



ML: DL Review Session 2

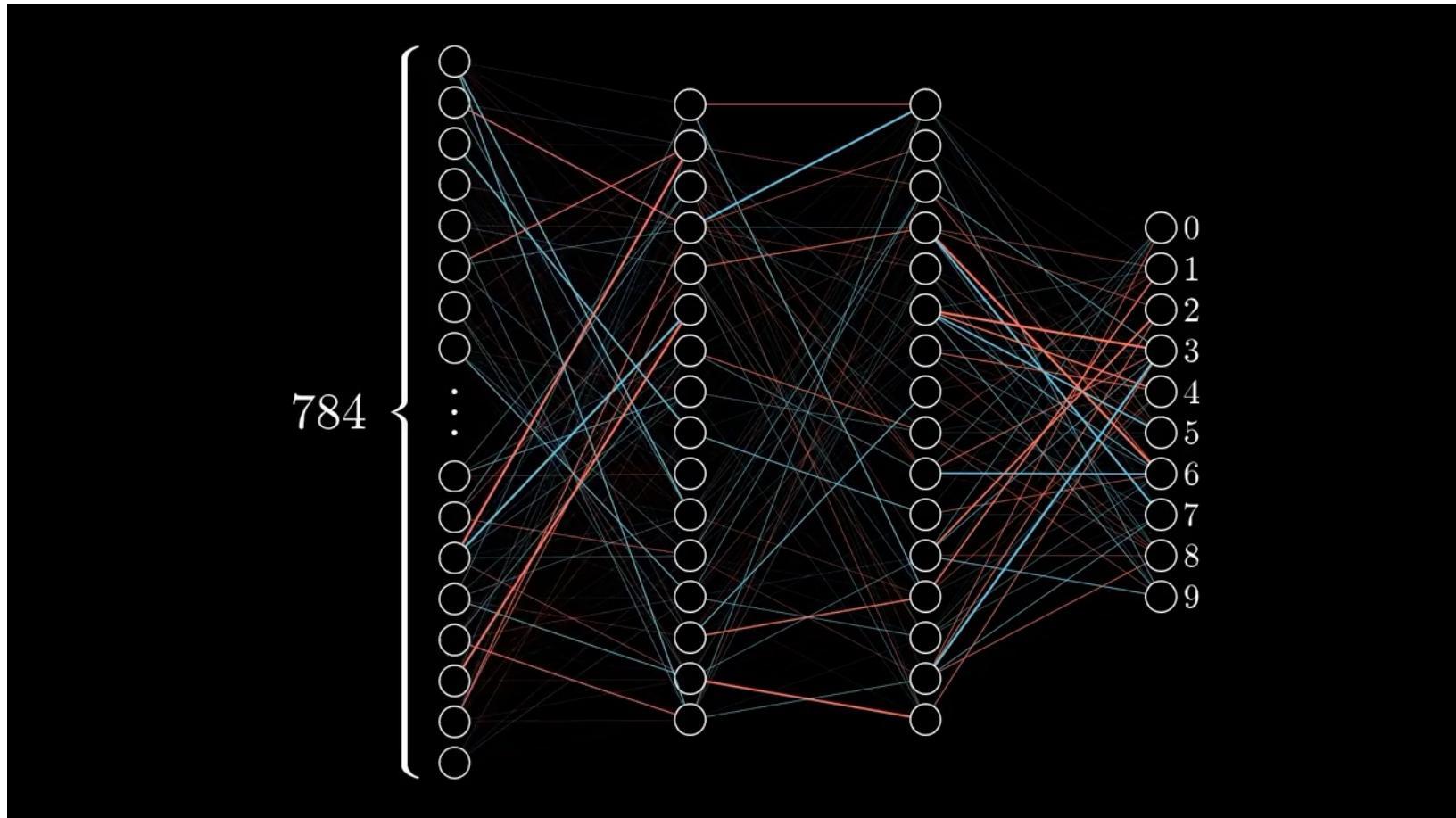
Agenda

- Computational Graph / Back-prop
 - Visualization for understanding
 - Practice
- Intro to NNs
 - Activation function
 - MLP walkthrough with google colab



Computational graph of a Neural Net

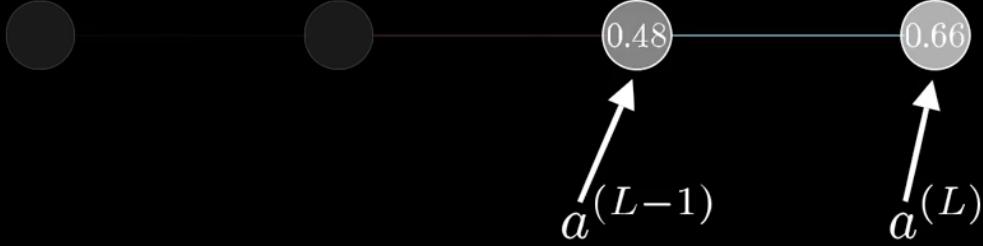




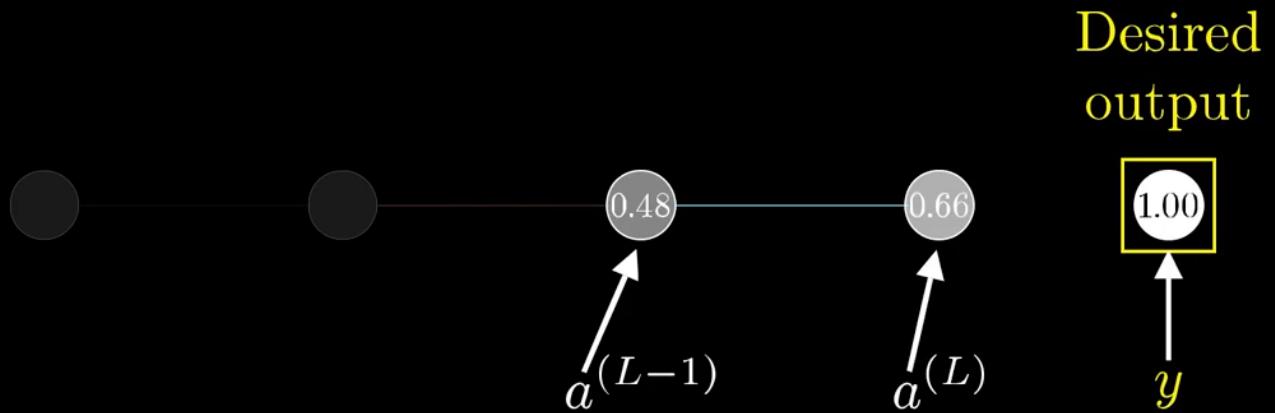
Seems very overwhelming...







