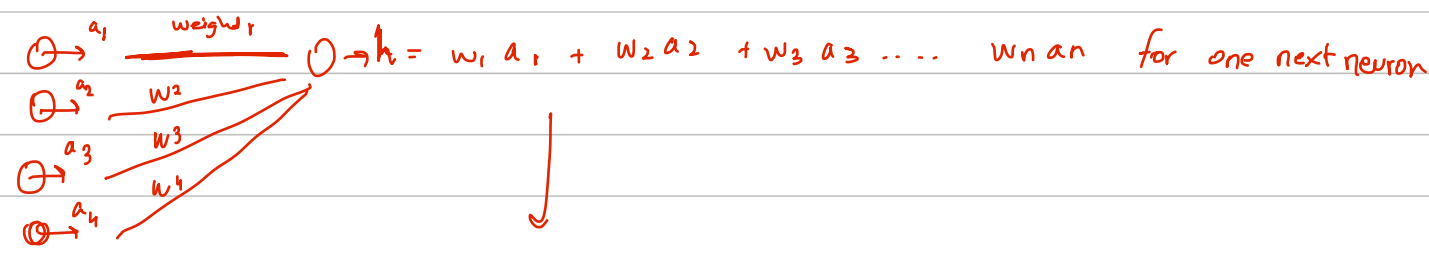


- MNIST  $\rightarrow$  784 neurons  $\rightarrow$  28x28 pixels  $\rightarrow$  28  $\times$  28  $\Rightarrow$  784
- final layer classification  $\rightarrow$  10 neurons for 10 digit classification
- hidden layers can be quite variable

$\hookrightarrow$  a combination of subcomponents  
 $\hookrightarrow$  comprised of subcomponents  
 pixels  
 edges  
 patterns  
 digits

