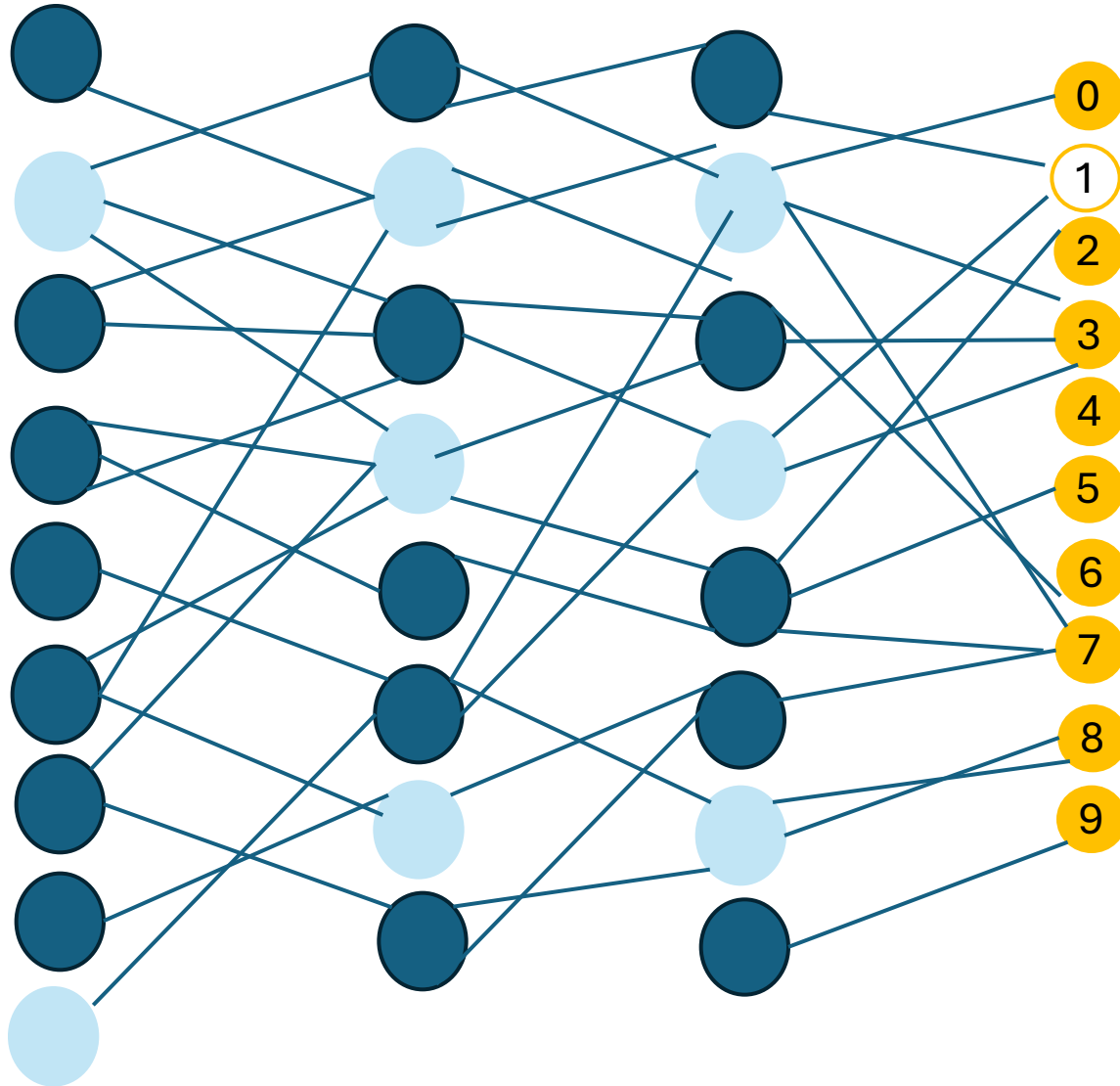
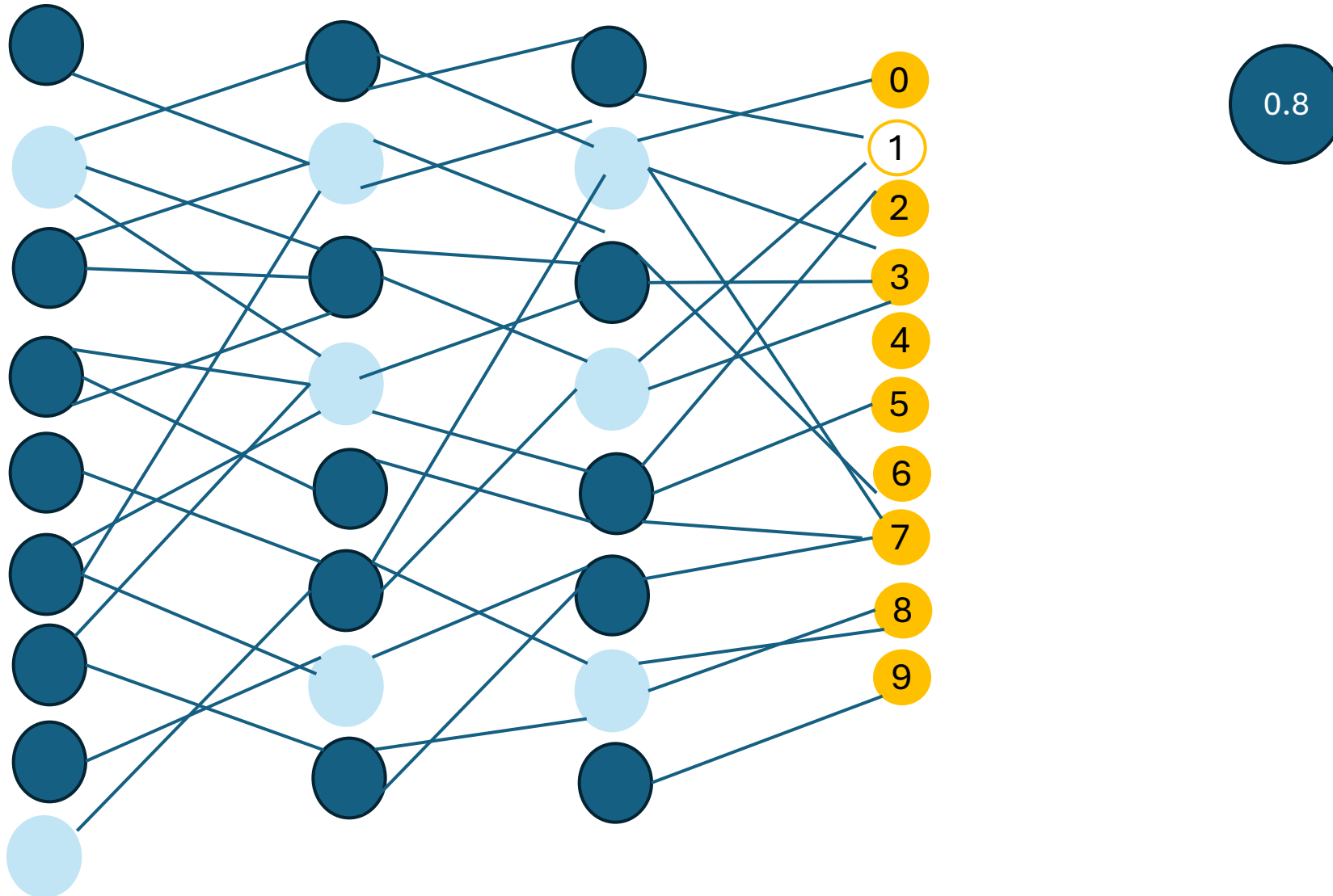