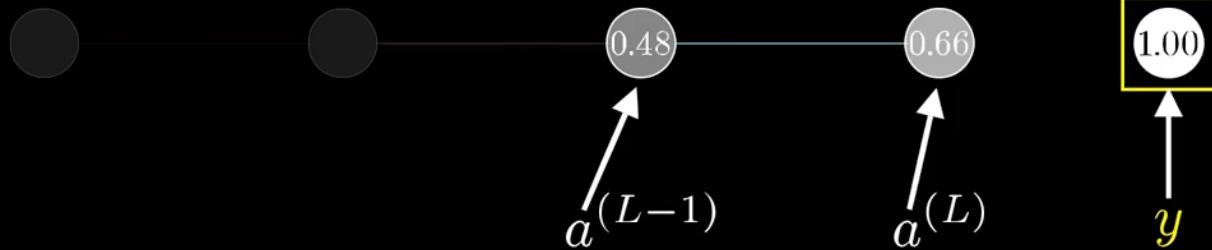
$$C_0(\dots) = (a^{(L)} - y)^2$$



$$\text{Cost} \rightarrow C_0(\dots) = (a^{(L)} - y)^2$$

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

Desired
output

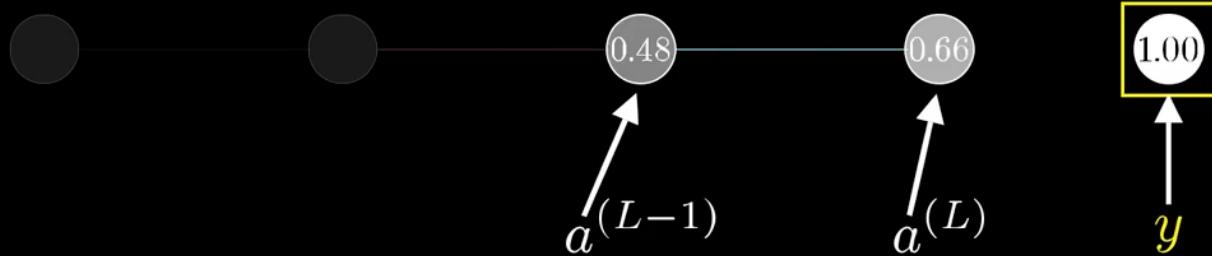


Cost $\rightarrow C_0(\dots) = (a^{(L)} - y)^2$

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

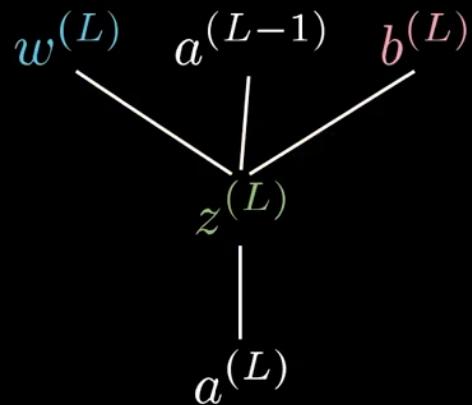
$$a^{(L)} = \sigma(z^{(L)})$$

Desired
output

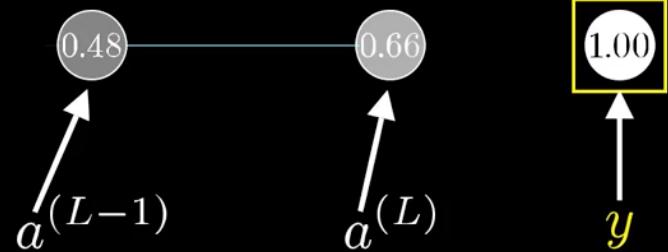


Cost $\rightarrow C_0(\dots) = (a^{(L)} - y)^2$

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

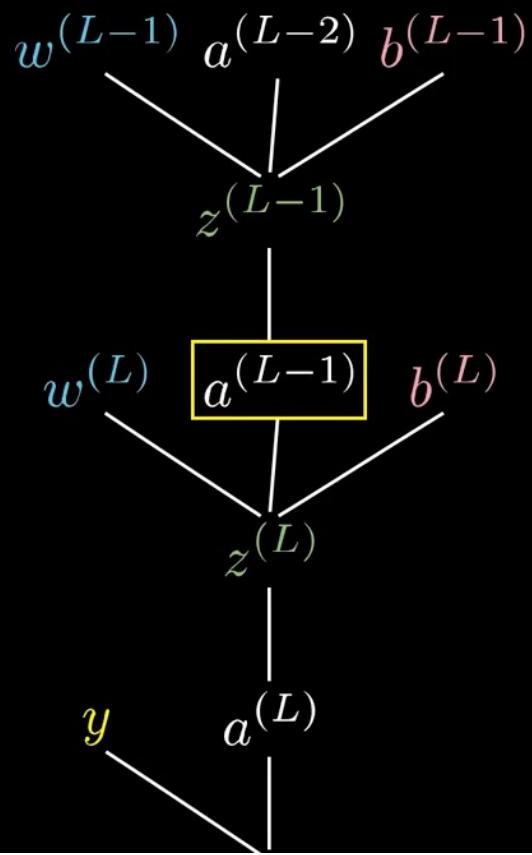


$$a^{(L)} = \sigma(z^{(L)})$$



Desired
output



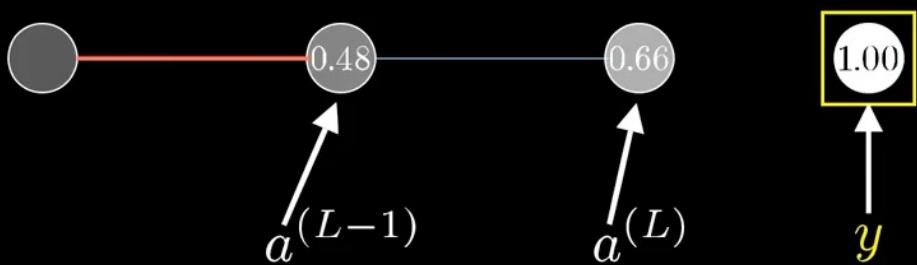


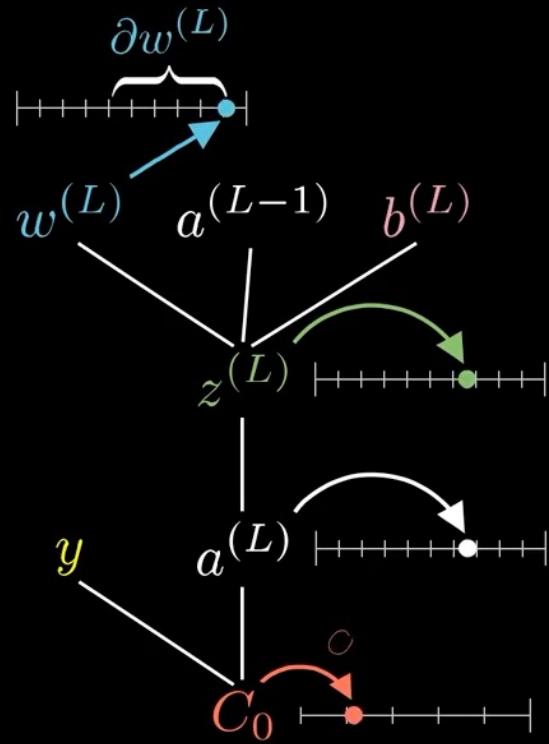
$$\text{Cost} \rightarrow C_0(\dots) = (a^{(L)} - y)^2$$

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

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

Desired
output

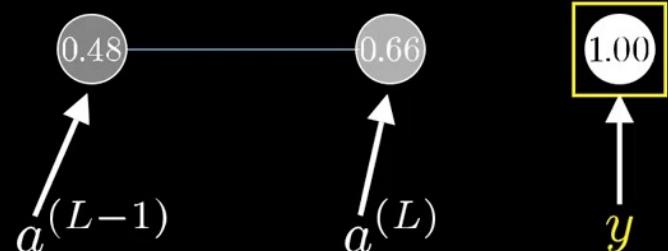




$$\text{Cost} \rightarrow C_0(\dots) = (a^{(L)} - y)^2$$

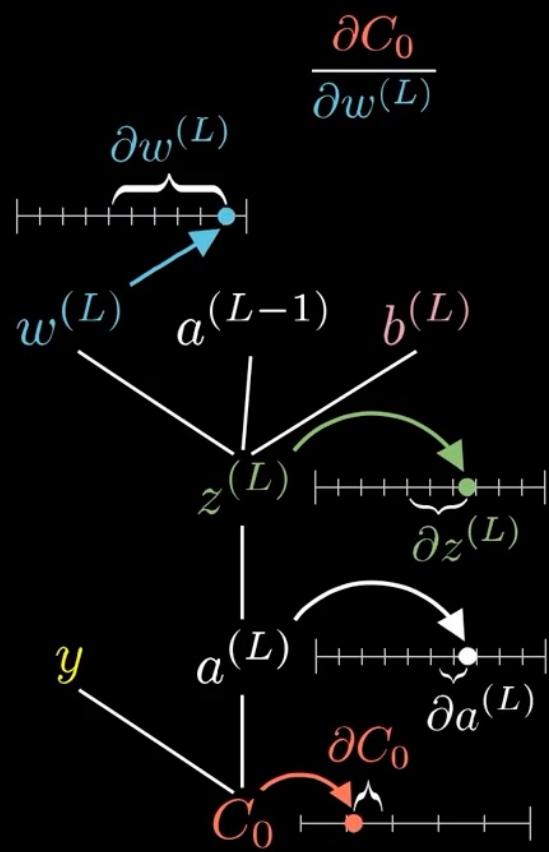
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

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



Desired
output

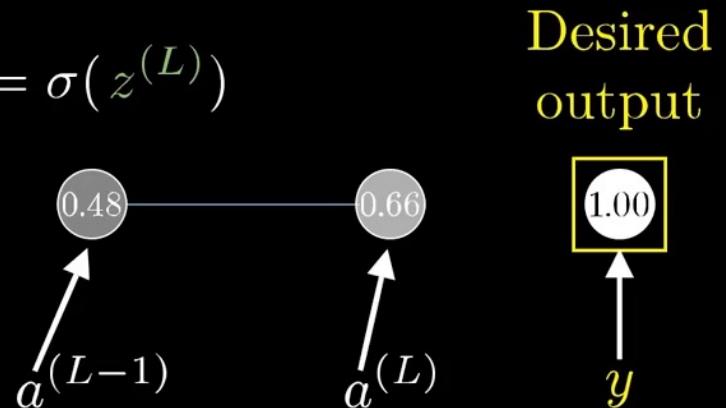


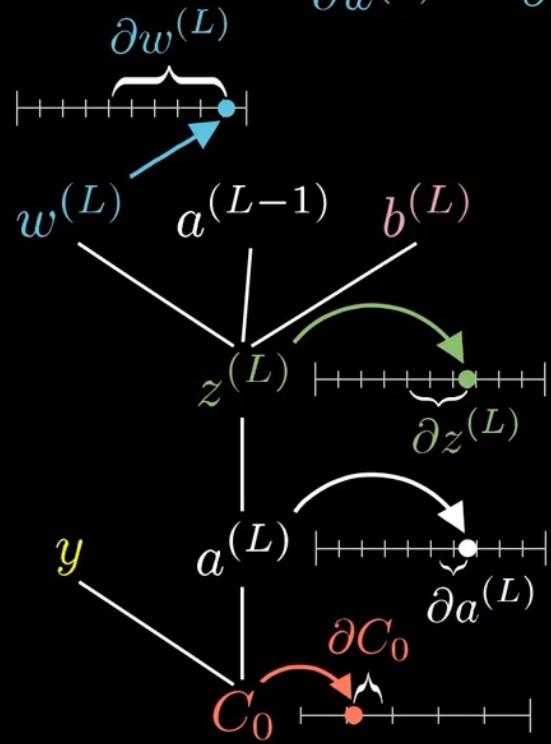


$$C_0(\dots) = (a^{(L)} - y)^2$$

$$z^{(L)} = \mathbf{w}^{(L)} a^{(L-1)} + b^{(L)}$$

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



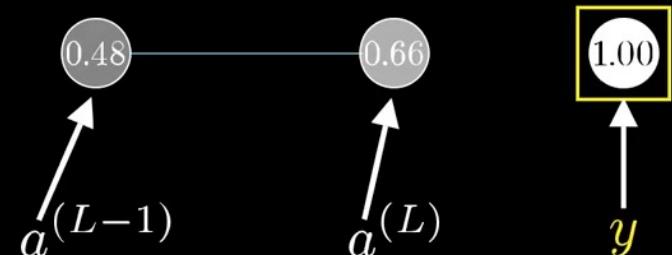


$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}} \quad C_0(\dots) = (a^{(L)} - y)^2$$

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

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

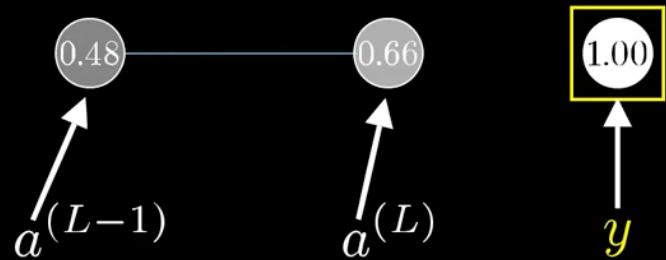
Desired output



$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

$$C_0 = (a^{(L)} - y)^2$$
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

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



$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

$$C_0 = (a^{(L)} - y)^2$$

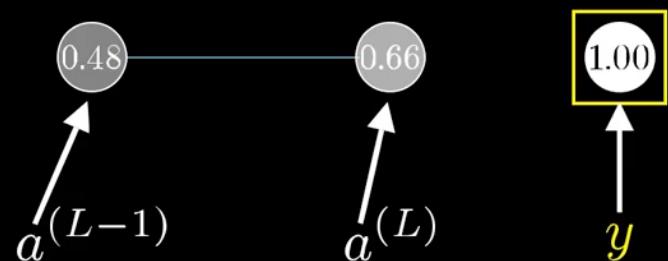
$$\frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