activation; an output of neuron circle/node

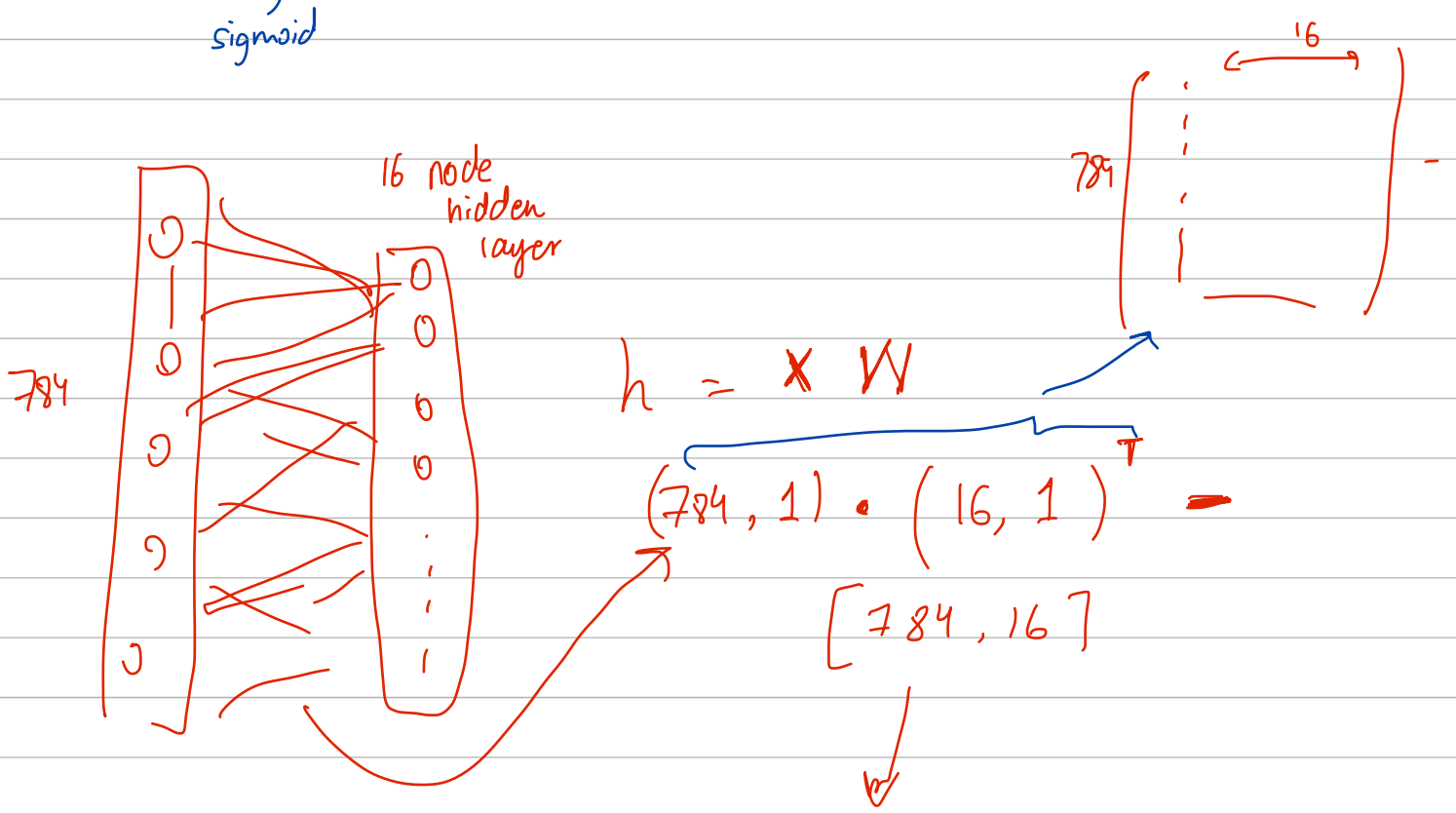


put  $h$  through an activation function to squeeze all values between 0 and 1

Sigmoid  $\sigma(x) = \frac{1}{1 + e^{-x}}$   $0 \leq \sigma(x) \leq 1$

bias: only activate when weighted sum  $>$  threshold

$\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 + \dots + w_n a_n + \text{bias})$   
 sigmoid



$$\sigma([784, 16])$$

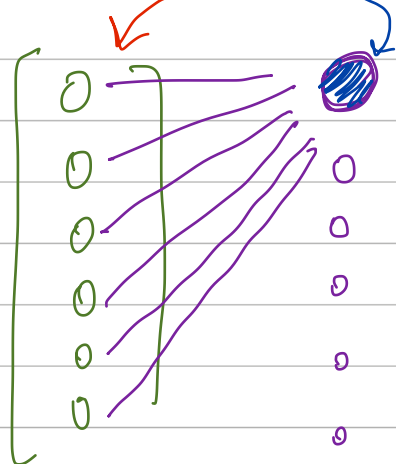
$$\Rightarrow [16, 1] - [16, 1]$$

activations, biases

Better:

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} = \begin{bmatrix} \\ \\ \\ \end{bmatrix}$$

all the <sup>weights</sup> connections between  $a_0 \dots a_n$  and one node in hidden layer.