# Introduction to Neural Network

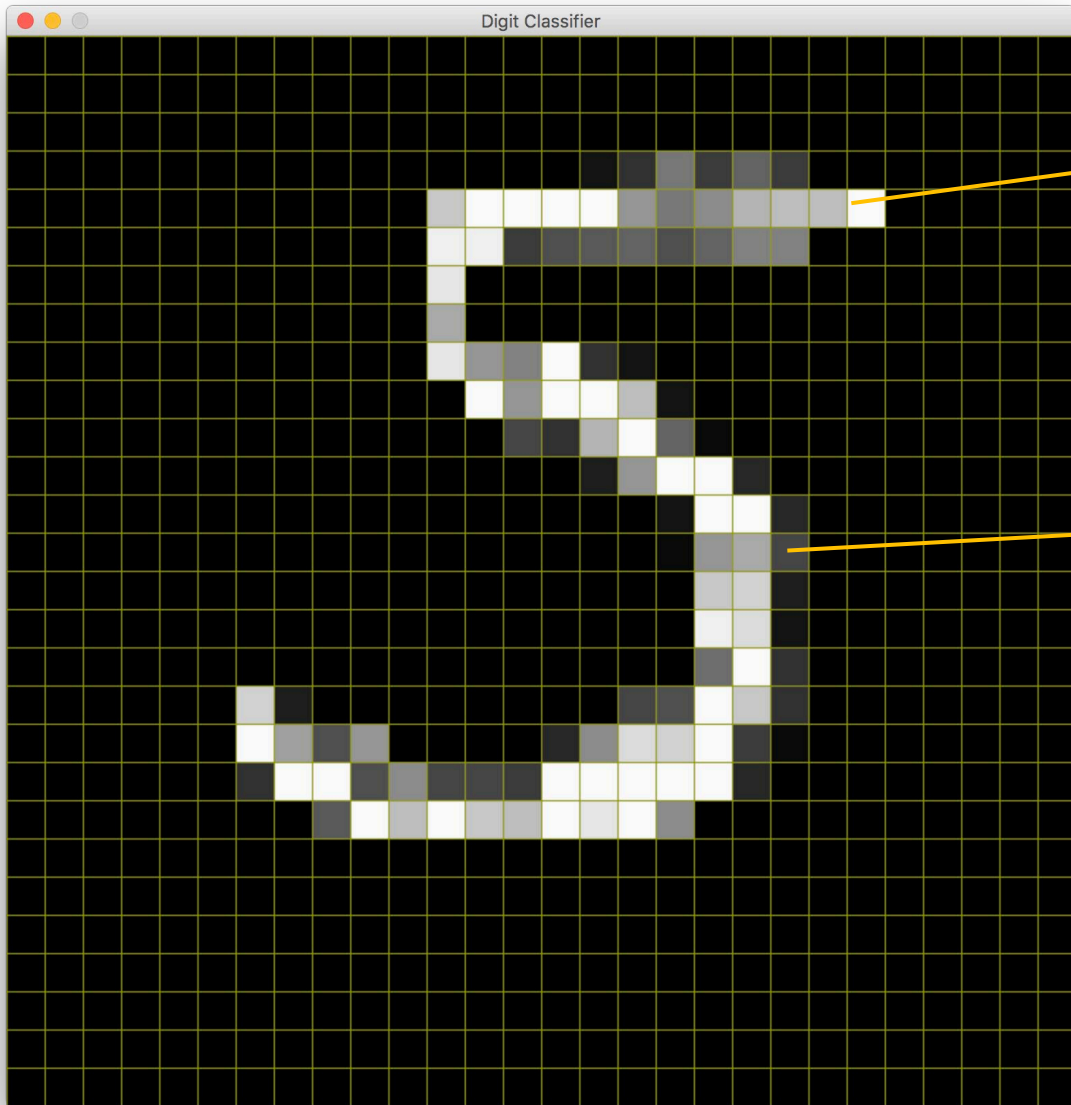
# Multilayer Perceptron



What is a Neuron → Thing that holds a number



## 28X28 pixels with (0-1) values **Input Layer**



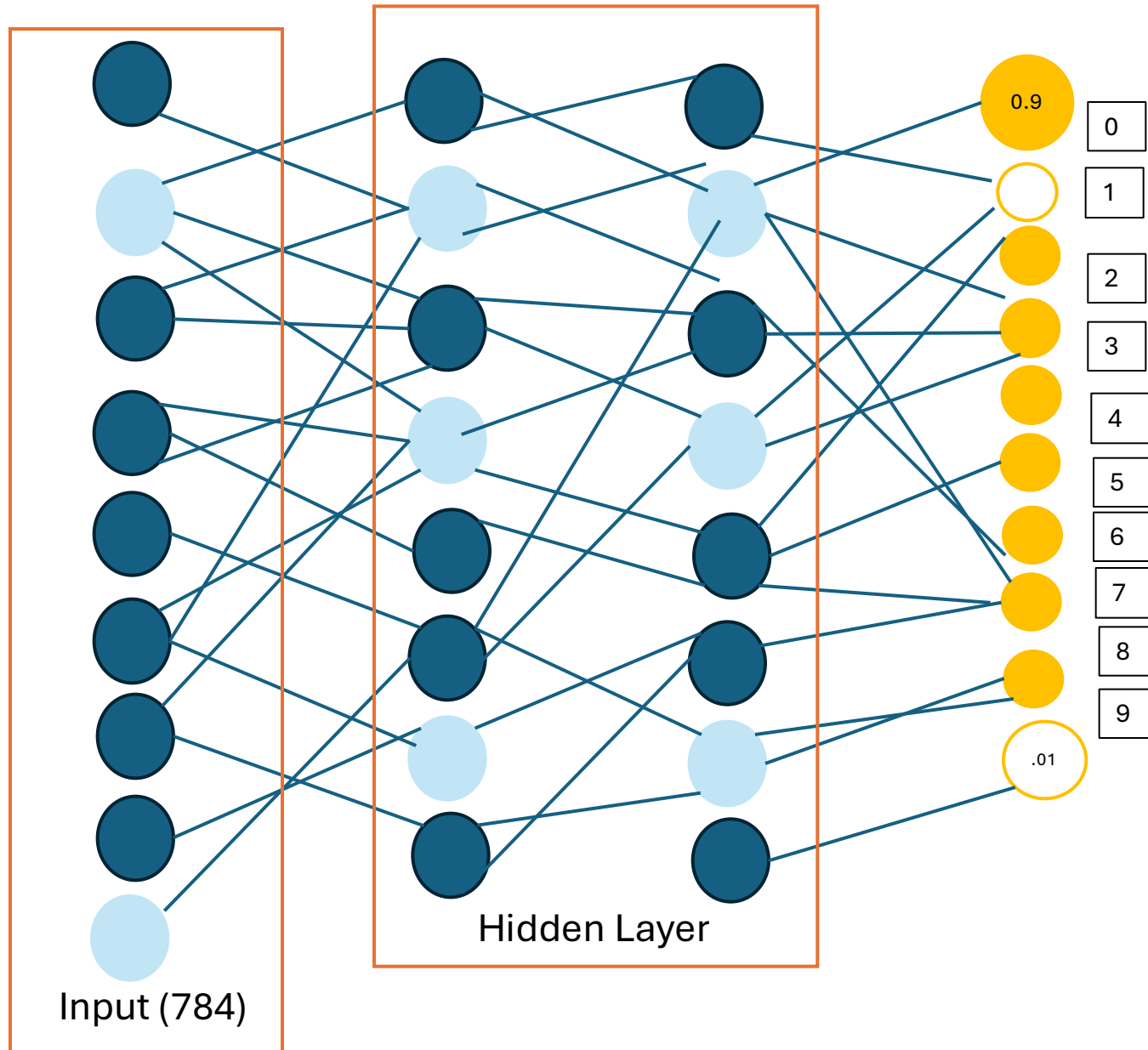
0.81

**Activation**

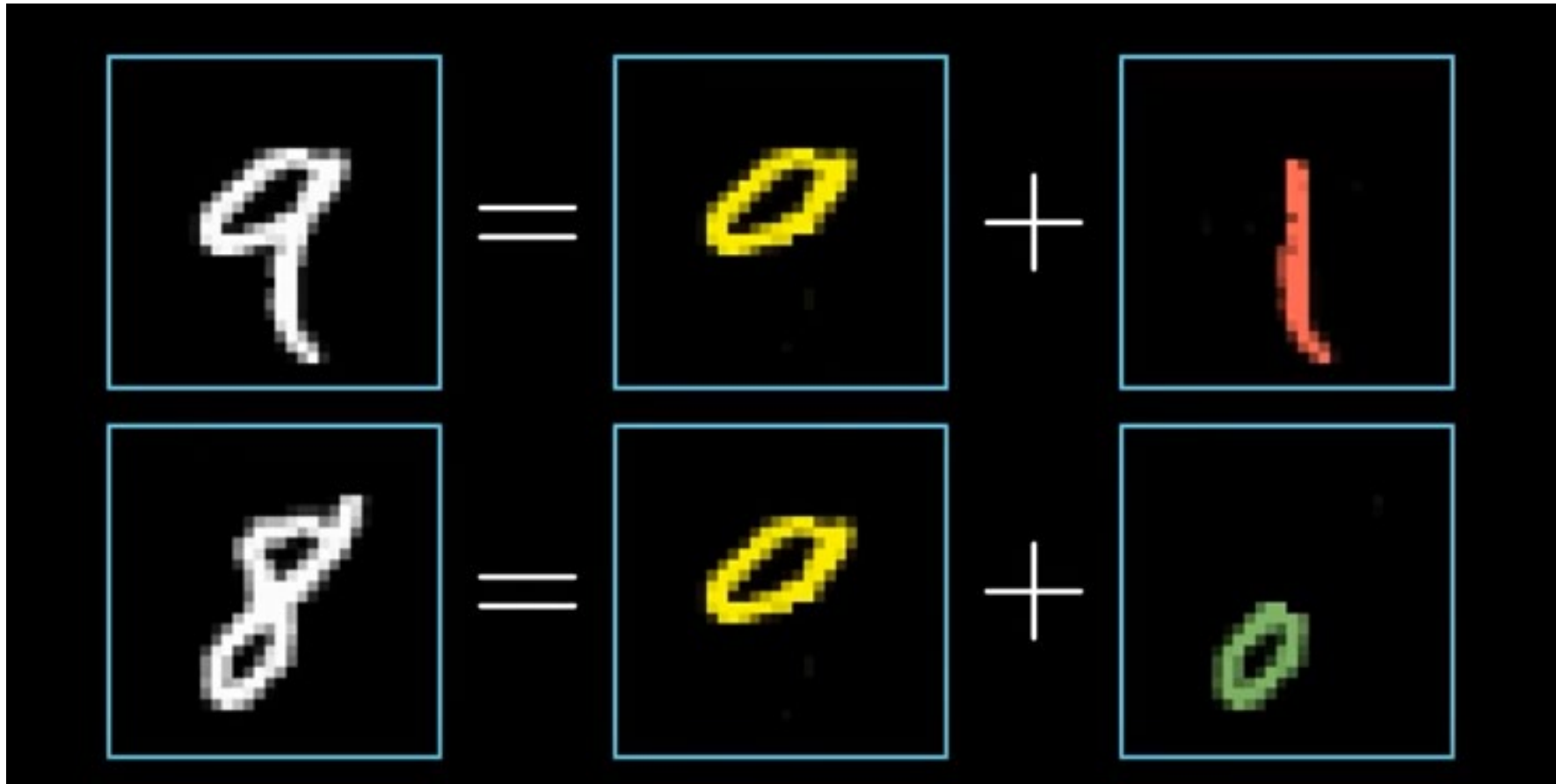
0.61

**Activation**

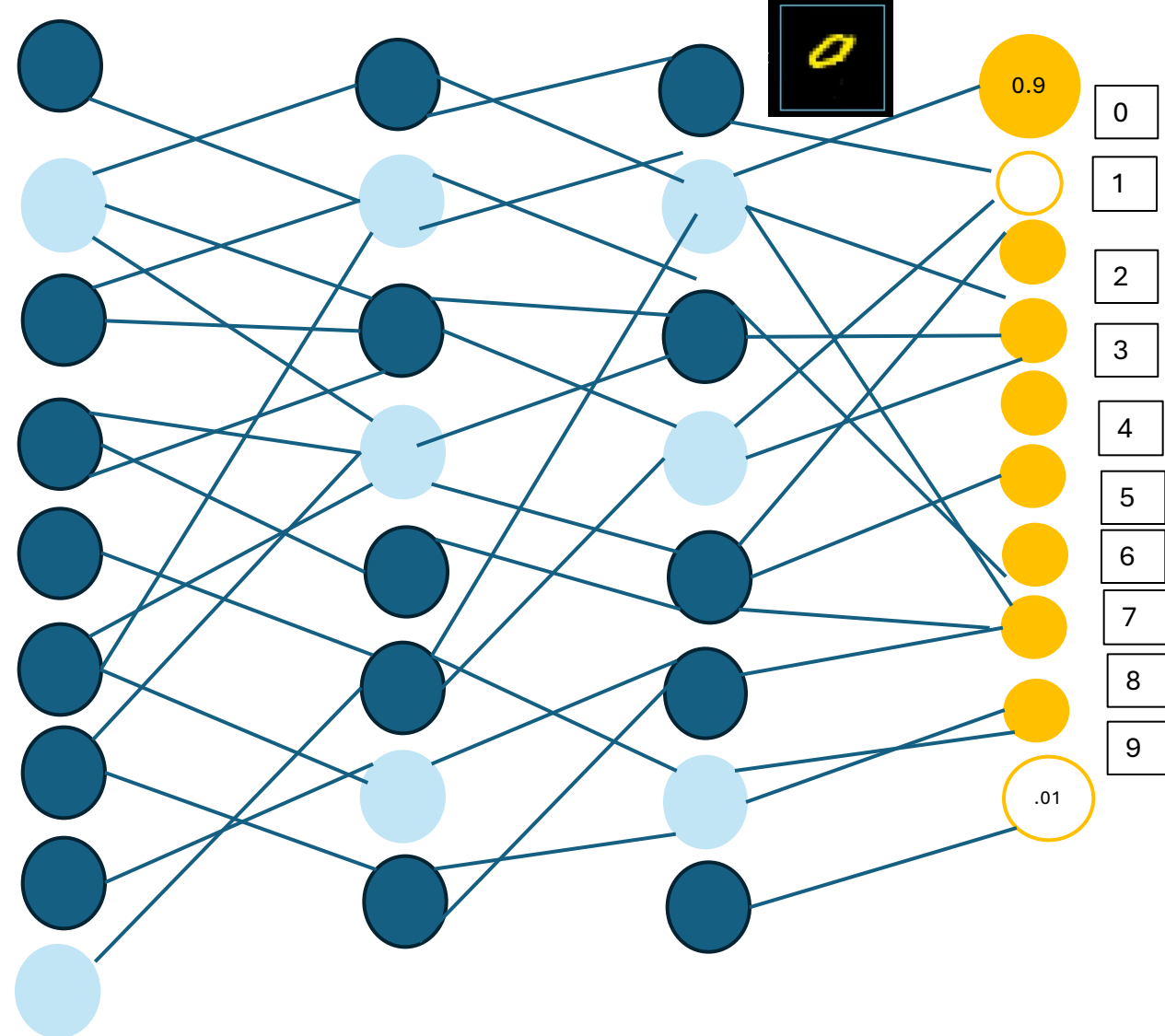
Activation in one layer influence activation in next layer