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

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

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

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

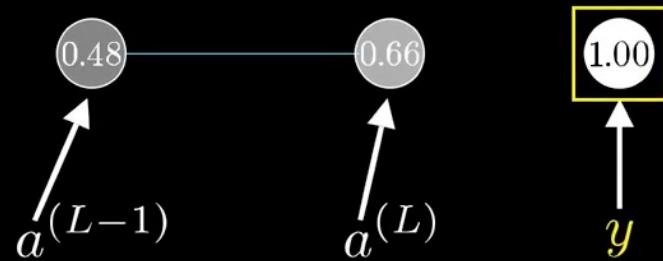


$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}} = a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

$$C_0 = (a^{(L)} - y)^2$$

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

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



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

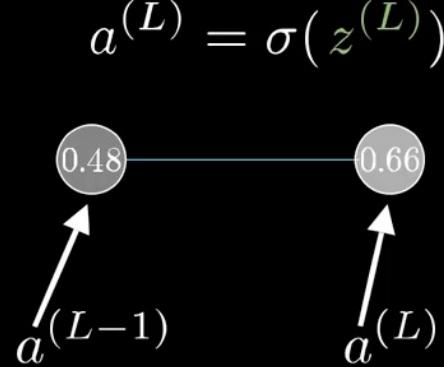
Average of all
training examples

$$C_0 = (a^{(L)} - y)^2$$

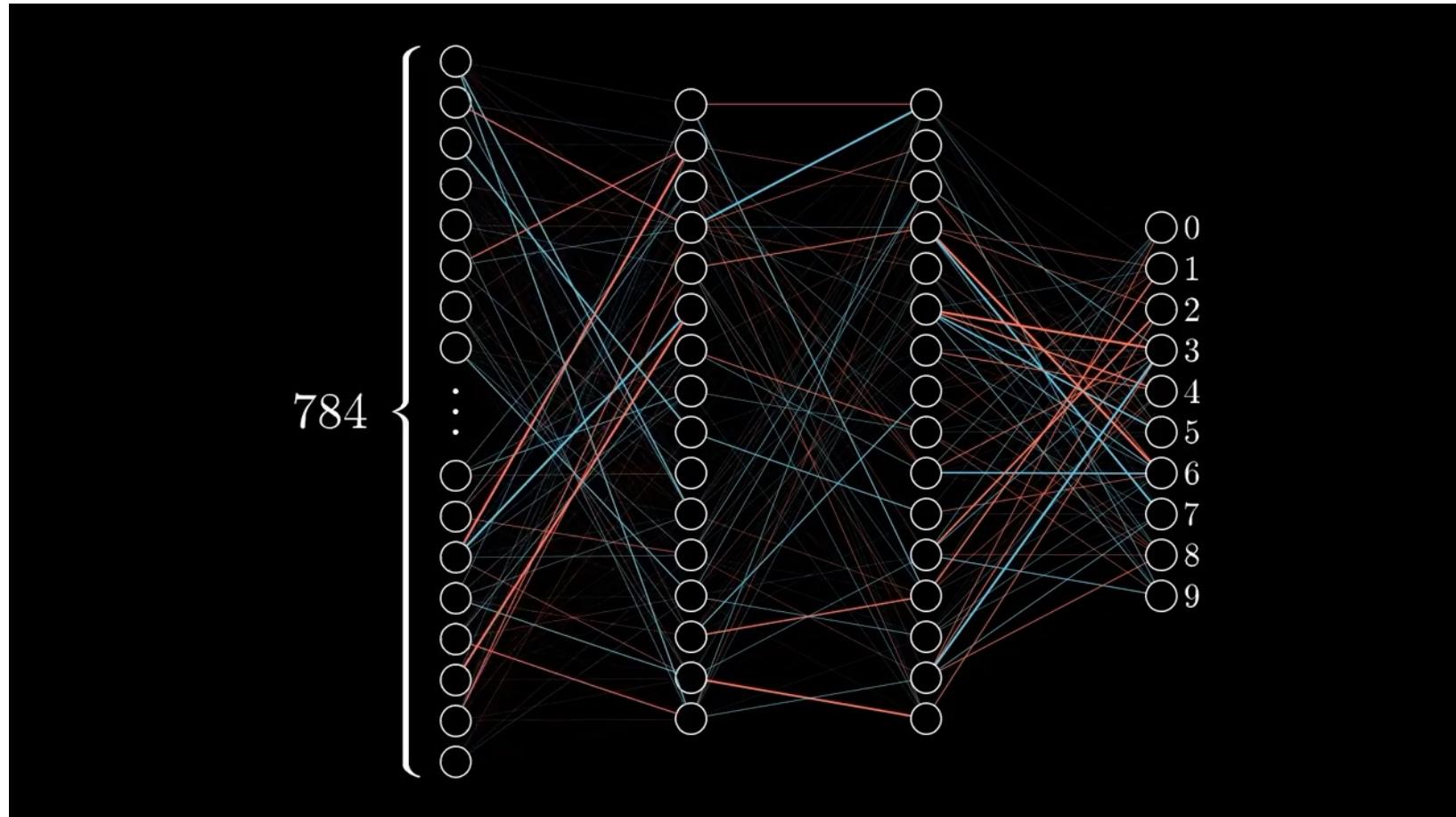
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

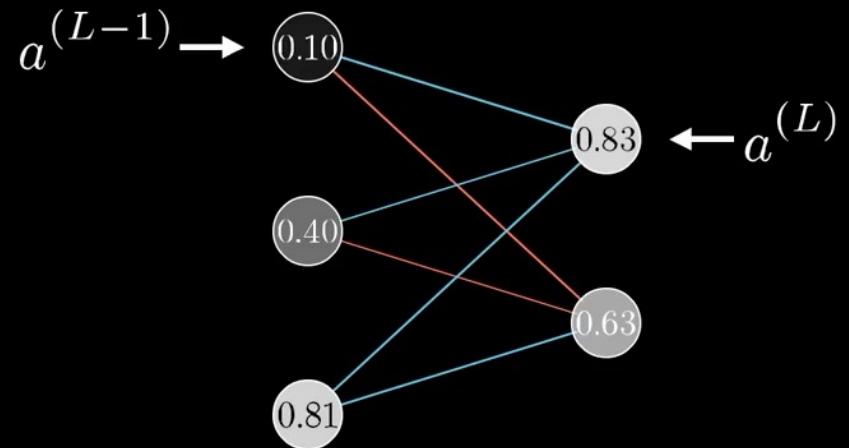
$$\underbrace{\frac{\partial C}{\partial w^{(L)}}}_{\text{Derivative of full cost function}} = \overbrace{\frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}}^{\text{Average of all training examples}}$$