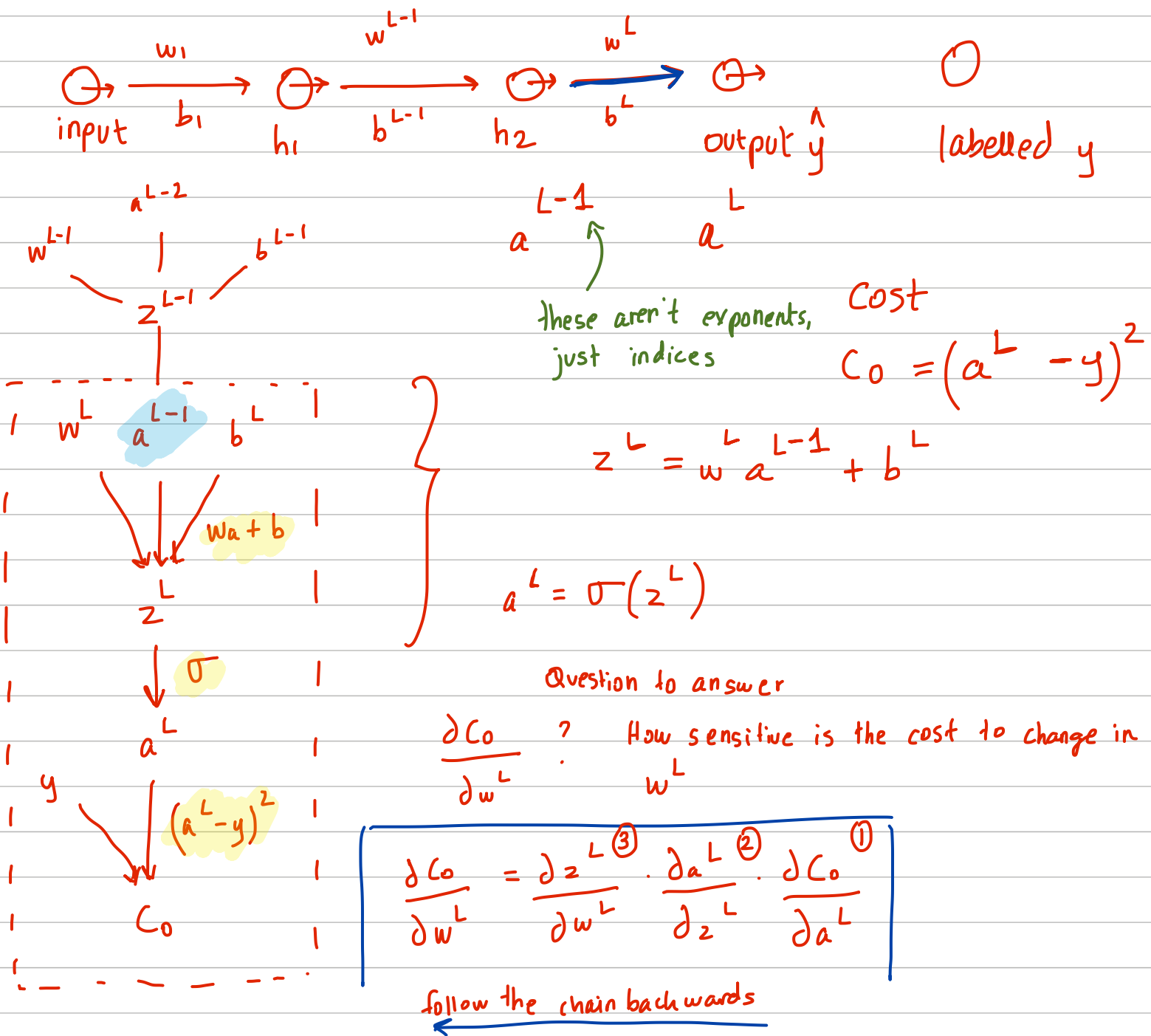


**Backprop:** for a single final neuron, fix weights  
**Gradient descent:** try to minimize cost  
 of all parameters  
 (all weights across all layers)

$$\downarrow C(w)$$

# Calculus :

Single node layers + single training example for cost function



$$C_0 = (a^L - y)^2 \quad (1)$$

$$\frac{\partial C_0}{\partial a^L} = 2(a^L - y)(1)$$

$$= 2(a^L - y)$$

$$(2) \quad a^L = \sigma(z^L) \quad \text{deriv. of sigmoid/relu etc}$$

$$\frac{\partial a^L}{\partial z^L} = \sigma'(z^L)$$

(3)

$$z^L = w^L a^{L-1} + b^L$$

$$\frac{\partial z^L}{\partial w^L} = a^{L-1}$$

the previous activation determines how strongly the weight nudges the layer

not exponents, just indices

Same neural net structure but all training examples

$$\frac{\partial C}{\partial w^L} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^L} \quad k == \text{the } k^{\text{th}} \text{ training example}$$

"stochastic"  
split examples  
into batches

But we also want to minimize cost for all weights,  
not just last layer

Gradient descent:

$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w^1} \\ \frac{\partial C}{\partial b^1} \\ \vdots \\ \frac{\partial C}{\partial w^L} \\ \frac{\partial C}{\partial b^L} \end{bmatrix}$$

gradient of the cost function

$w =$  new weights & biases

all the "nudges"

all weights & biases

gradient  $\nabla C$

Similarly, for a single example

stays same

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial z^L}{\partial b^L} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial C_0}{\partial a^L}$$