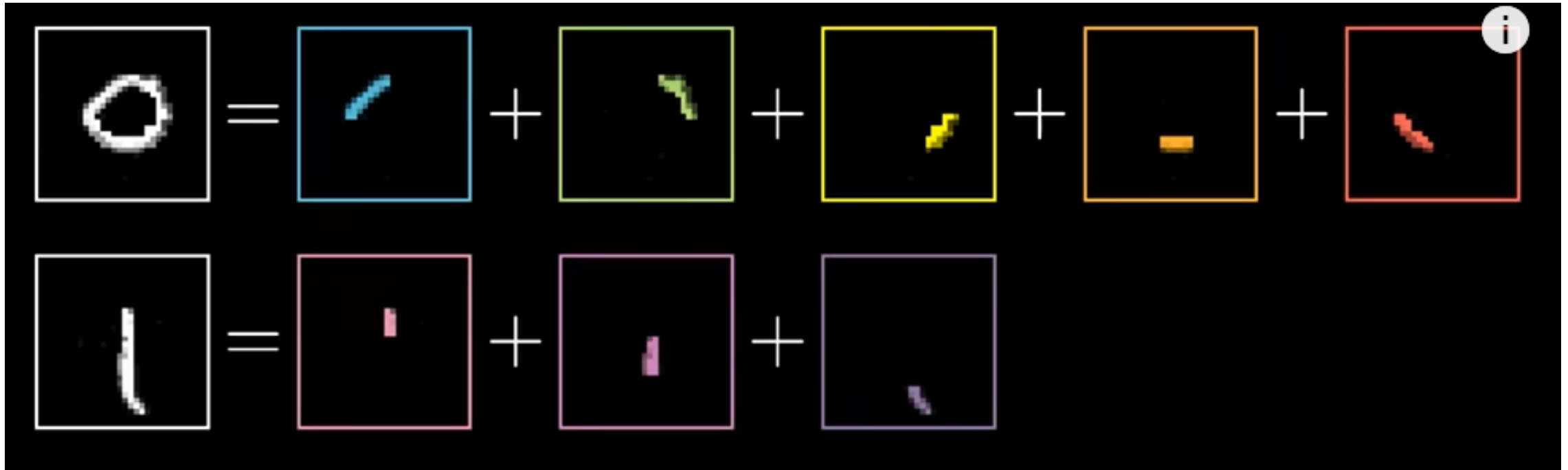
What are the middle layers doing



Activation in one layer influence activation in next layer



What are the middle layers doing





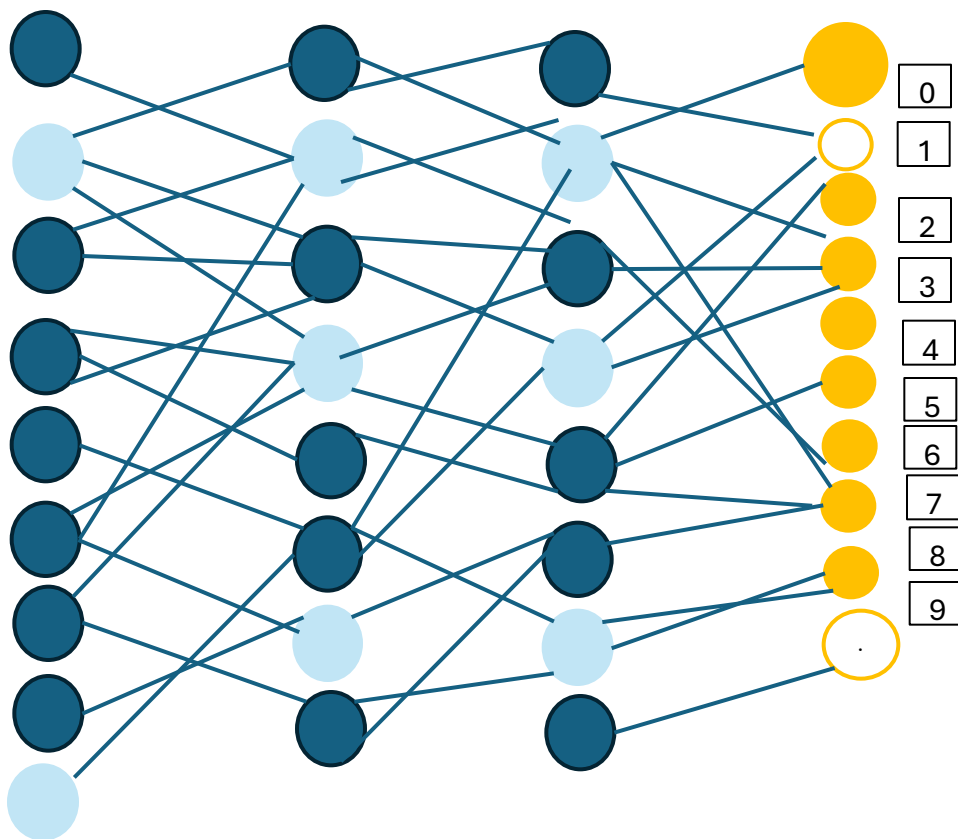
Picture

Edge

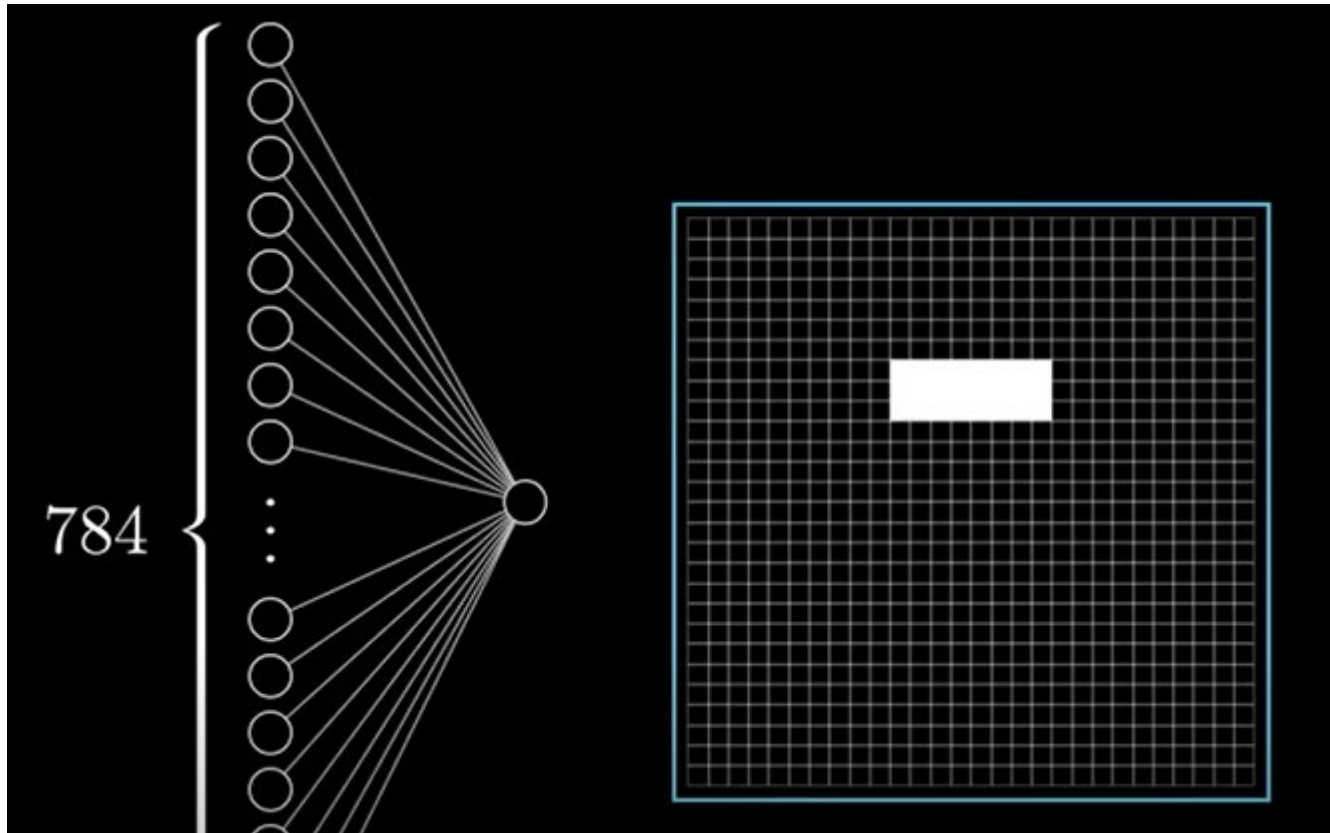
Pattern



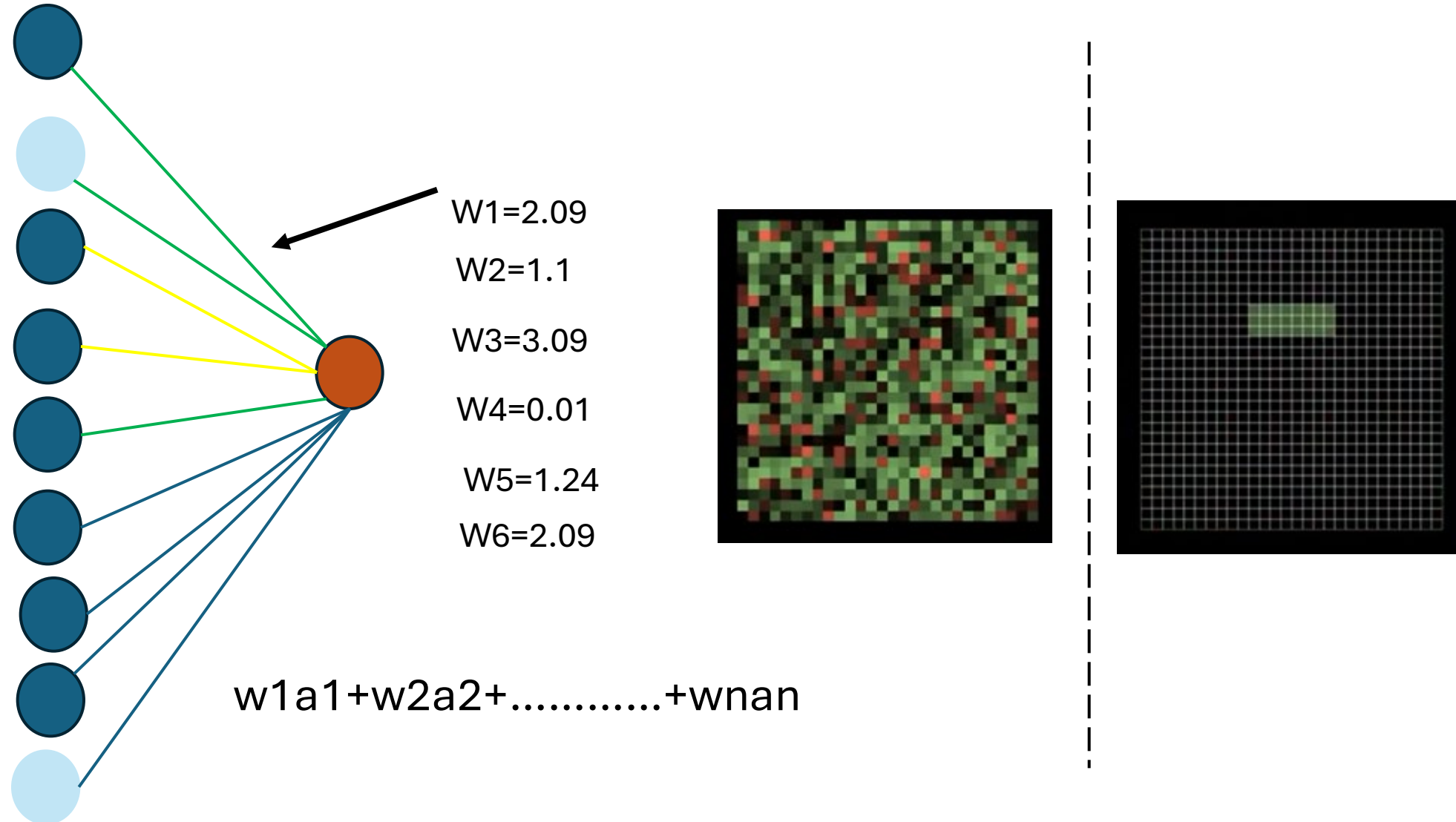
Lights up all the neurons  