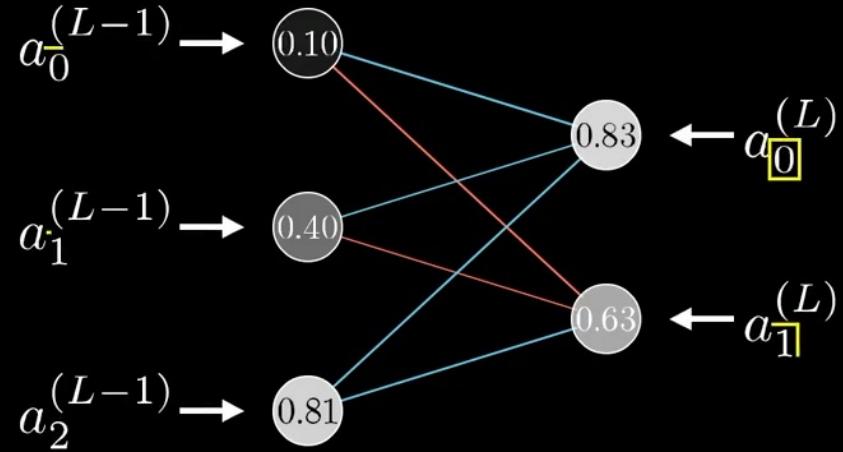
Derivative of
full cost function

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


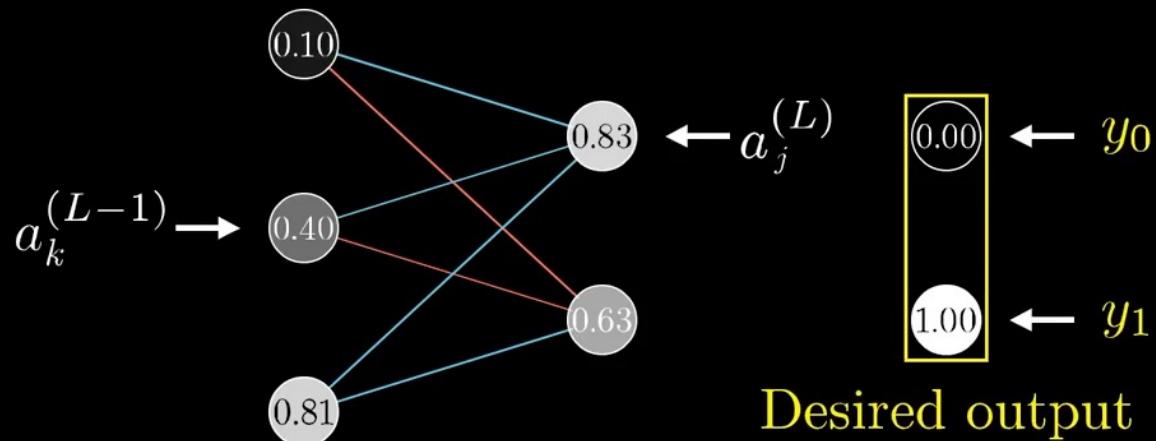




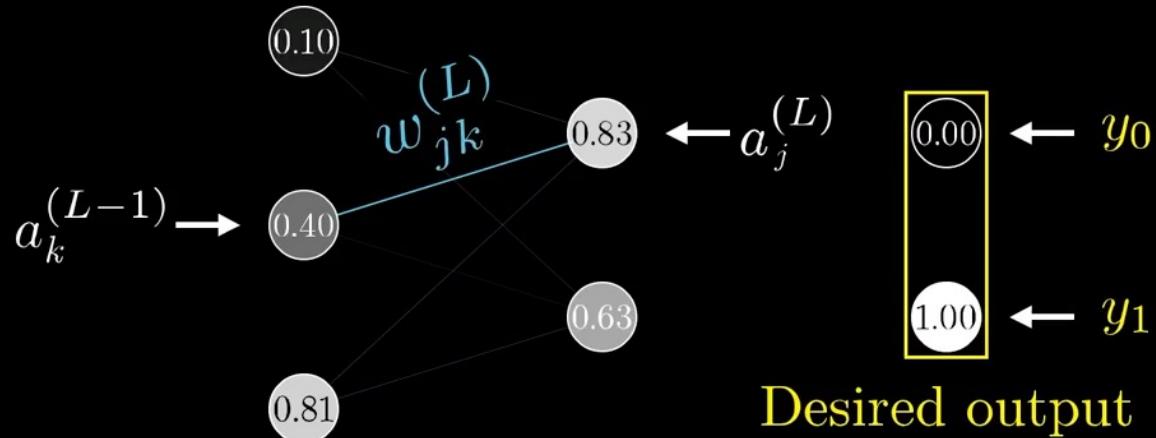




$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

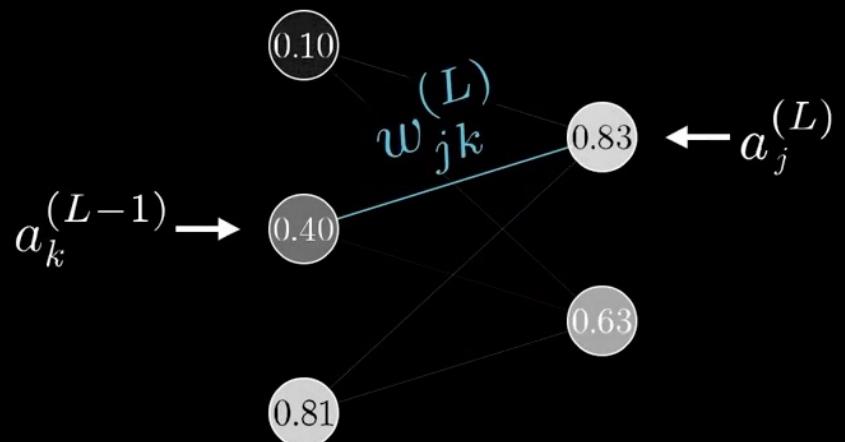


$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$



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

$$a_j^{(L)} = \sigma(z_j^{(L)})$$



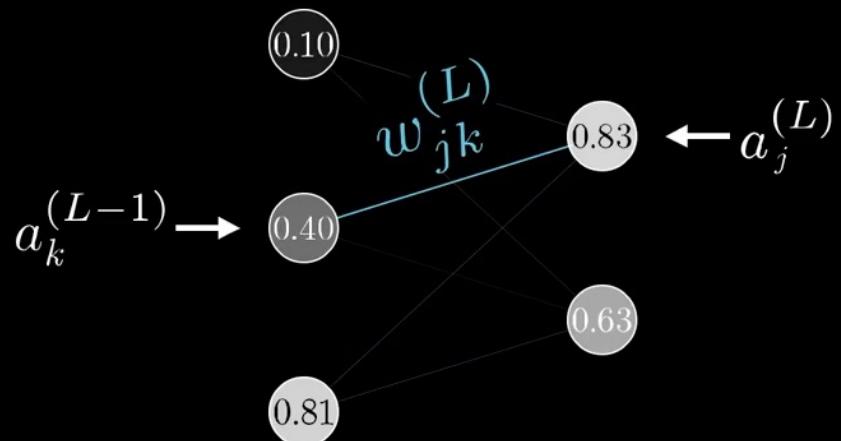
$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$



$$\frac{\partial C_0}{\partial w_{jk}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}$$

$$z_j^{(L)} = \dots + w_{jk}^{(L)} a_k^{(L-1)} + \dots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

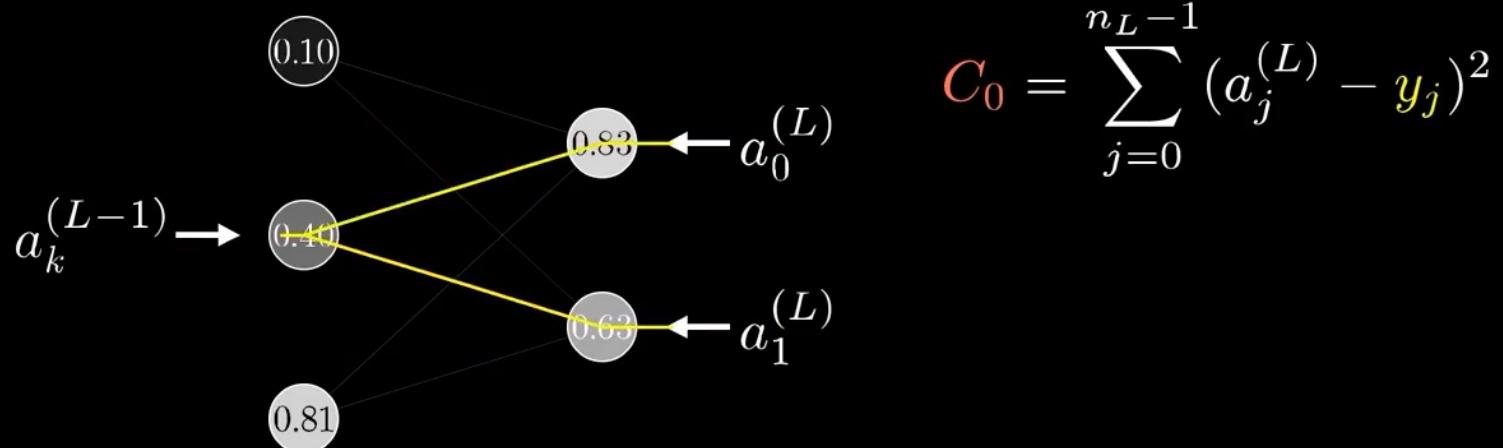


$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$



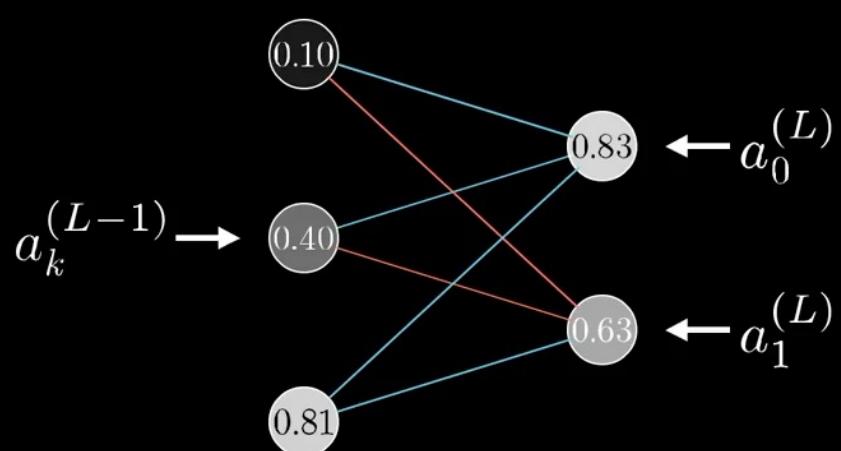
$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}} \quad z_j^{(L)} = \dots + w_{jk}^{(L)} a_k^{(L-1)} + \dots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$



$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \underbrace{\sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}}_{\text{Sum over layer L}} \quad z_j^{(L)} = \dots + w_{jk}^{(L)} a_k^{(L-1)} + \dots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

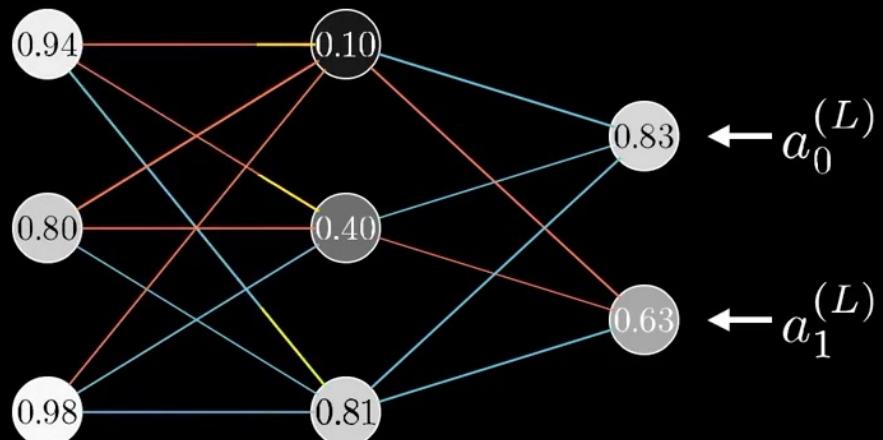


$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$



$$\boxed{\frac{\partial C_0}{\partial a_k^{(L-1)}}} = \underbrace{\sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}}_{\text{Sum over layer L}} \quad z_j^{(L)} = \dots + w_{jk}^{(L)} a_k^{(L-1)} + \dots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$



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



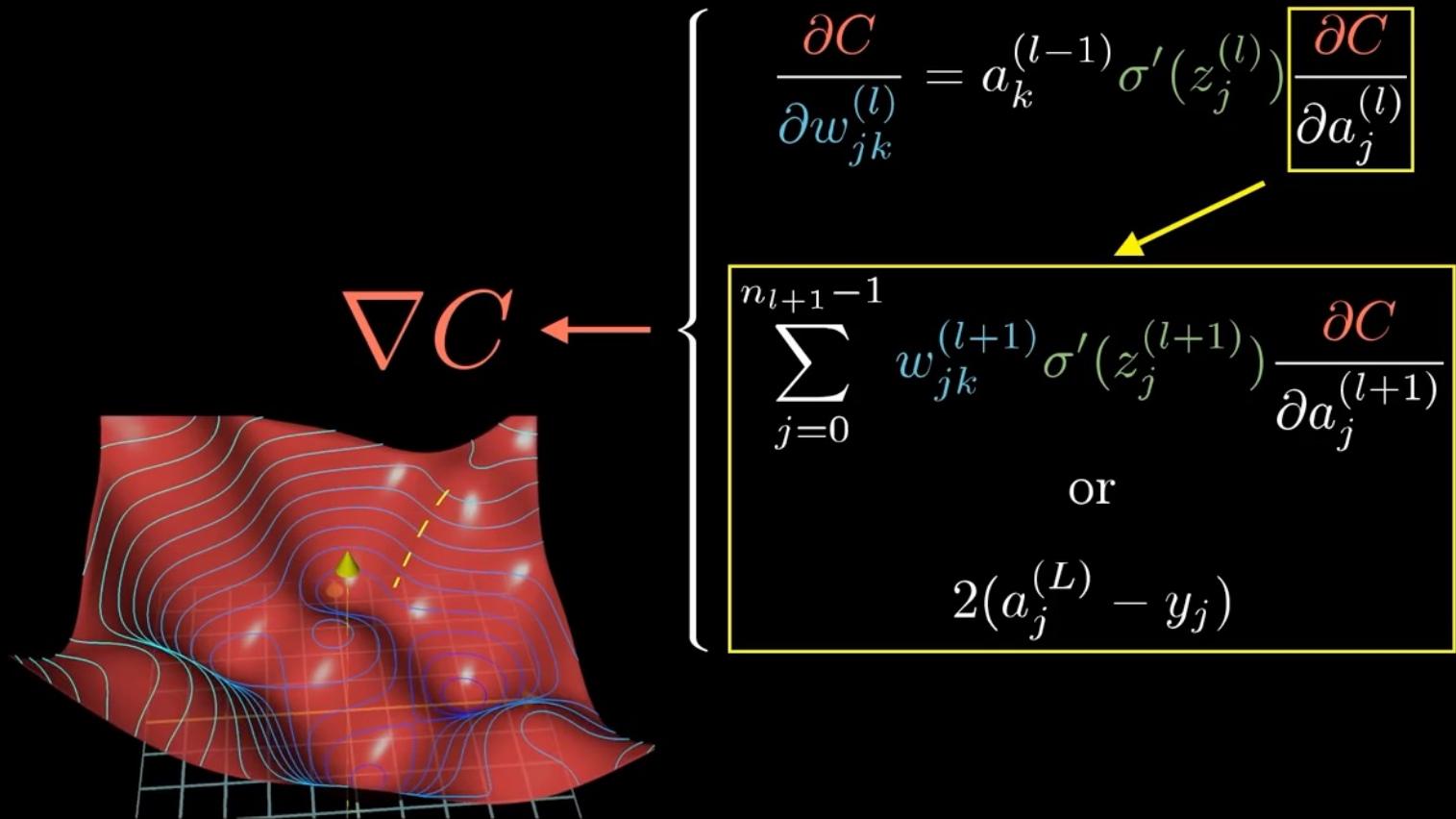
$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}} = a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

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

$C_0 = (a^{(L)} - y)^2$
 $z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$
 $a^{(L)} = \sigma(z^{(L)})$

The diagram illustrates a single layer of a neural network. It shows two input nodes labeled $a^{(L-1)}$ with values 0.48 and 0.66. These inputs feed into two output nodes labeled $a^{(L)}$, which have values 0.66 and 1.00. The output node with value 1.00 is highlighted with a yellow box. Arrows indicate the flow of information from the input layer to the output layer.