work backwards

$$= (1) \cdot \sigma'(z^L) \cdot 2(a^L - y)$$

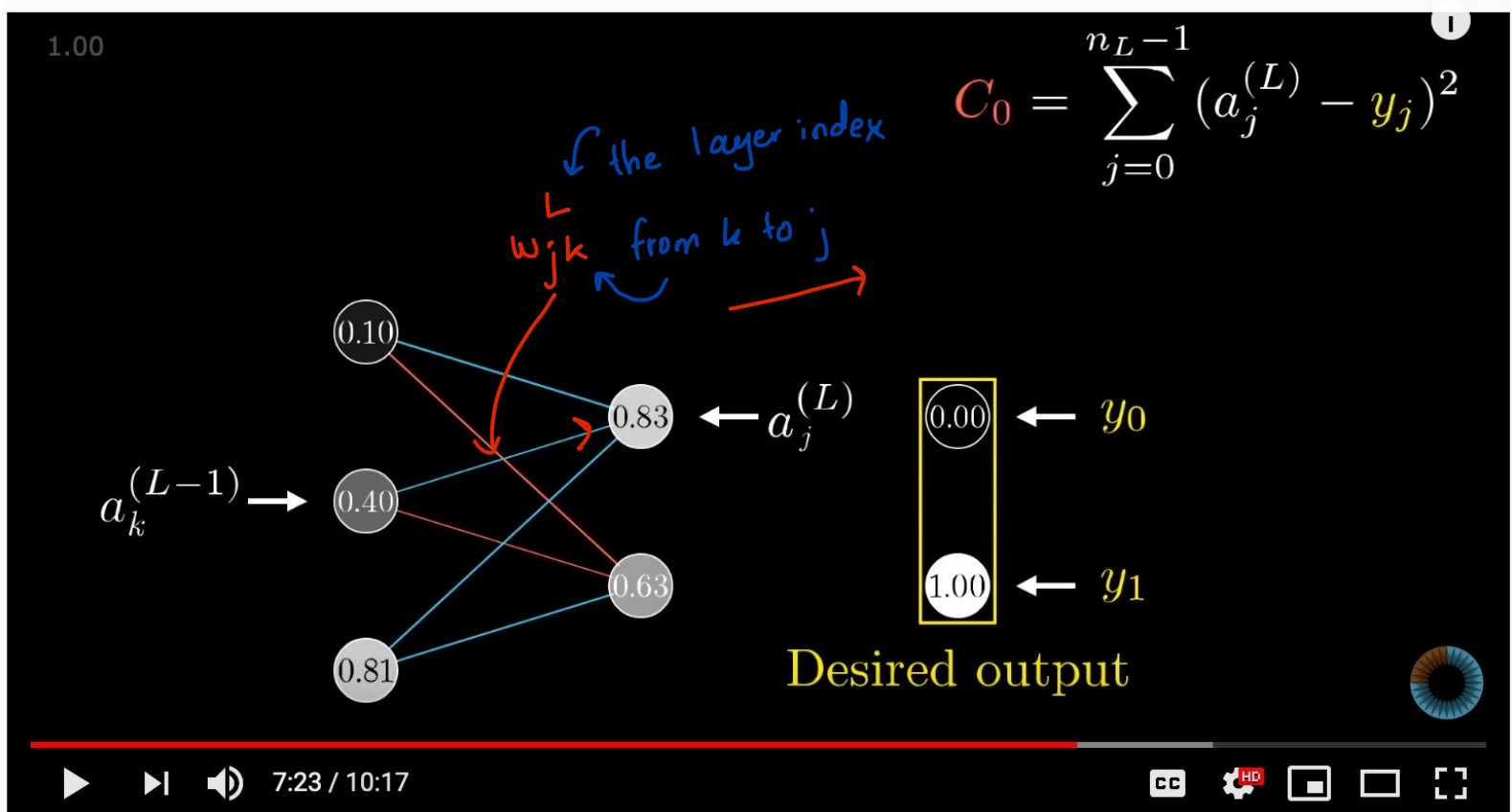
the changed part

Similarly for a single example stays same

$$\frac{dC_0}{da^{L-1}} = \frac{dz^L}{da^{L-1}} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial C_0}{\partial a^L}$$

$$= \underbrace{w^L}_{\text{the changed part}} \cdot \sigma'(z^L) \cdot 2(a^L - y)$$

keep going until you hit  $w^1$  (the backprop part)

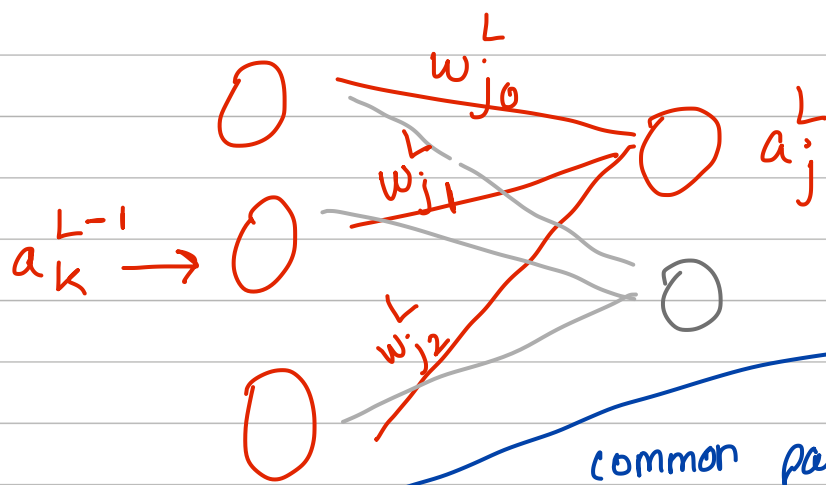


# Multi layer network with one training example

$$C_0 = \sum_{j=0}^{n_{L-1}} (a_j^L - y_j)^2$$

for all  $j$  in the output layer

$$z_j^L = \underbrace{w_{j0}^L a_0^{L-1}} + \underbrace{w_{j1}^L a_1^{L-1}} + \underbrace{w_{j2}^L a_2^{L-1}} + \underbrace{b_j^L}$$