that are recognizing the  
different regions



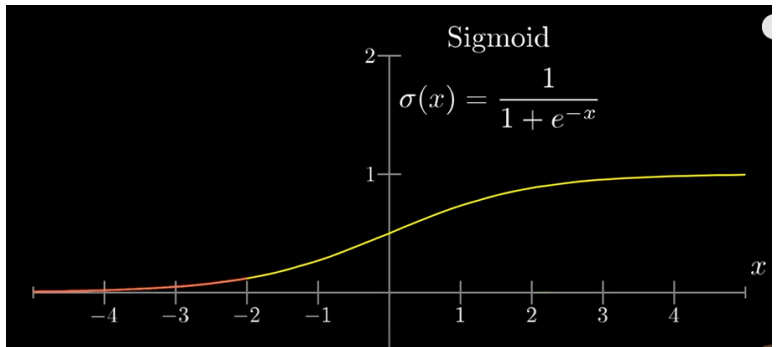
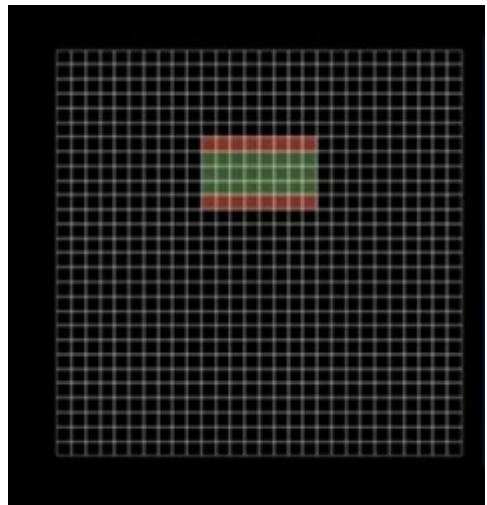
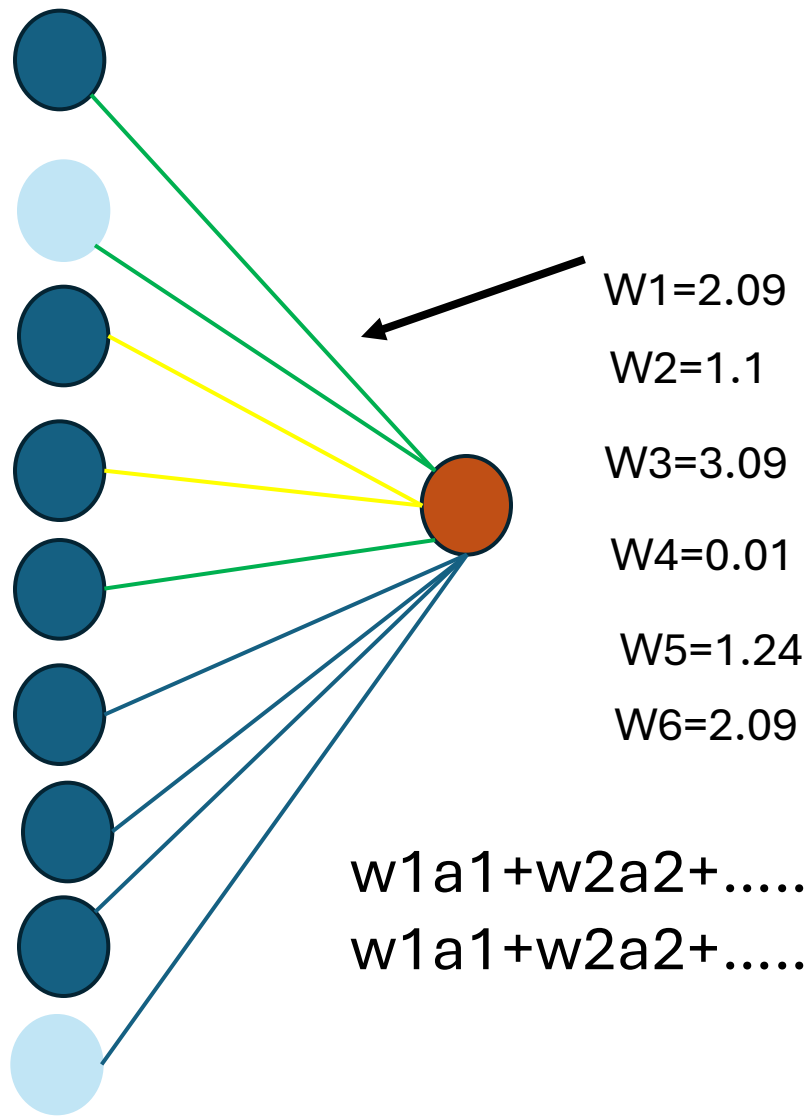
Whether a picture has an edge in this region ?  
What Parameters should the network have?



# Weights

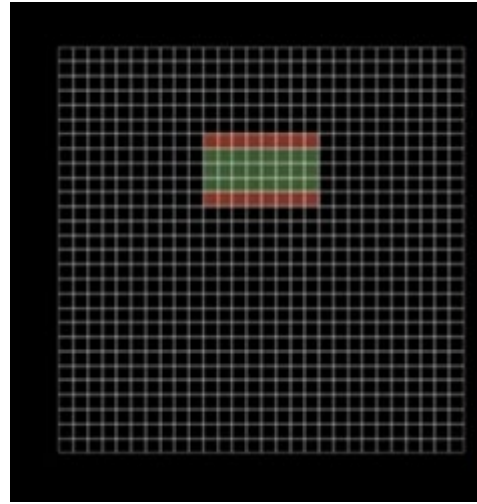
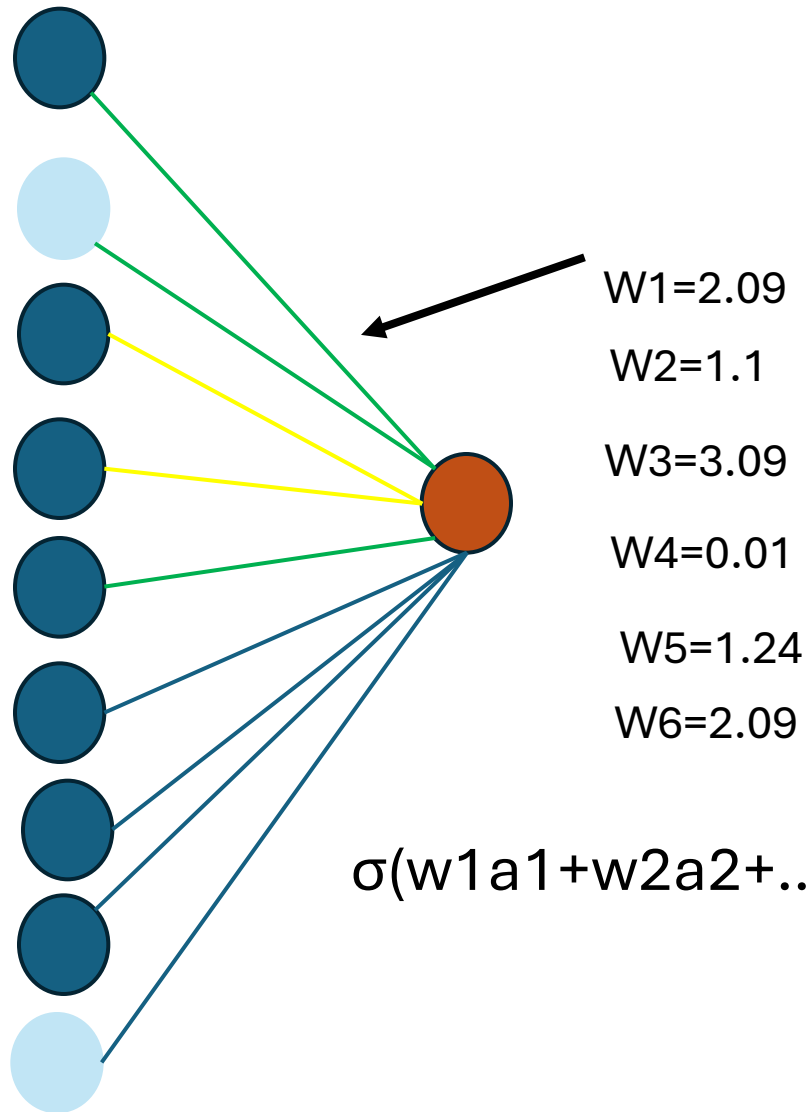