A little recap about Backprop...



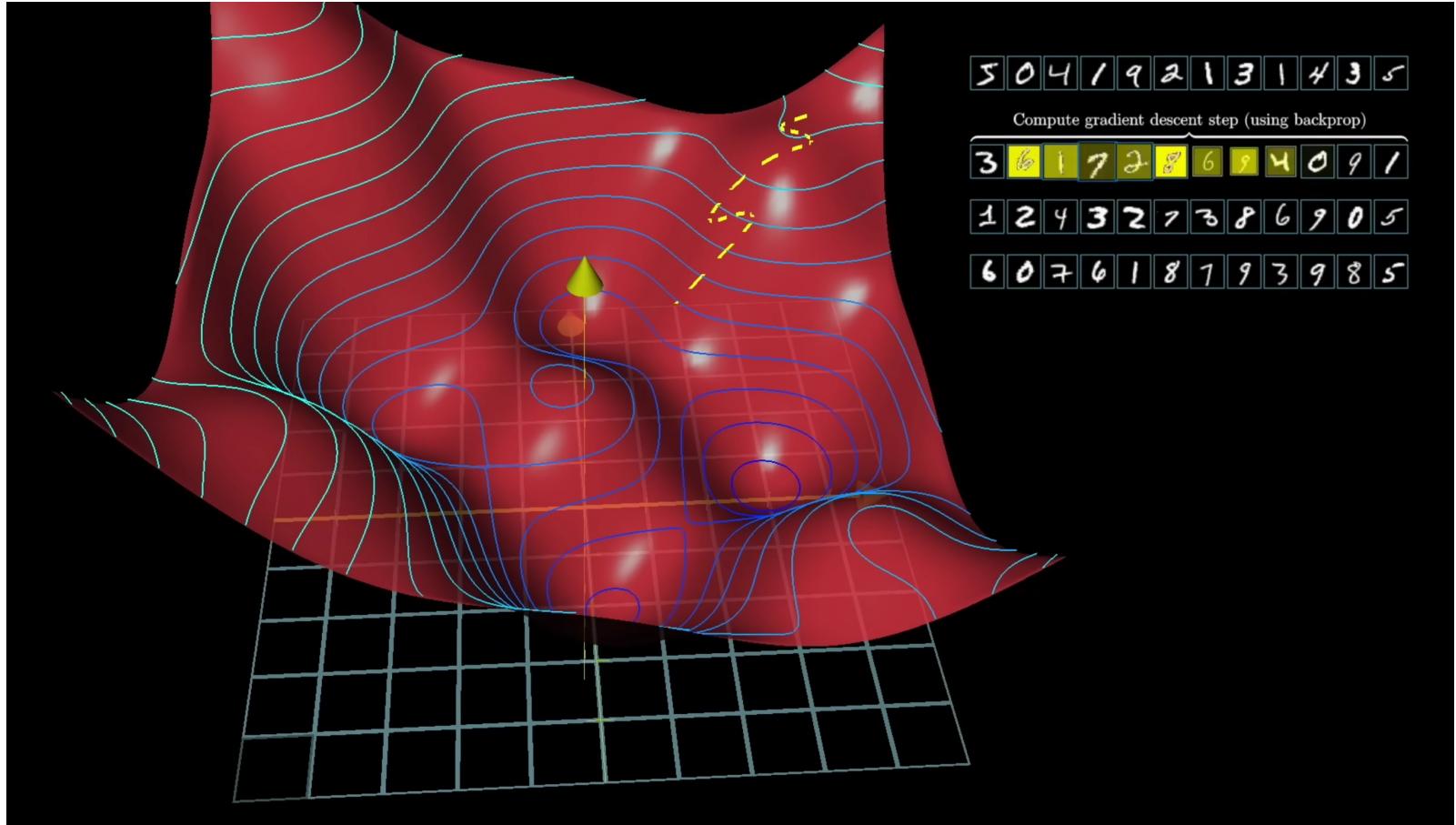
1	1	0	6	4	0	4	9	8	1	3	6
3	4	5	2	9	7	9	9	7	6	4	8
5	3	0	3	9	9	8	1	4	0	1	7
2	8	3	3	6	0	4	1	5	8	5	6
0	9	7	9	2	7	9	2	1	1	3	3



“Mini-batches”

5	0	4	1	9	2	1	3	1	4	3	5
3	6	1	7	2	8	6	9	4	0	9	1
1	2	4	3	2	7	3	8	6	9	0	5
6	0	7	6	1	8	7	9	3	9	8	5





5 0 4 1 9 2 1 3 1 4 3 5

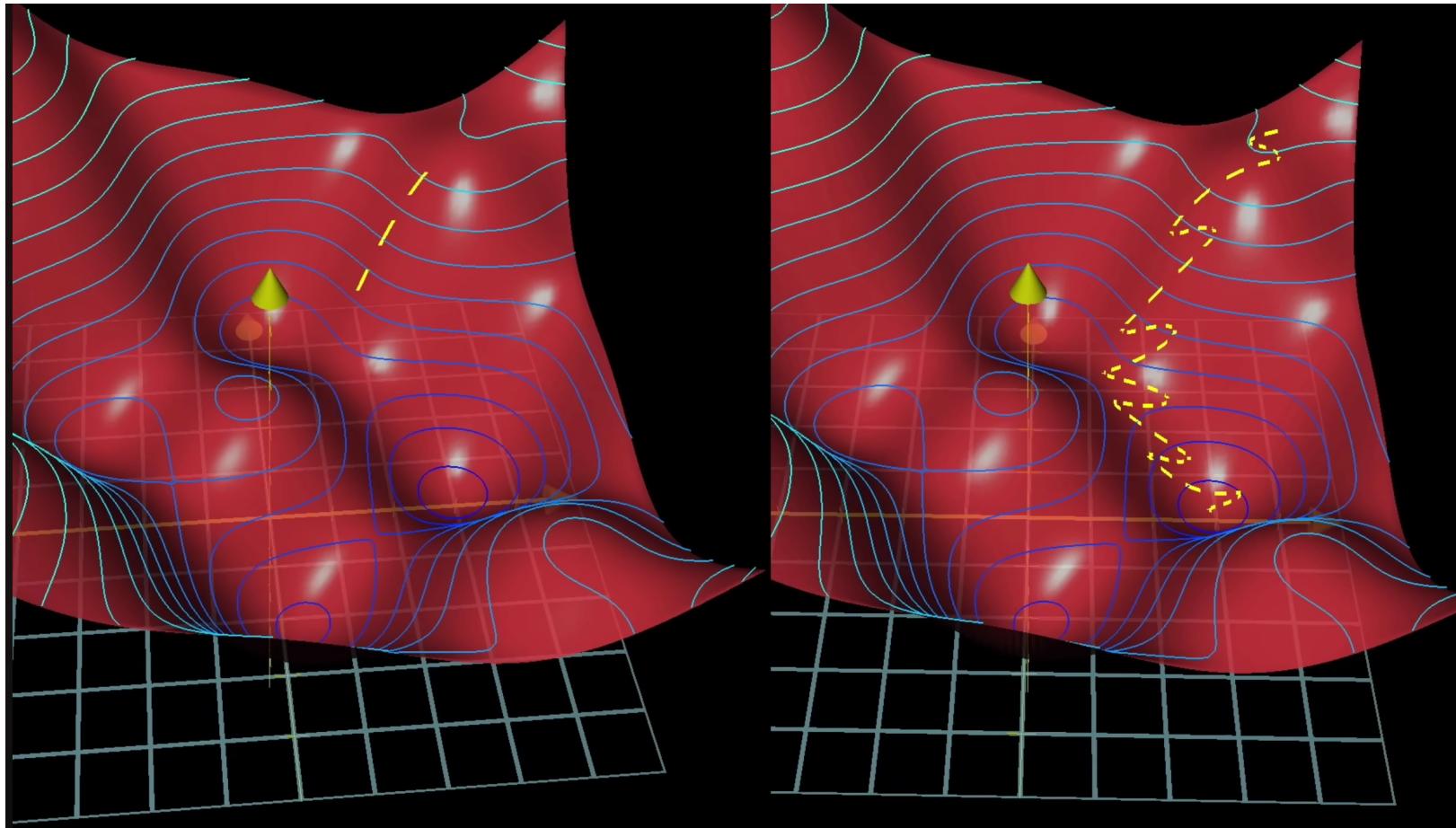
Compute gradient descent step (using backprop)

3 6 1 7 2 8 6 9 4 0 9 1

1 2 4 3 2 7 3 8 6 9 0 5

6 0 7 6 1 8 7 9 3 9 8 5





Summary in words

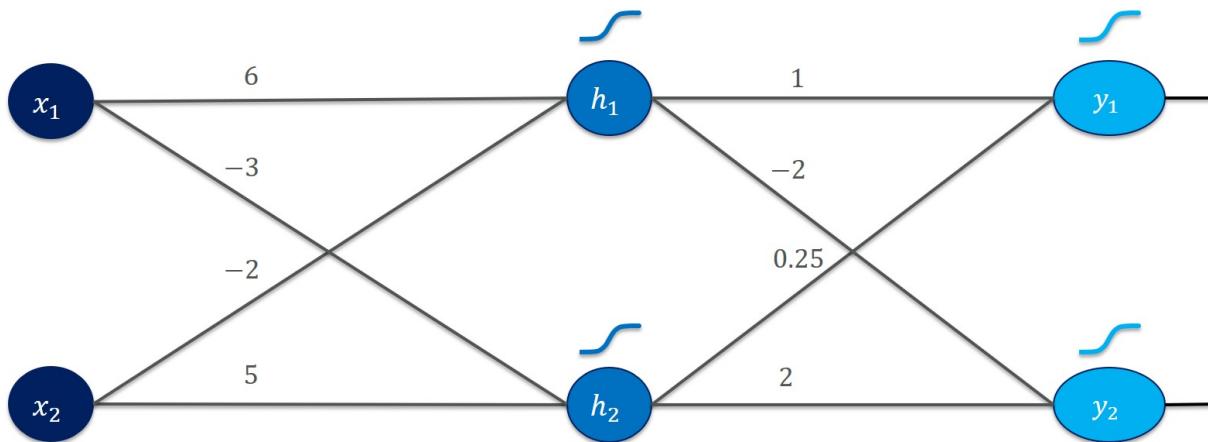
- Backprop determine how a single training example nudge the weights and biases
- Gradient descent step involve doing this for x1000000 of times
 - But that computationally slow...
- Introduce mini-batches
 - Stochastic gradient descent

