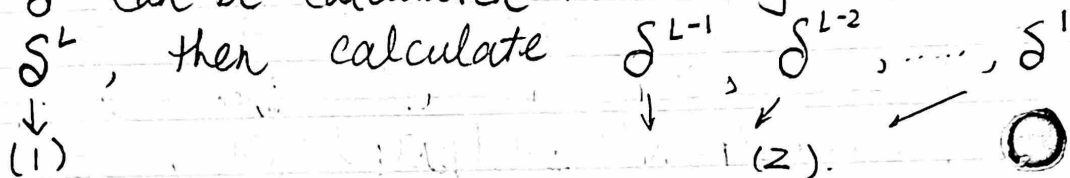
$$\boxed{\delta^l = ((W^{l+1})^T \cdot \delta^{l+1}) \odot \phi'(z^l)} \quad (2)$$

Idea:

$(W^{l+1})^T \delta^{l+1}$  : propagate the error from  $l+1$  backwards to  $l$ .

$\odot \phi'(z^l)$  : pass the error in  $l$  into the activation func.

w/ (1) & (2),  $\delta^l$  can be calculated in all layers  $l$ .

First get  $\delta^L$ , then calculate  $\delta^{L-1}, \delta^{L-2}, \dots, \delta^1$   


Rate of Change in Cost w.r.t any bias in the network.

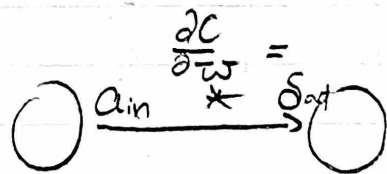
$$\boxed{\frac{\partial C}{\partial b_j^l} = \delta_j^l} \quad (3) \quad \text{or} \quad \frac{\partial C}{\partial b} = \delta$$

for the same neuron in the same layer

Rate of change in Cost w.r.t any weight in the network

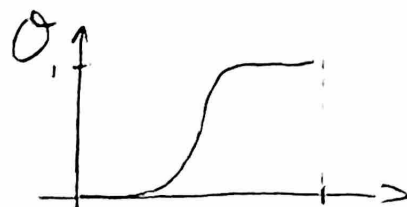
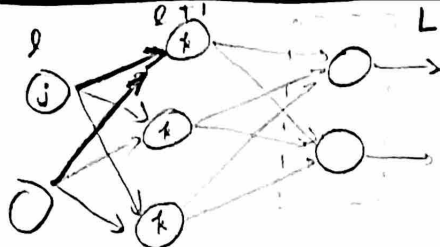
$$\boxed{\frac{\partial C}{\partial W_{jk}^l} = a_k^{l-1} \delta_j^l} \quad (4) \quad \text{or}$$

$$\frac{\partial C}{\partial W} = a_{in} \delta_{out}$$



i.e.,  $\frac{\partial C}{\partial W} = (\text{activation input to } W) * (\text{error output of } W)$

meaning, Weights output from low-activation neurons learn slowly.



$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L) \quad (1)$$

$$\delta_j^L = ((W^{L+1})^T \delta^{L+1}) \odot \sigma'(z_j^L) \quad (2)$$

$$\frac{\partial C}{\partial b_j^L} = \delta_j^L \quad (3)$$

$$\frac{\partial C}{\partial w_{jk}^L} = a_k^{L-1} \delta_j^L \quad (4)$$

(1) (2):  $\sigma'(z_j^L) \approx 0$  when  $\sigma(z_j^L) \approx 1$  or  $\sigma(z_j^L) \approx 0$   
 weight in the final layer will learn slowly if the output neuron has low activation  $\sigma(z_j^L) \approx 0 \rightarrow$  saturated.  
 The error  $\delta_j^L$  will be small when the neuron is saturated

(2) proof

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} \quad (\text{definition})$$

$$\delta_j^L = \sum_k \left( \frac{\partial C}{\partial z_k^{L+1}} \right) \frac{\partial z_k^{L+1}}{\partial z_j^L}$$



$$\rightarrow = \sum_k \delta_k^{L+1} \frac{\partial z_k^{L+1}}{\partial z_j^L} = \sum_k \delta_k^{L+1} w_{kj}^{L+1} \sigma'(z_j^L)$$

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

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

after differentiating, all terms not equal to 0  $\rightarrow 0$

fit-partial, back propagation:

$$\text{error} = \text{layer-output}[-1] - y = a_0^3 - y \left( = \frac{\partial C}{\partial a_0^3} = \frac{\partial C}{\partial a_0^3} = \frac{\partial C}{\partial a_0^3} \right)$$

$$\text{layer-gradient} = \begin{bmatrix} \delta^3 = \text{error} * \text{sigmoid-derivative}(\text{layer-output}[-2].\text{dot}(\text{self.W}[-1])), \\ \delta^2 = \delta^3 (W^3)^T \odot \text{sigmoid-derivative}(\text{layer-output}[-1].\text{dot}(\text{self.W}[-1])), \\ \delta^1 = \delta^2 (W^2)^T \odot \text{sigmoid-derivative}(\text{layer-output}[0].\text{dot}(\text{self.W}[0])) \end{bmatrix}$$

$$\delta^L = \delta^3 = \delta_0^3 = \frac{\partial C}{\partial a_0^3} \cdot \sigma'(z_0^3)$$

$$= \text{error} * \text{sigmoid-derivative}([a_0^2, a_1^2] \cdot \begin{bmatrix} W_{00}^3 \\ W_{01}^3 \end{bmatrix})$$

$$= \text{error} * \text{sigmoid-derivative}(\text{layer-output}[-2].\text{dot}(\text{self.W}[-1])) \rightarrow \text{scalar}$$

for layer in (2, 1)