# Weights Continued ...



$w_1a_1 + w_2a_2 + \dots + w_n a_n$  max for +ve  
 $w_1a_1 + w_2a_2 + \dots + w_n a_n$  min for -ve } (0-1)

Bias

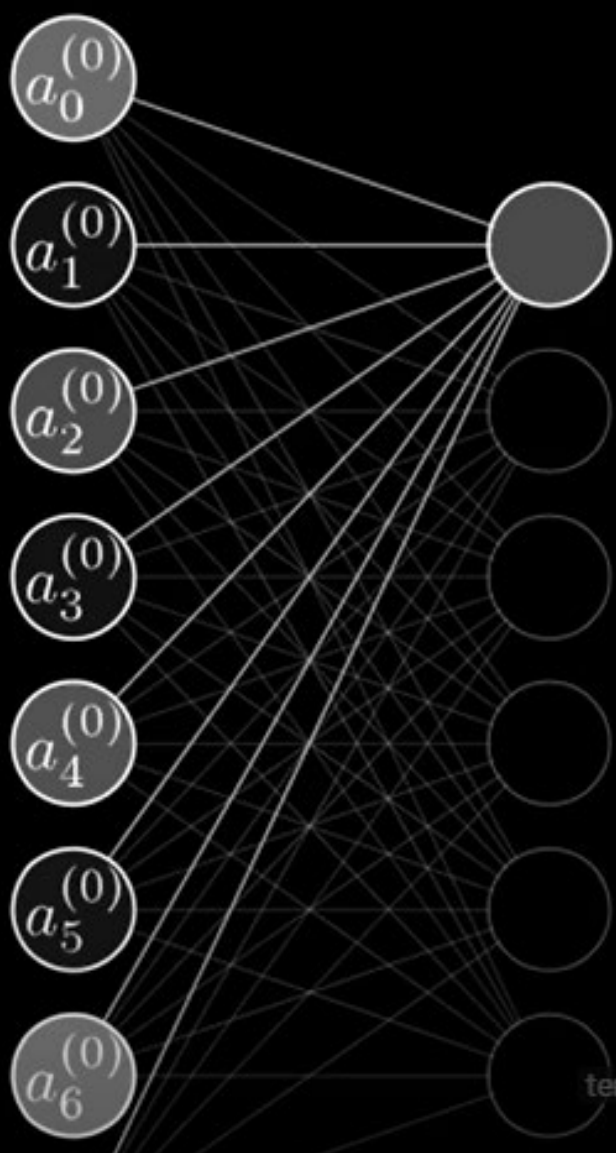


$$\sigma(w_1a_1+w_2a_2+\dots+w_n a_n)$$

But you only want to activate meaningfully when the weighted sum >10

$$\sigma(w_1a_1+w_2a_2+\dots+w_n a_n - 10)$$

*"Bias"*



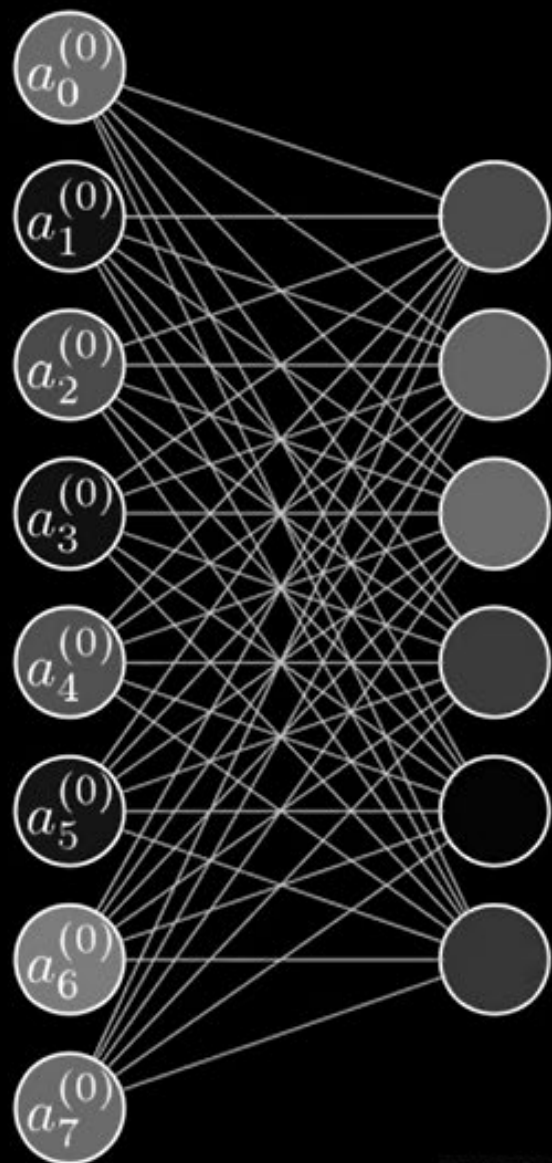
# Sigmoid

$$a_0^{(1)} = \sigma \left( \underbrace{w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \dots + w_{0,n} a_n^{(0)}}_{\text{Bias}} + b_0 \right)$$

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

terms in the matrix vector product of everything we have on the left here.





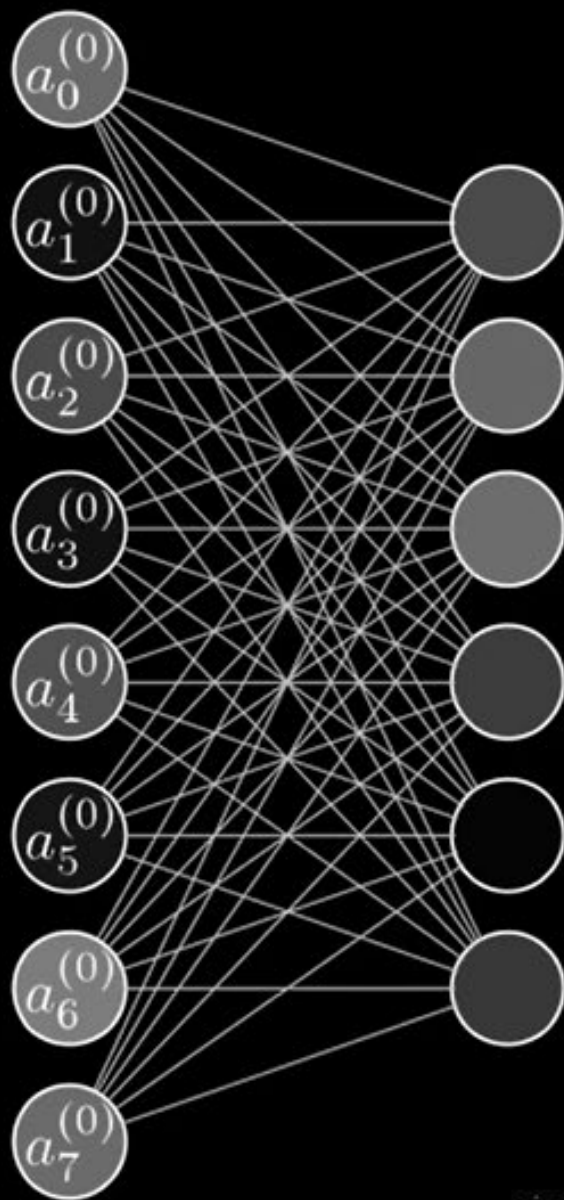
$$\sigma(\mathbf{W}\mathbf{a}^{(0)} + \mathbf{b})$$

$$\sigma \left( \begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$



you can communicate the full transition of activations from one layer to the next in an



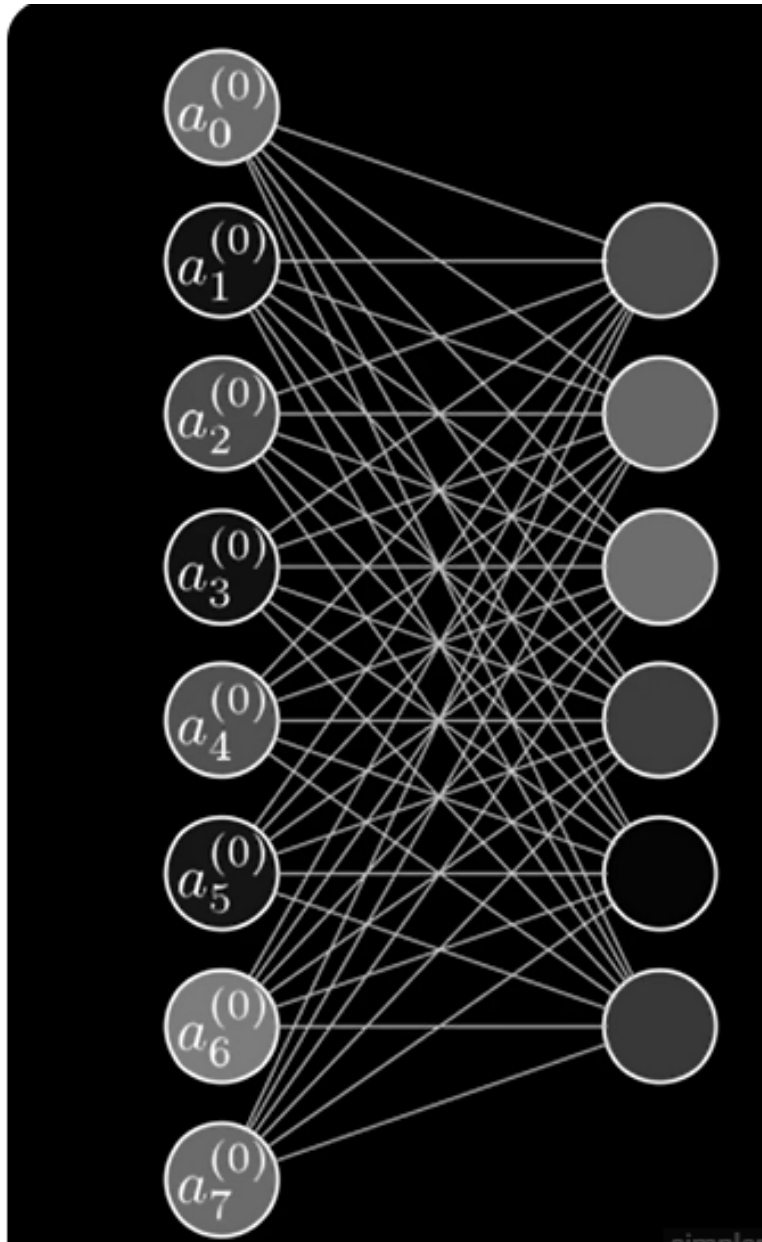


$$\mathbf{a}^{(1)} = \sigma(\mathbf{W}\mathbf{a}^{(0)} + \mathbf{b})$$

```
class Network(object):  
    def __init__(self, *args, **kwargs):  
        #...yada yada, initialize weights and biases...  
  
    def feedforward(self, a):  
        """Return the output of the network for an input vector a"""  
        for b, w in zip(self.biases, self.weights)
```