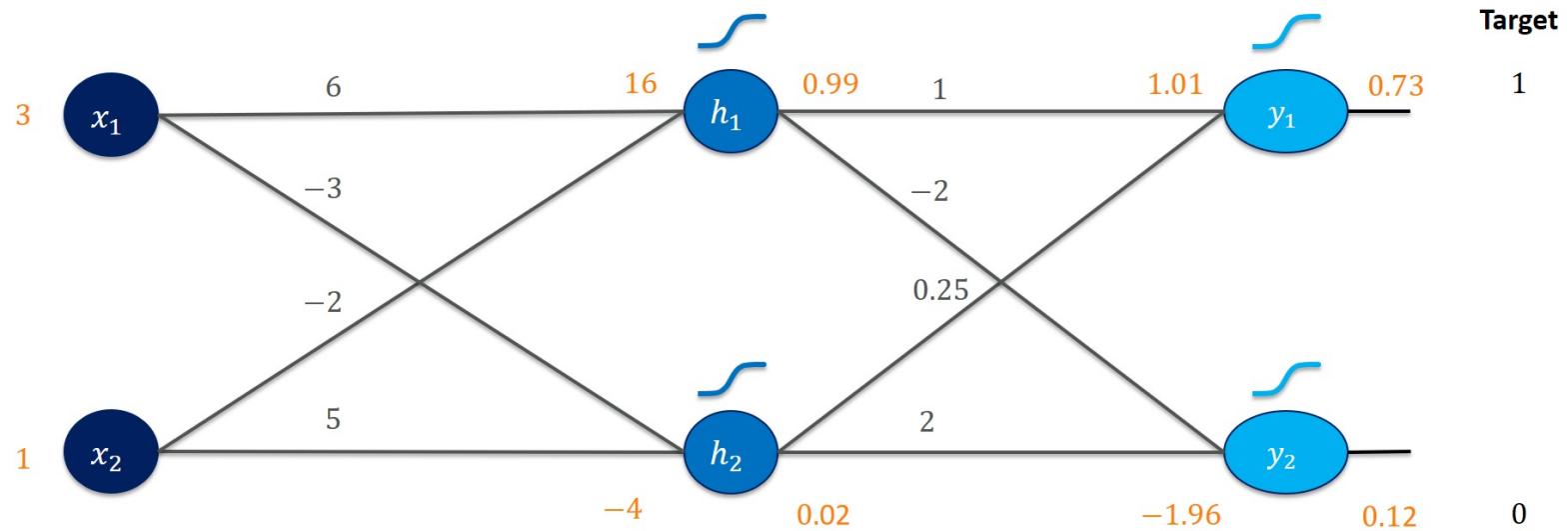


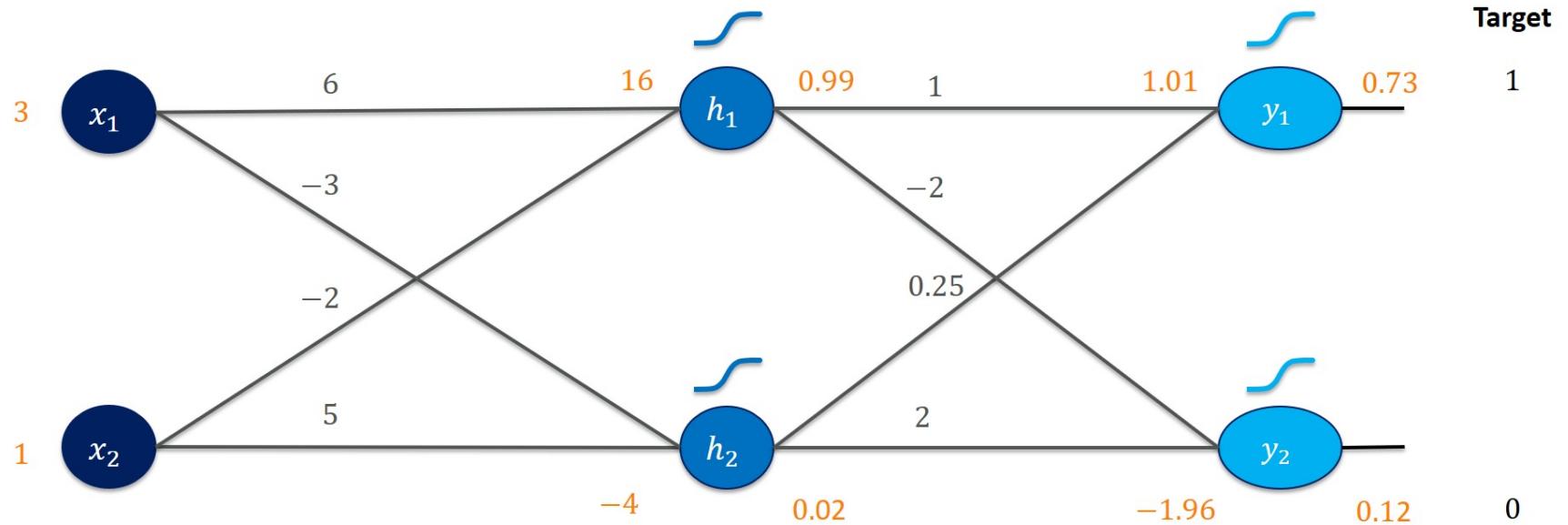
Examples work through



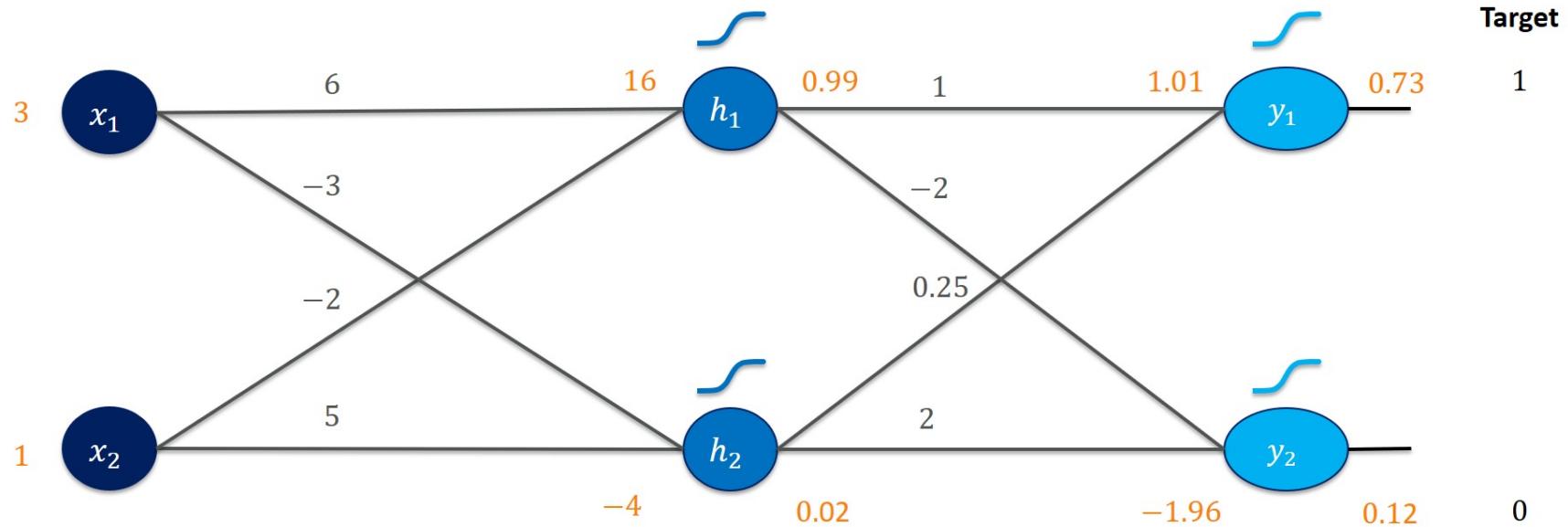
	Inputs		Targets	
Training data	x_1	x_2	t_1	t_2
	3	1	1	0
	-1	4	0	1