$$\frac{\partial C_0}{\partial w_{jk}^L} = \underbrace{\frac{\partial C_0}{\partial a_j^L}}_{\text{common part}} \times \underbrace{\frac{\partial a_j^L}{\partial z_j^L}}_{\text{variable part}} \times \underbrace{\frac{\partial z_j^L}{\partial w_{jk}^L}}$$

change in cost for training example 0 =  $\frac{\text{change in cost for training e.g. 0}}{\text{change in activation of } j \text{ node in } L \text{ layer}} \times \frac{\text{change in activation of } j \text{ node in } L \text{ layer}}{\text{change in linear combination (i.e. pre-squishing with } \sigma \text{) of node } j \text{ in } L \text{ layer}} \times \frac{\text{(you get it!)}}{\text{the weight}}$

change in weight connecting node  $k$  in  $L-1$  layer to node  $j$  in  $L$  layer

very similar setup for

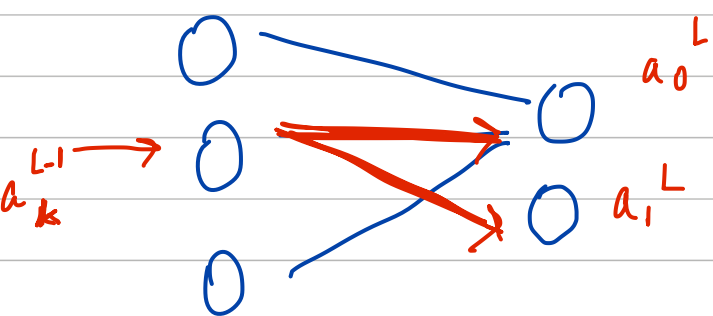
$$\frac{\partial C_0}{\partial b_j^L}$$

and for  $\frac{\partial C_0}{\partial a_k^{L-1}}$

$$\frac{\partial C_0}{\partial a_k^{L-1}} = \sum_{j=0}^{n_{L-1}} \underbrace{\frac{\partial z_j^L}{\partial a_k^{L-1}}}_{\text{changed part}} \underbrace{\frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial C_0}{\partial a_j^L}}_{\text{same as before}}$$

$$\frac{\partial C_0}{\partial b_j^L} = \underbrace{\frac{\partial z_j^L}{\partial b_j^L}}_{\text{changed part}} \underbrace{\frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial C_0}{\partial a_j^L}}_{\text{same as before}}$$

→ why the sum?



the previous activation  $a_k^{L-1}$  not only influences  $a_0^L$  but also every other node in layer  $L$  like  $a_1^L$

1.00

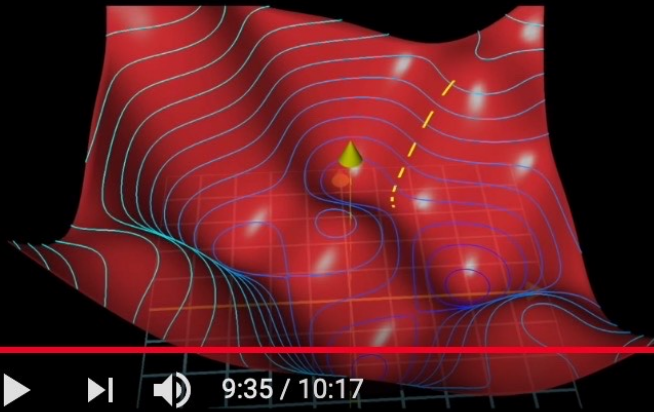
$$\frac{\partial C}{\partial w_{jk}^{(l)}} = a_k^{(l-1)} \sigma'(z_j^{(l)}) \boxed{\frac{\partial C}{\partial a_j^{(l)}}}$$

 $\nabla C$ 

$$\sum_{j=0}^{n_{l+1}-1} w_{jk}^{(l+1)} \sigma'(z_j^{(l+1)}) \frac{\partial C}{\partial a_j^{(l+1)}}$$

or

$$2(a_j^{(L)} - y_j)$$



9:35 / 10:17