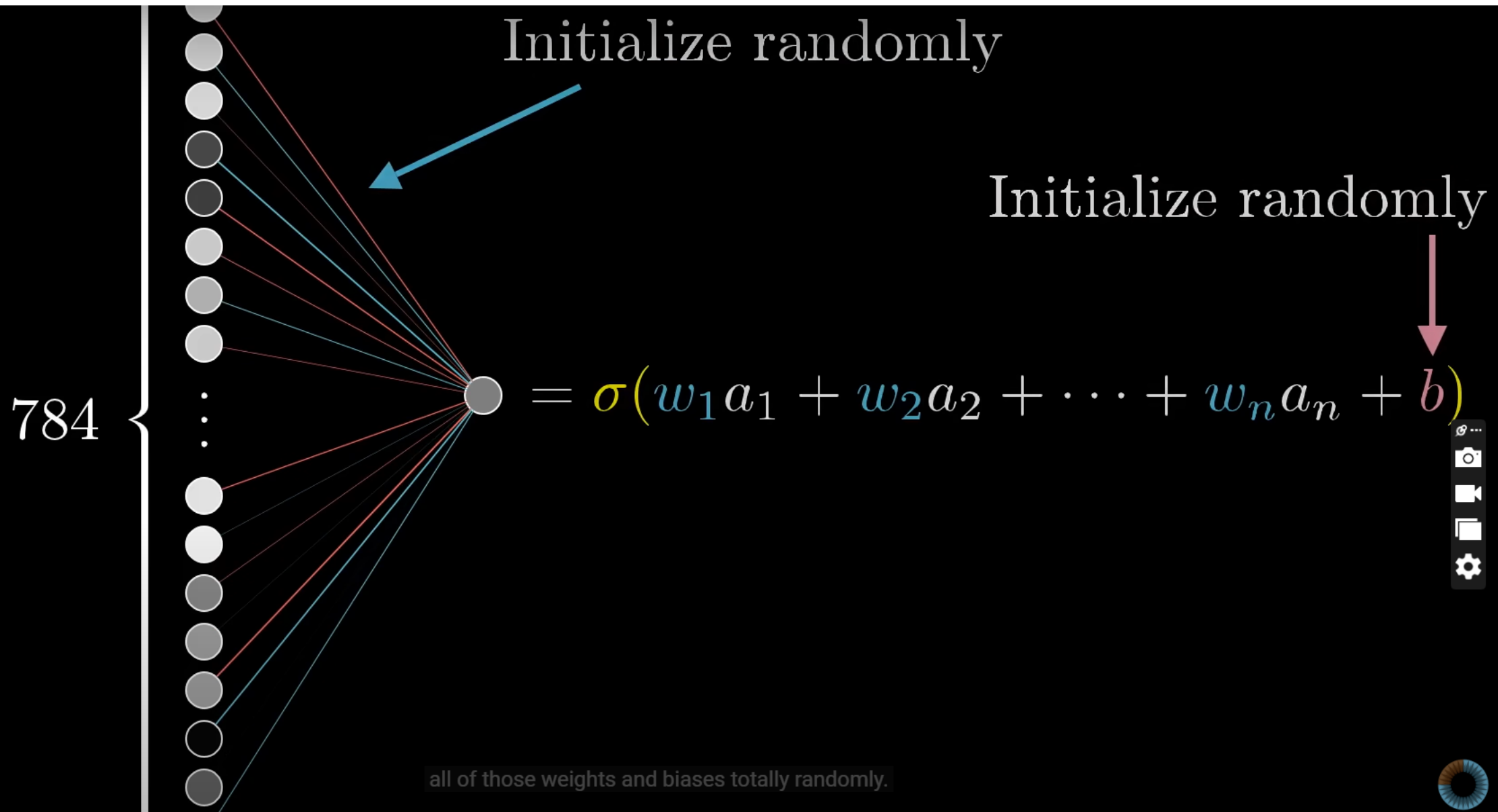




$784 \times 16 + 16 \times 16 + 16 \times 10$ : weights

$16 + 16 + 10$ : biases  $\sim 13000$

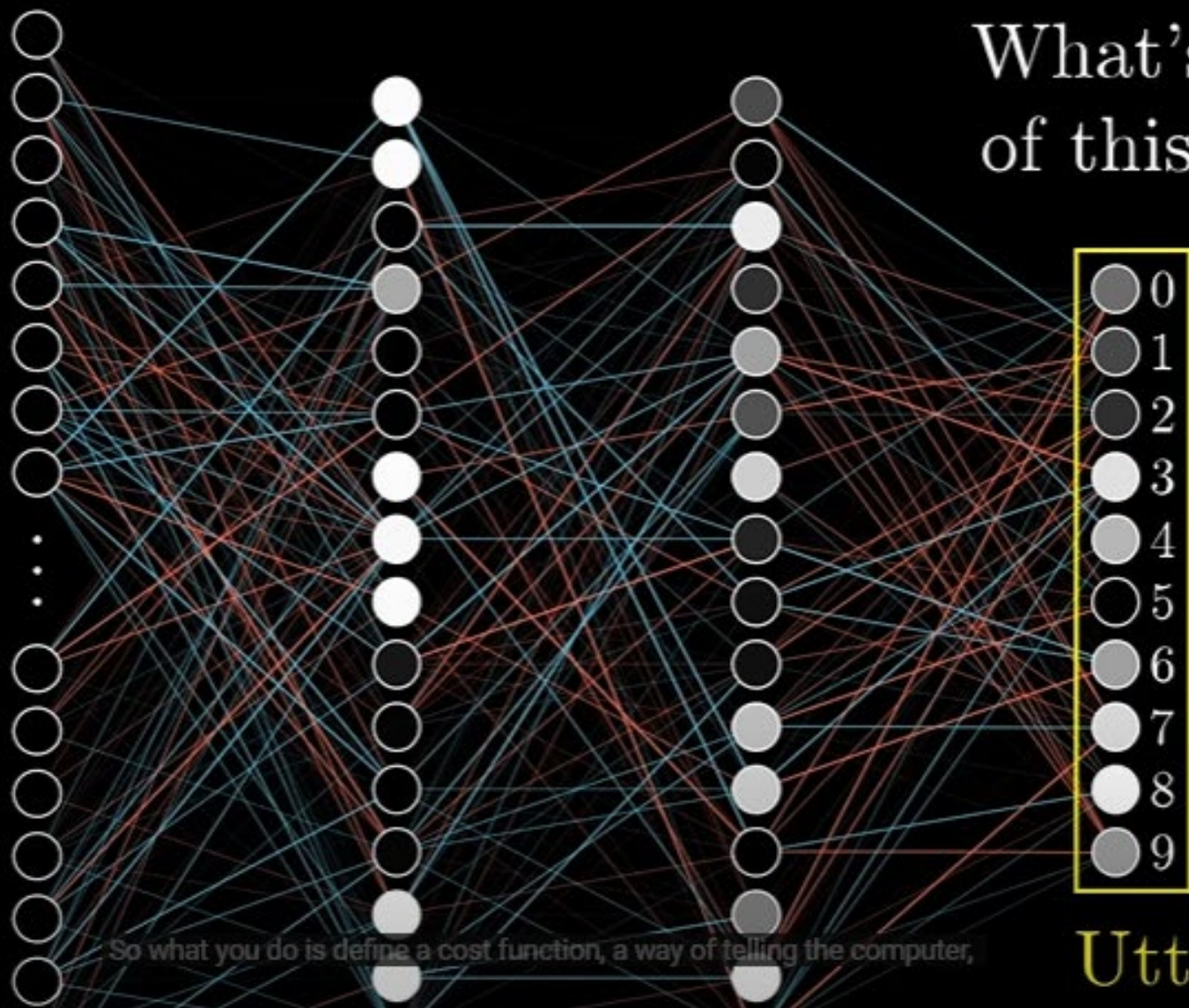
# How u train a NN?



# Cost Function



784



What's the “cost”<sup>i</sup> of this difference?

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

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



So what you do is define a cost function, a way of telling the computer,

Utter trash



Cost Function is small when the NN classify accurately

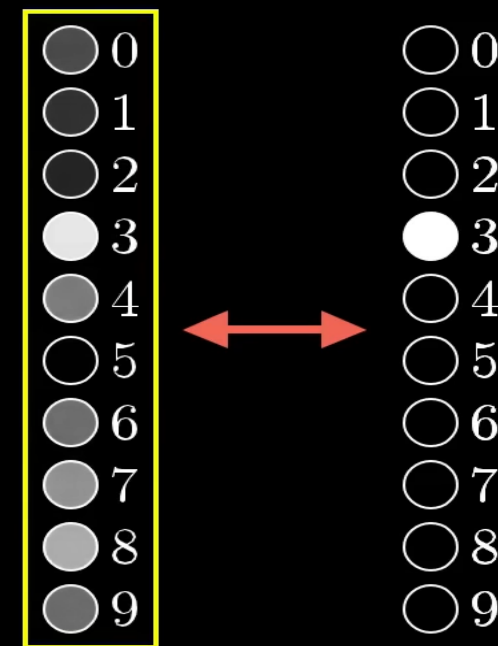
Cost of

3

1.61

$$\begin{aligned} &0.0865 \leftarrow (0.29 - 0.00)^2 + \\ &0.0390 \leftarrow (0.20 - 0.00)^2 + \\ &0.0214 \leftarrow (0.15 - 0.00)^2 + \\ &0.0077 \leftarrow (0.91 - 1.00)^2 + \\ &0.2550 \leftarrow (0.51 - 0.00)^2 + \\ &0.0002 \leftarrow (0.01 - 0.00)^2 + \\ &0.1930 \leftarrow (0.44 - 0.00)^2 + \\ &0.3472 \leftarrow (0.59 - 0.00)^2 + \\ &0.4698 \leftarrow (0.69 - 0.00)^2 + \\ &0.1879 \leftarrow (0.43 - 0.00)^2 \end{aligned}$$

What's the “cost” of this difference?