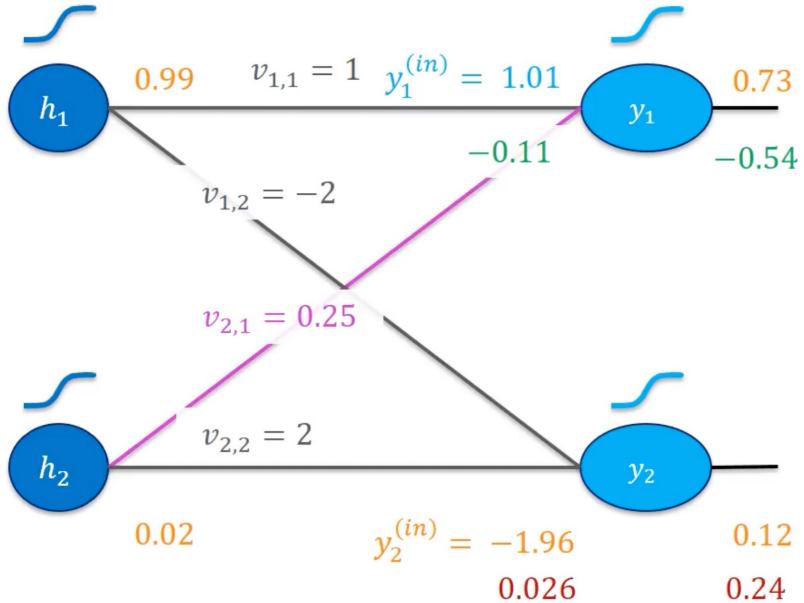


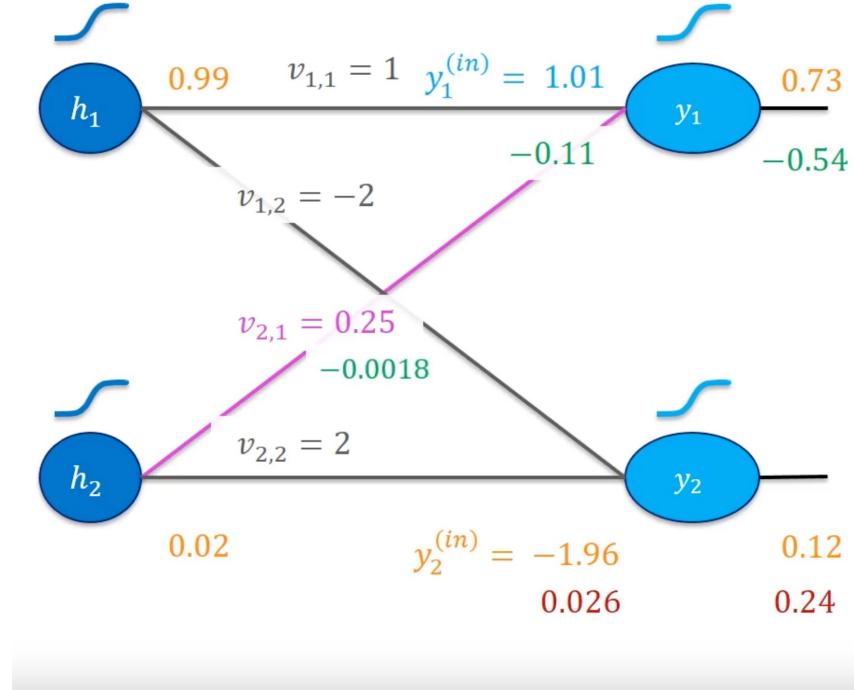




$$L(y_1, y_2) = (y_1 - t_1)^2 + (y_2 - t_2)^2 = e_1^2 + e_2^2$$







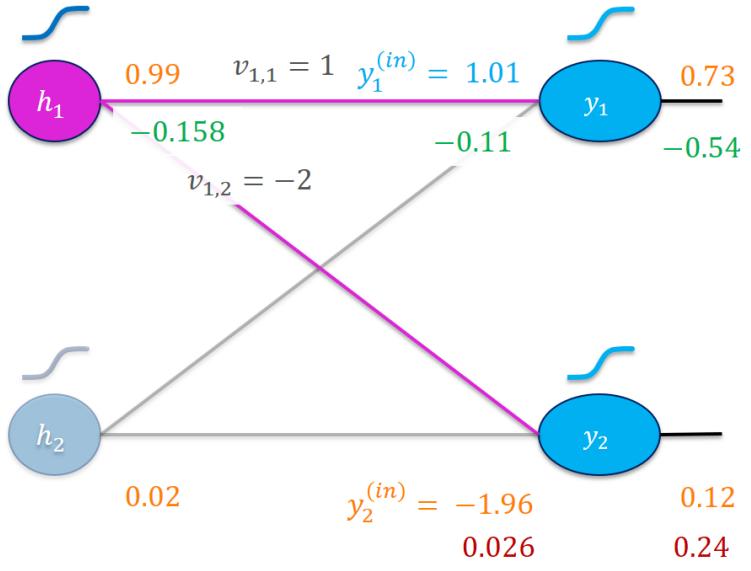
$$y_1^{(in)} = v_{1,1} \cdot h_1 + v_{2,1} \cdot h_2$$

$$\frac{\partial L}{\partial y_1^{(in)}} = -0.11$$

$$\frac{\partial L}{\partial v_{2,1}} = \frac{\partial L}{\partial y_1^{(in)}} \cdot \underbrace{\frac{\partial y_1^{(in)}}{\partial v_{2,1}}}_{h_2}$$

$$= -0.11 \cdot 0.02 \\ = -0.0018$$





$$y_1^{(in)} = v_{1,1} \cdot h_1 + v_{2,1} \cdot h_2$$

$$y_2^{(in)} = v_{1,2} \cdot h_1 + v_{2,2} \cdot h_2$$

$$\frac{\partial L}{\partial y_1^{(in)}} = -0.11$$

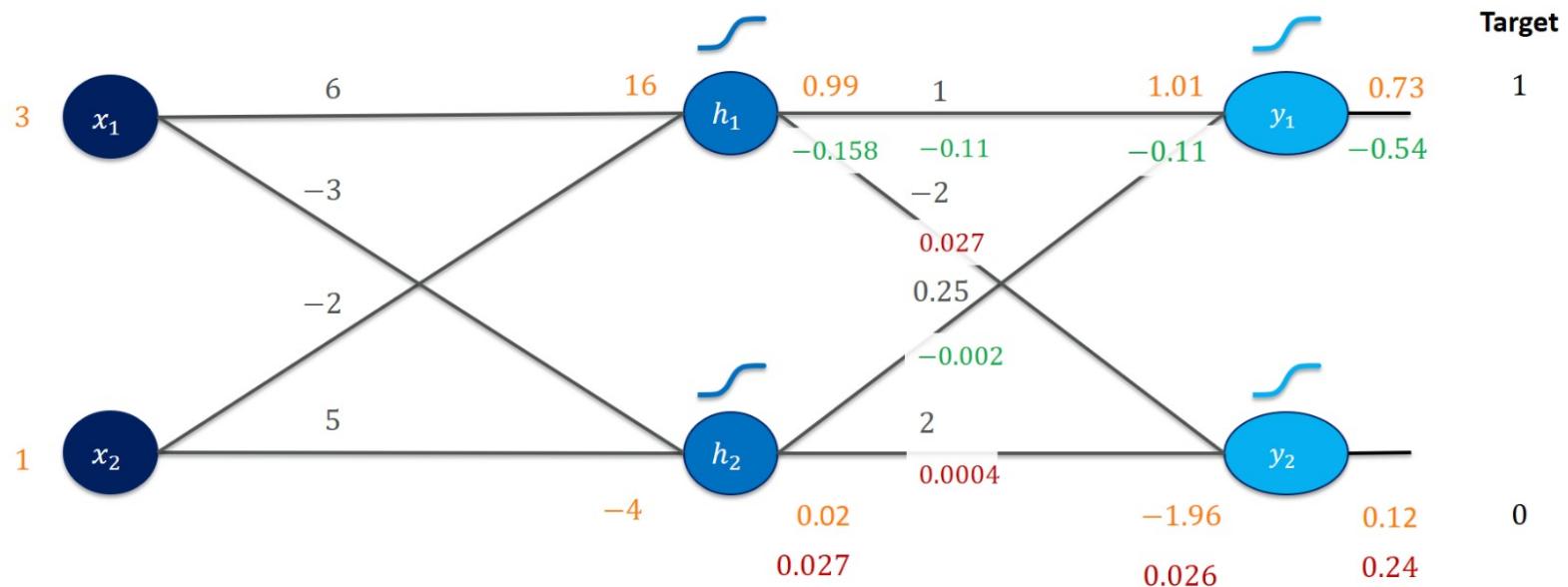
$$\frac{\partial L}{\partial y_2^{(in)}} = 0.026$$

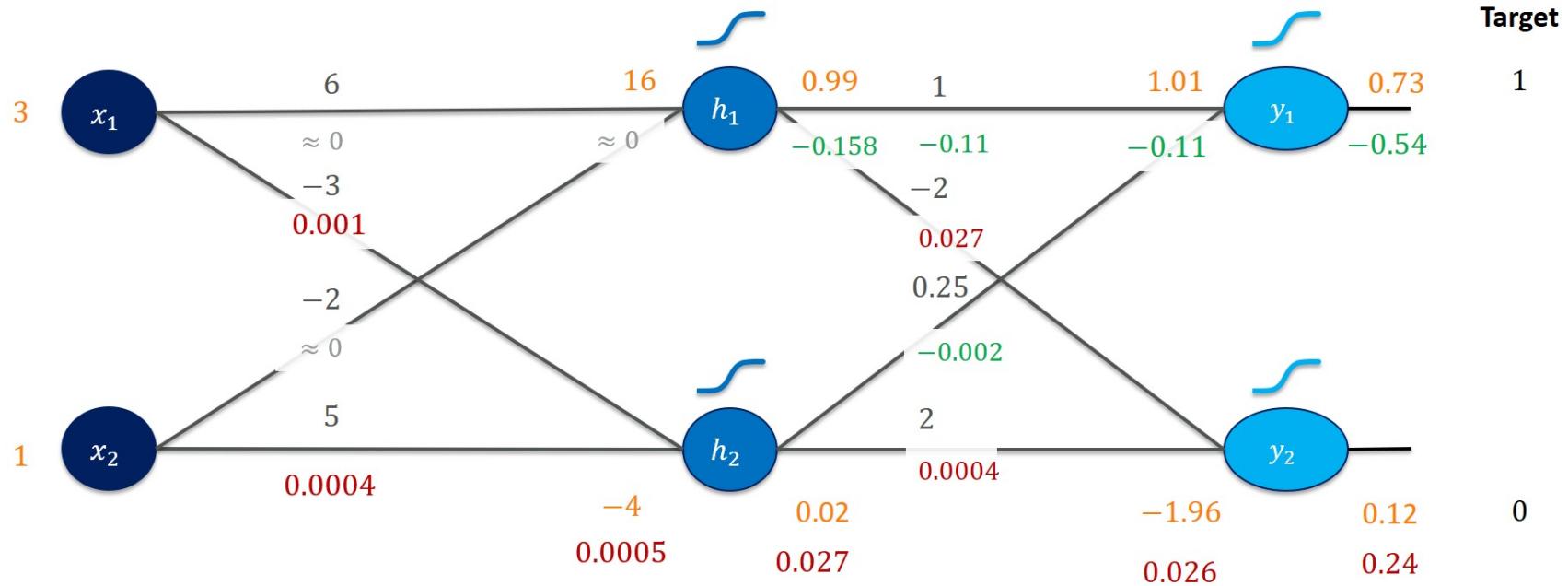
$$\frac{\partial L}{\partial h_1} = \underbrace{\frac{\partial L}{\partial y_1^{(in)}} \cdot \frac{\partial y_1^{(in)}}{\partial h_1}}_{v_{1,1}} + \underbrace{\frac{\partial L}{\partial y_2^{(in)}} \cdot \frac{\partial y_2^{(in)}}{\partial h_1}}_{v_{1,2}}$$

$$= -0.11 \cdot 1 + 0.026 \cdot (-2)$$

$$= -0.158$$







$$\nabla_w^{(1)} L = \boxed{0 \quad 0 \quad 0.0001 \quad 0.00004 \quad -0.11 \quad 0.027 \quad -0.002 \quad 0.0004}$$



Activation function

- Why is activation function useful?