• layer = 2

$$\text{current-gradient} (= \delta^2 = ((W^3)^T \delta^3) \odot \sigma'(z^2))$$

$$= (\text{layer-gradient}[-1]).\text{dot}(\text{self.W}[\text{layer}].T) * \text{sigmoid-derivative}([a_0^2, a_1^2, a_2^2] \cdot \begin{bmatrix} W_{00}^2 & W_{10}^2 \\ W_{01}^2 & W_{11}^2 \\ W_{02}^2 & W_{12}^2 \end{bmatrix})$$

$$= (\text{layer-gradient}[-1]).\text{dot}(\text{self.W}[\text{layer}].T) * \text{sigmoid-derivative}(\text{layer-output}[\text{layer}-1].\text{dot}(\text{self.W}[\text{layer}-1])) \rightarrow 1 \times 2$$

• layer = 1

$$\text{current-gradient} (= \delta^1 = ((W^2)^T \delta^2) \odot \sigma'(z^1))$$

$$= (\text{layer-gradient}[-1]).\text{dot}(\text{self.W}[\text{layer}].T) * \text{sigmoid-derivative}([x_0^0, x_1^0, x_2^0] \cdot \begin{bmatrix} W_{00}^1 & W_{10}^1 & W_{20}^1 \\ W_{01}^1 & W_{11}^1 & W_{21}^1 \\ W_{02}^1 & W_{12}^1 & W_{22}^1 \end{bmatrix})$$

$$= (\text{layer-gradient}[-1]).\text{dot}(\text{self.W}[\text{layer}].T) * \text{sigmoid-derivative}(\text{layer-output}[\text{layer}-1].\text{dot}(\text{self.W}[\text{layer}-1]))$$

Weight update:

$$W^1: \begin{bmatrix} W_{00}^1 \\ W_{01}^1 \\ W_{02}^1 \end{bmatrix} = \begin{bmatrix} W_{00}^1 \\ W_{01}^1 \\ W_{02}^1 \end{bmatrix} - \alpha \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \delta_0^1 \rightarrow \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} [\delta_0^1 \delta_1^1 \delta_2^1]$$

$$\begin{bmatrix} W_{10}^1 \\ W_{11}^1 \\ W_{12}^1 \end{bmatrix} = \begin{bmatrix} W_{10}^1 \\ W_{11}^1 \\ W_{12}^1 \end{bmatrix} - \alpha \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \delta_1^1$$

$$\begin{bmatrix} W_{20}^1 \\ W_{21}^1 \\ W_{22}^1 \end{bmatrix} = \begin{bmatrix} W_{20}^1 \\ W_{21}^1 \\ W_{22}^1 \end{bmatrix} - \alpha \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \delta_2^1$$

$$W^2: \begin{bmatrix} W_{00}^2 \\ W_{01}^2 \\ W_{02}^2 \end{bmatrix} = \begin{bmatrix} W_{00}^2 \\ W_{01}^2 \\ W_{02}^2 \end{bmatrix} - \alpha \begin{bmatrix} a_0^1 \\ a_1^1 \\ a_2^1 \end{bmatrix} \delta_0^2 \quad \begin{bmatrix} a_0^1 \\ a_1^1 \\ a_2^1 \end{bmatrix} [\delta_0^2 \delta_1^2]$$

$$\begin{bmatrix} W_{10}^2 \\ W_{11}^2 \\ W_{12}^2 \end{bmatrix} = \begin{bmatrix} W_{10}^2 \\ W_{11}^2 \\ W_{12}^2 \end{bmatrix} - \alpha \begin{bmatrix} a_0^1 \\ a_1^1 \\ a_2^1 \end{bmatrix} \delta_1^2$$

$$W^3: \begin{bmatrix} W_{00}^3 \\ W_{01}^3 \end{bmatrix} = \begin{bmatrix} W_{00}^3 \\ W_{01}^3 \end{bmatrix} - \alpha \begin{bmatrix} a_0^2 \\ a_1^2 \end{bmatrix} \delta_0^3 \quad \begin{bmatrix} a_0^2 \\ a_1^2 \end{bmatrix} [\delta_0^3]$$

$$W^3 = W^3 - \alpha * \text{layeroutput}[2].T \cdot \text{dot}(\text{layer-gradient}[2])$$

$$W^2 = W^2 - \alpha * \text{layeroutput}[1].T \cdot \text{dot}(\text{layer-gradient}[1])$$

$$W^1 = W^1 - \alpha * \text{layeroutput}[0].T \cdot \text{dot}(\text{layer-gradient}[0])$$