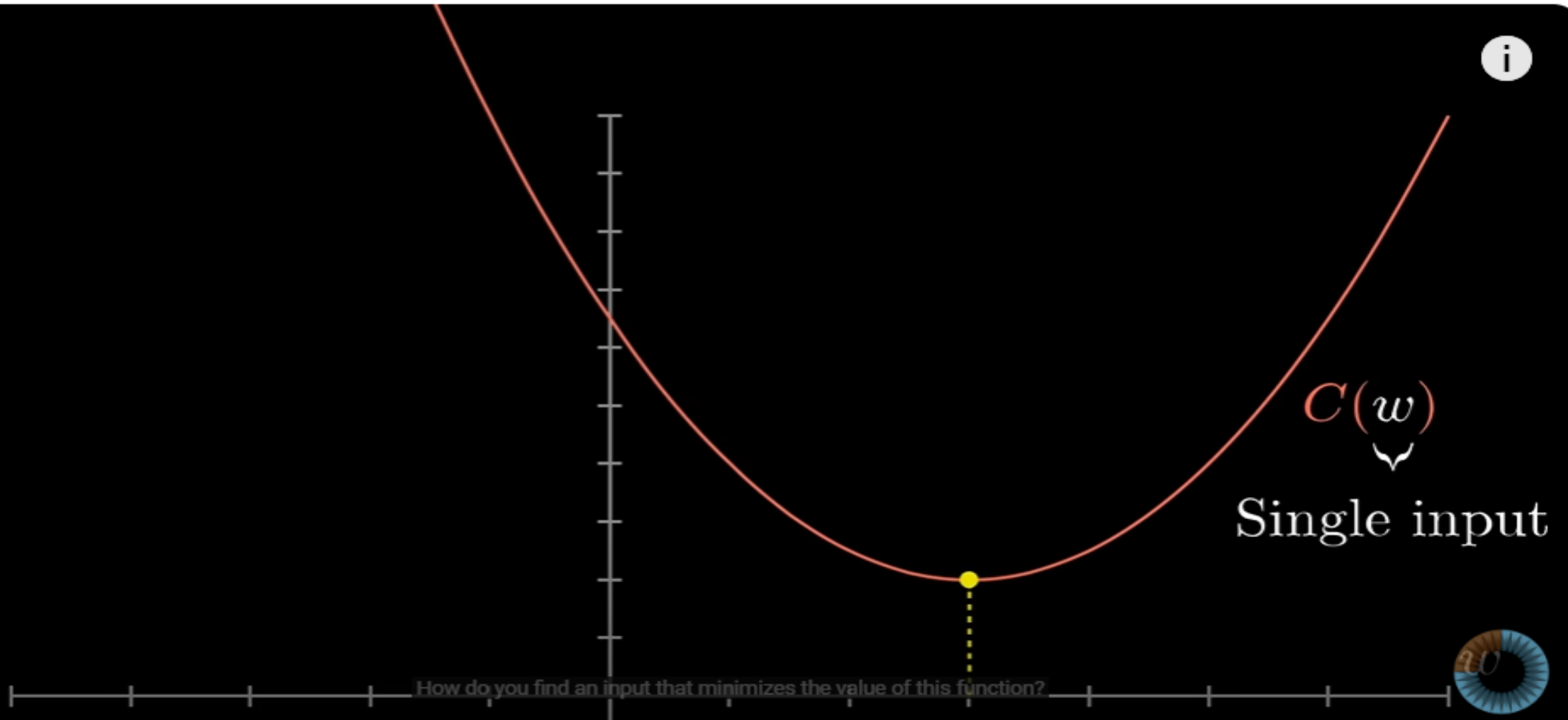


Utter trash

Notice this sum is small when the network confidently classifies the image correctly,

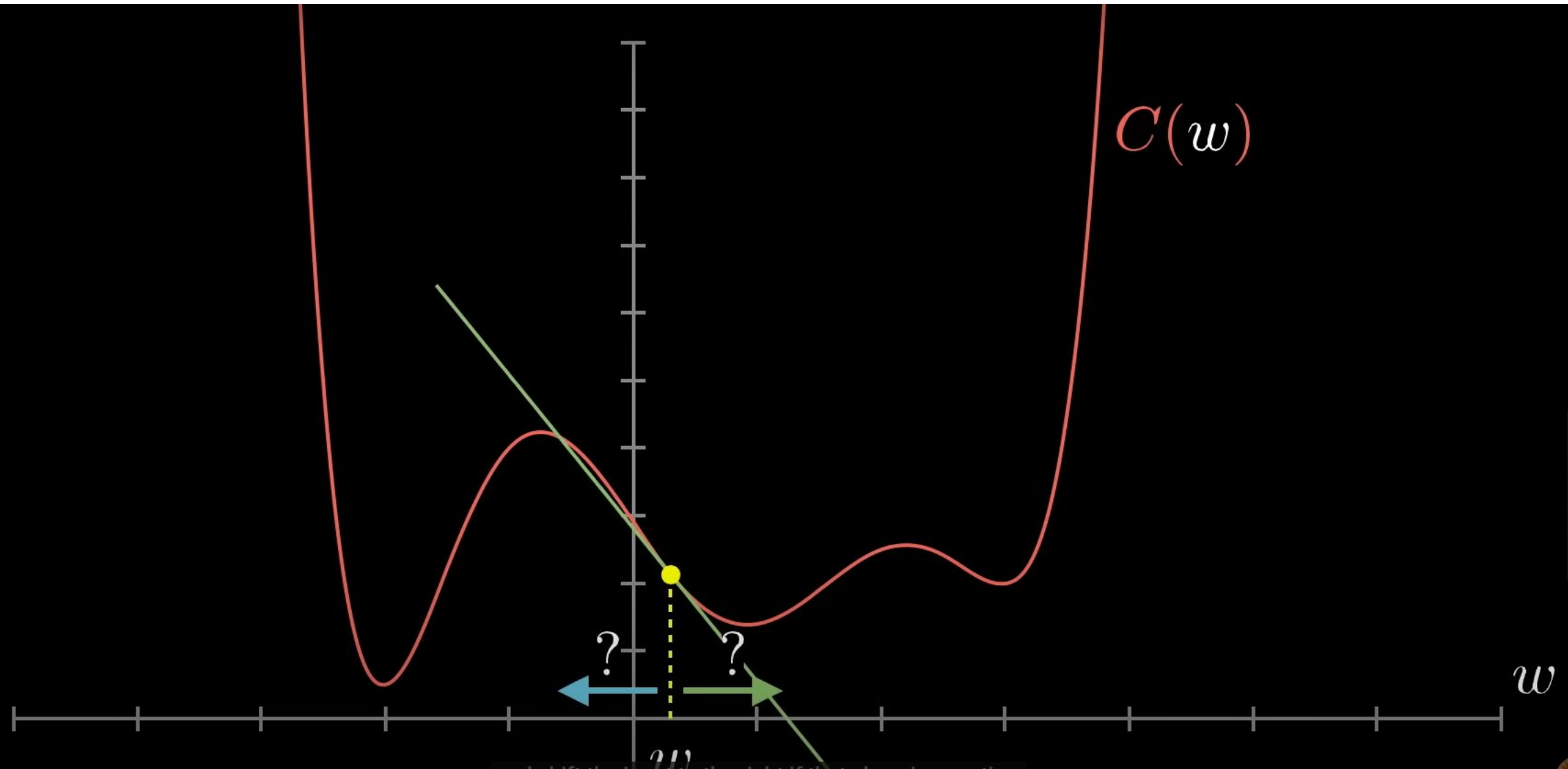


Let the cost function for all 13000 parameters be defined as  $C(w)$ . We want to minimize this fx.

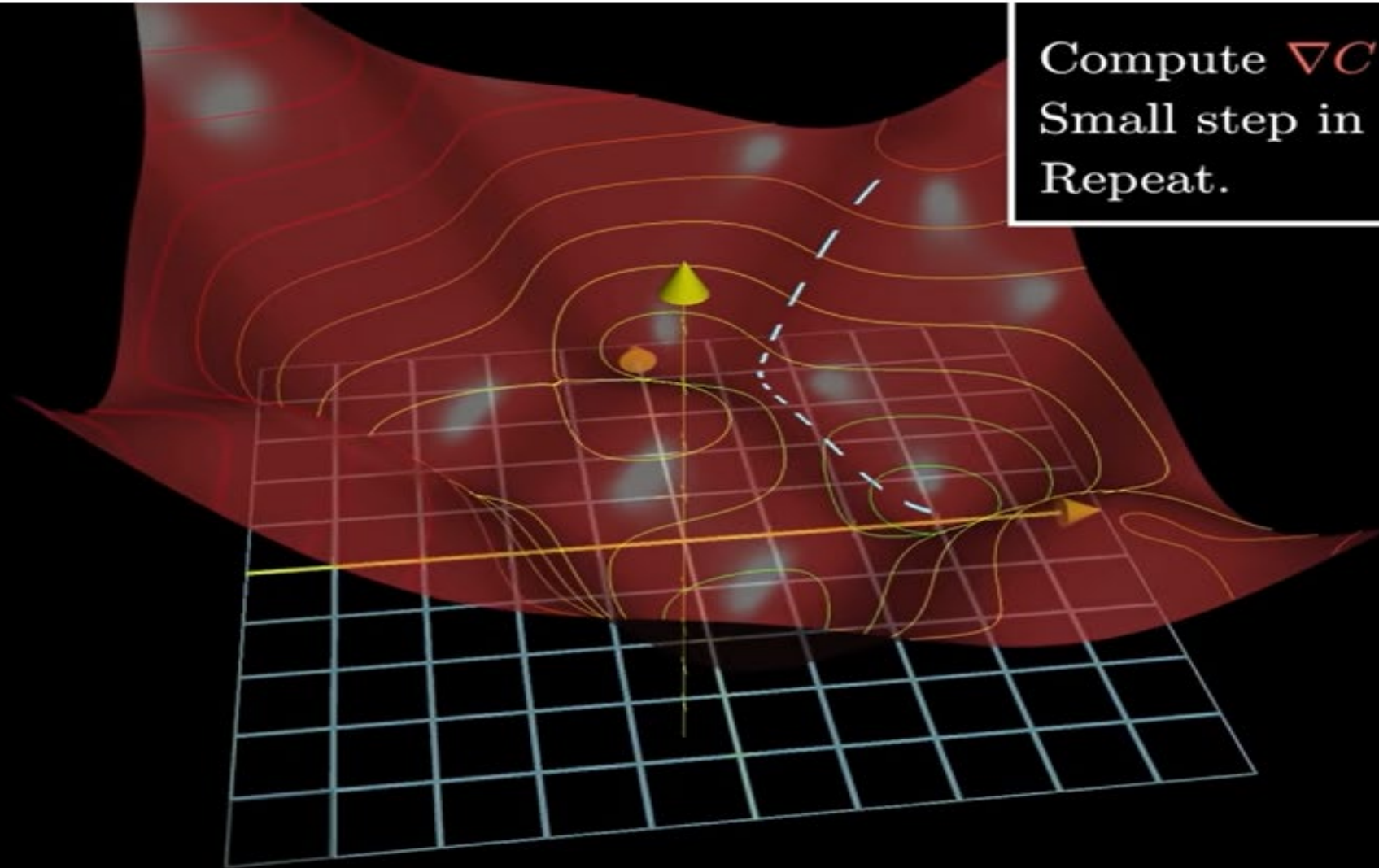




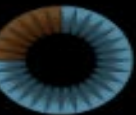
Find the slope of the  $C(w)$   $\rightarrow$  shift right (-ve slope)  $\rightarrow$  shift left (+ve slope)



Find the slope of the  $C(w)$   $\rightarrow$  shift right (-ve slope)  $\rightarrow$  shift left (+ve slope)



Compute  $\nabla C$   
Small step in  $-\nabla C$  direction  
Repeat.