Activation function

- Why is activation function useful?
 - Non-linearity
 - *Great for deep learning!*
 - Gradient Propagation
 - *Ensure back propagation*
 - Regularization
 - *How (Implicit) Regularization of ReLU Neural Networks Characterizes the Learned Function (Jakob, 2023)*

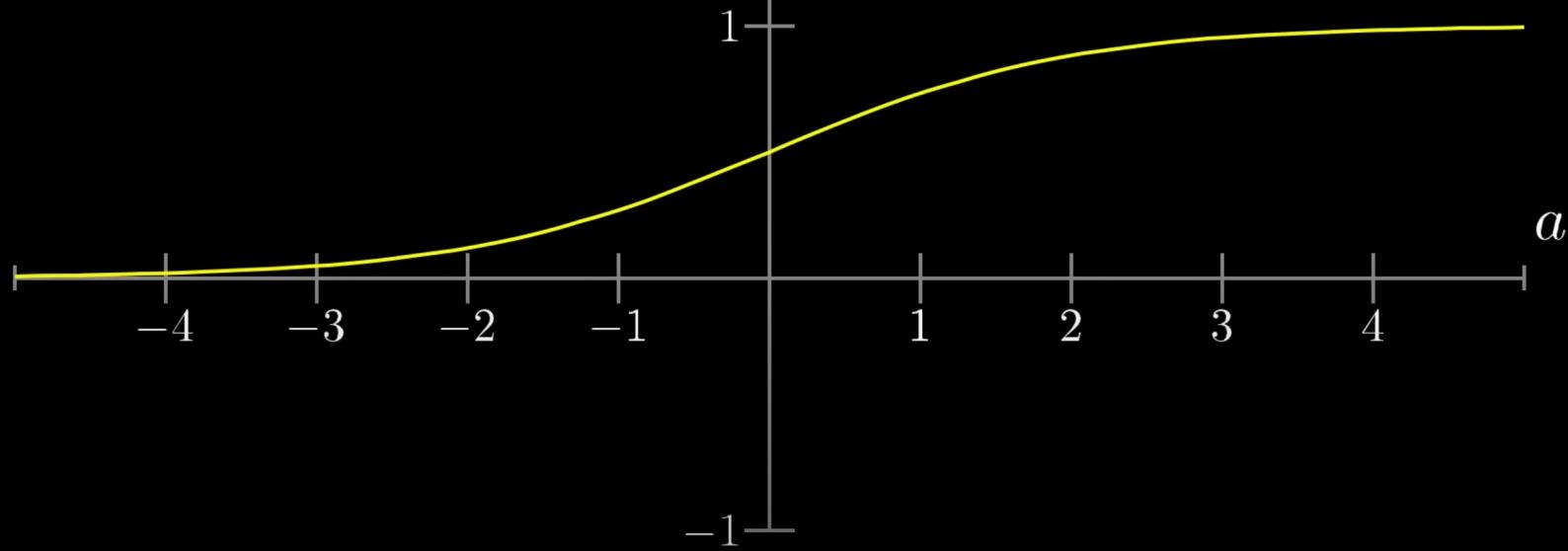


ReLU vs. Sigmoid



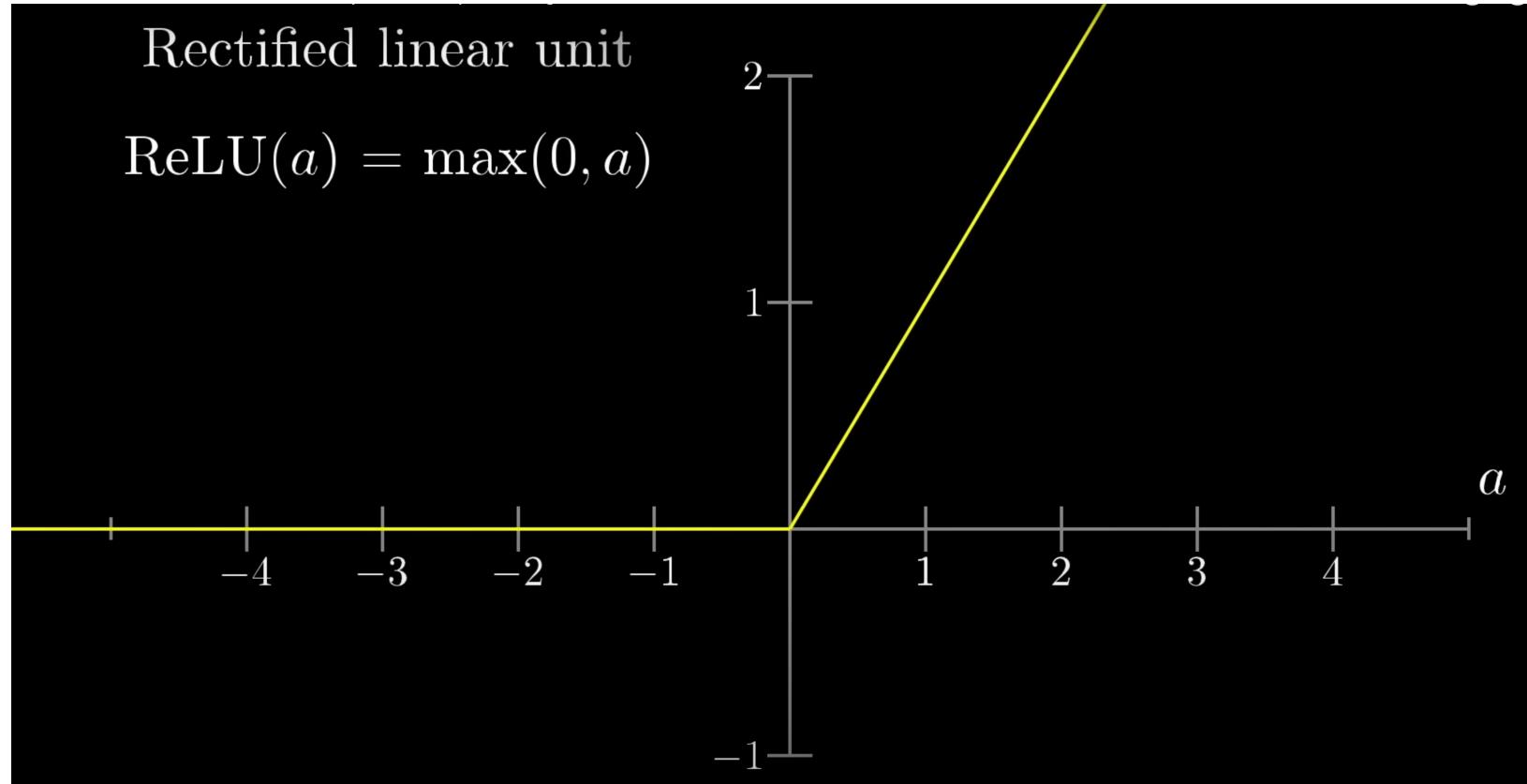
Sigmoid

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$



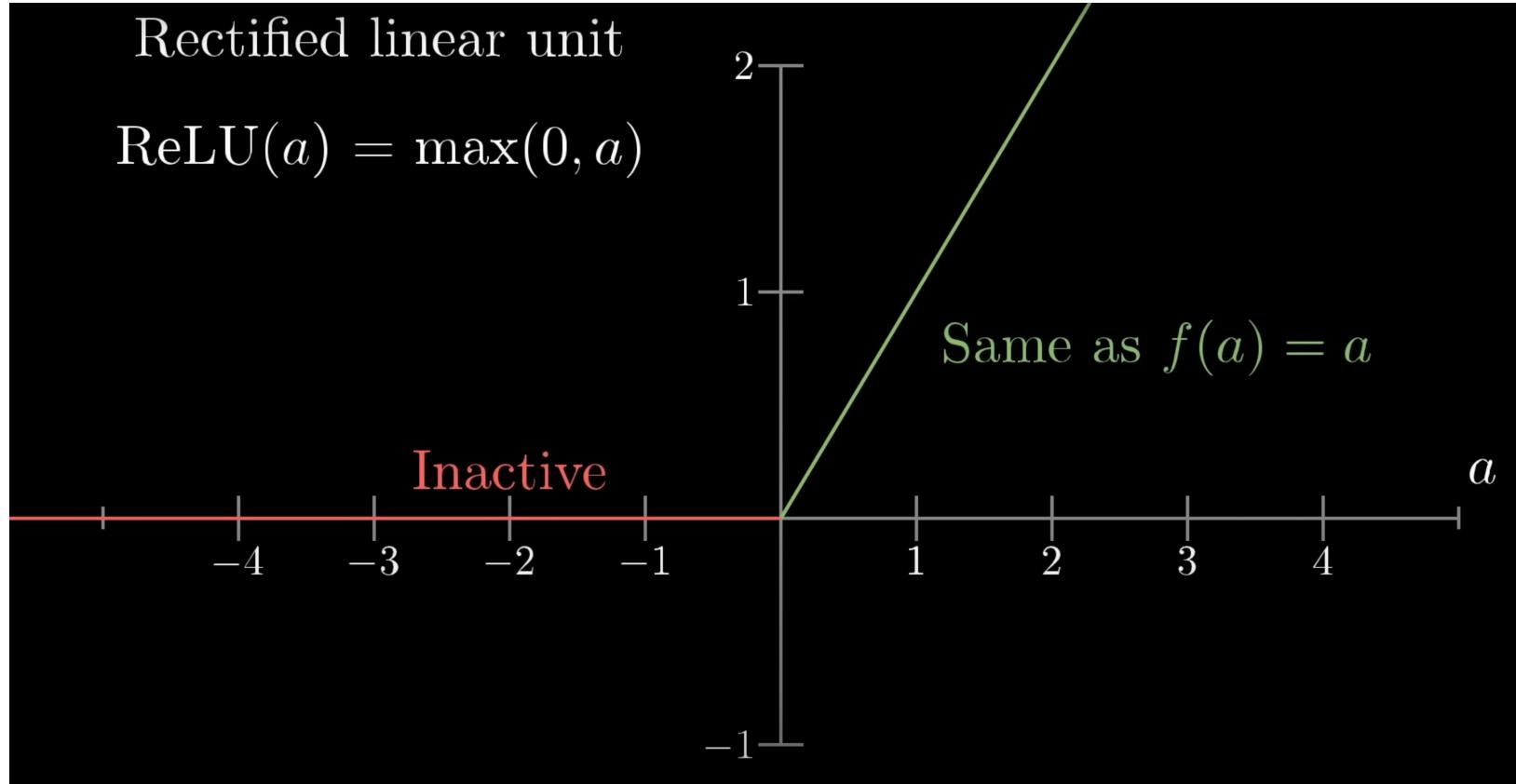
Rectified linear unit

$$\text{ReLU}(a) = \max(0, a)$$



Rectified linear unit

$$\text{ReLU}(a) = \max(0, a)$$



$\text{Max}(0, x)$

ReLU

$$\frac{1}{1 + e^{-x}}$$

Softmax



Let's see how MLP works in code