Which nudges will rapidly decrease the cost function

i

13,002 weights and biases

How to nudge all  
weights and biases

$$\vec{\mathbf{W}} = \begin{bmatrix} 2.43 \\ -1.12 \\ 1.47 \\ \vdots \\ -0.76 \\ 3.50 \\ 2.03 \end{bmatrix}$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$





Remember !!

1. Cost function is average on all the training samples
2. When network is learning it's minimizing a cost function.
3. The process of nudging an input with +ve gradient of cost function is *Gradient Descent*. *It's a way to converge to local minima of cost function.*

$$\vec{W} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{13,000} \\ w_{13,001} \\ w_{13,002} \end{bmatrix}$$

To exit full screen, press Esc

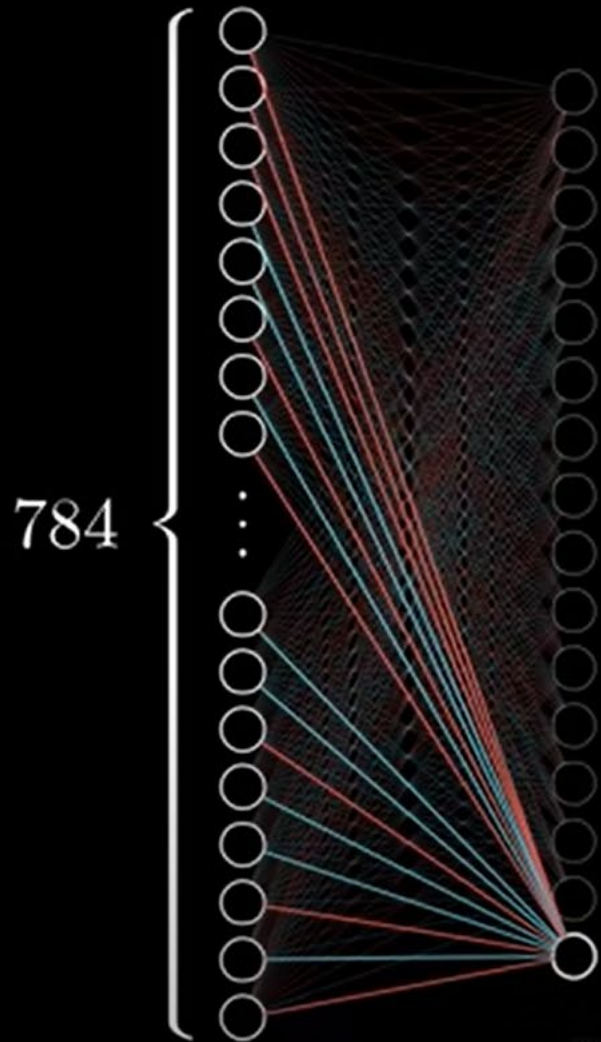
$$-\nabla C(\vec{W}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix}$$

$w_0$  should increase somewhat  
 $w_1$  should increase a little  
 $w_2$  should decrease a lot  
 $w_{13,000}$  should increase a lot  
 $w_{13,001}$  should decrease somewhat  
 $w_{13,002}$  should increase a little

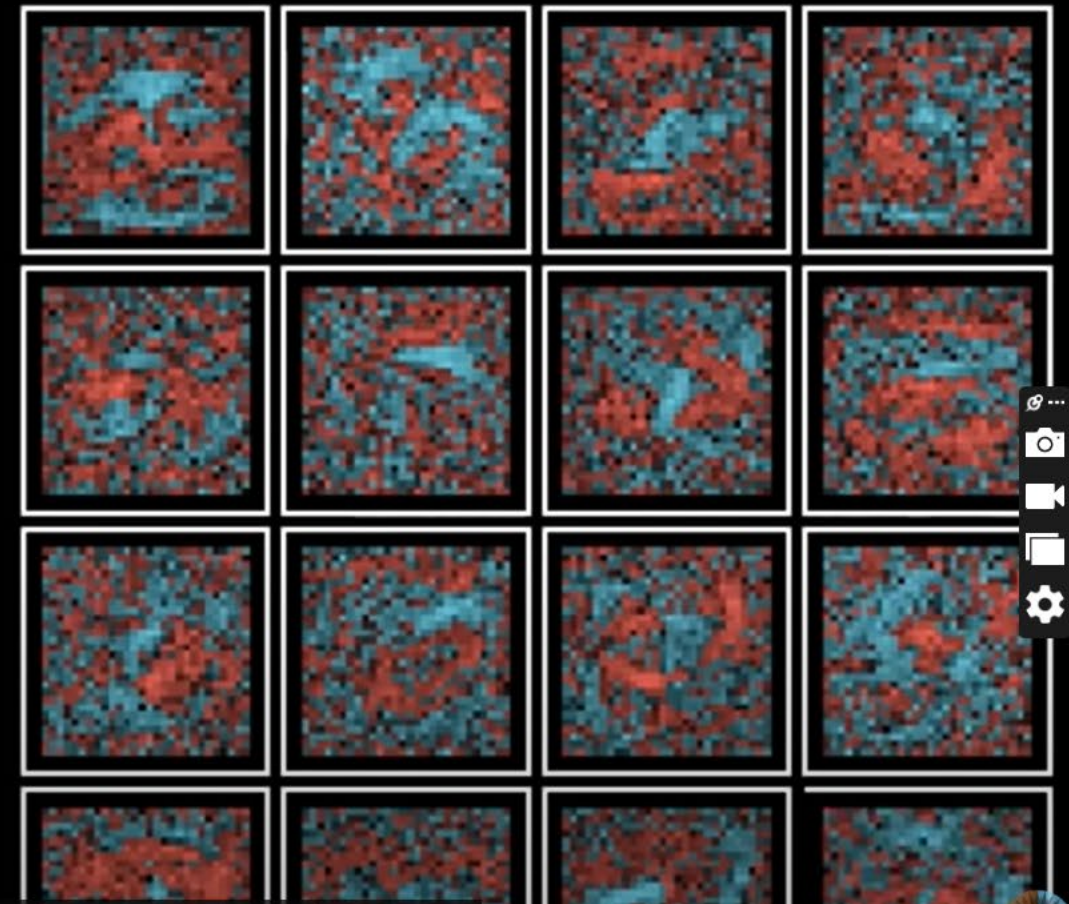
- Signs mentioned if the input vector should be nudged up (+ve) or down (-ve).
- Relative magnitude mentions which changes matter more.

*In short, which changes to which weights matters the most*

So, is my Multilayer Perceptron learning edges and patterns? No, its learning some random patterns based on finding the local minima.

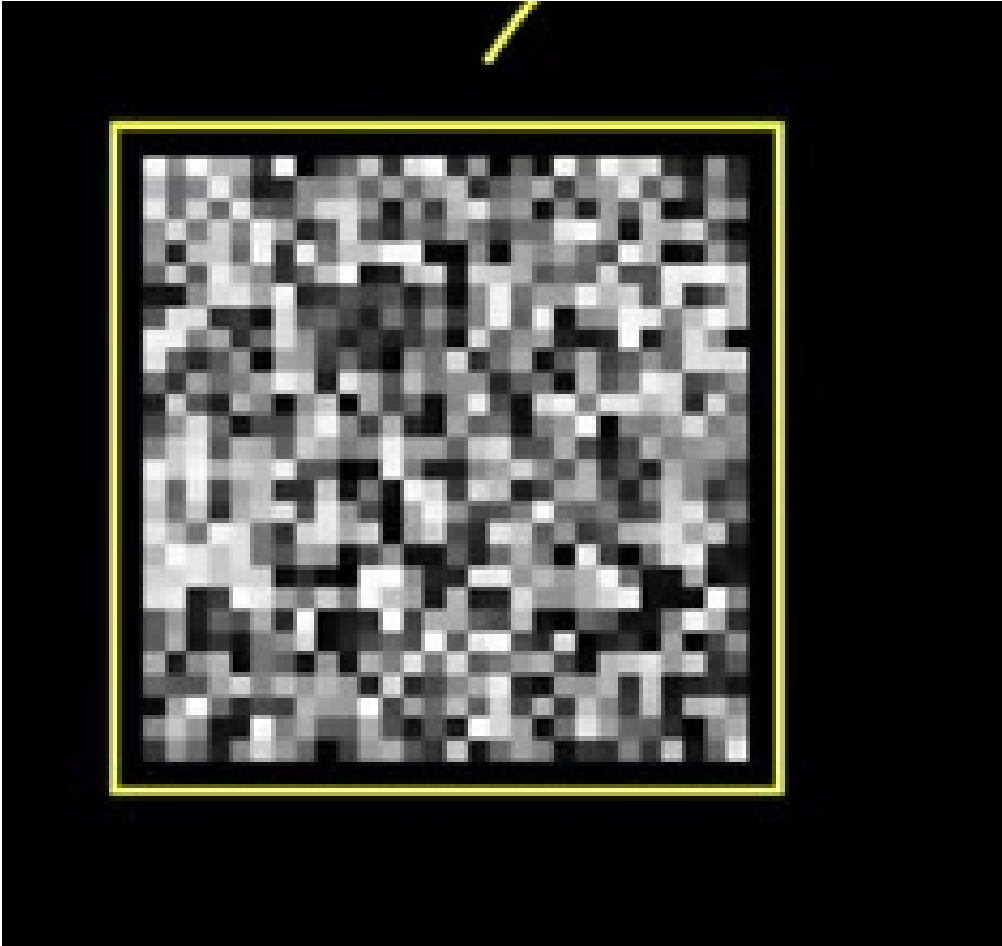


What second layer  
neurons look for



And to really drive this point home, watch what happens when you input a random image.

What will happen if you input a random image ?



Confidently gives you random answer because the network can only differentiate it has no idea how to draw a digit?