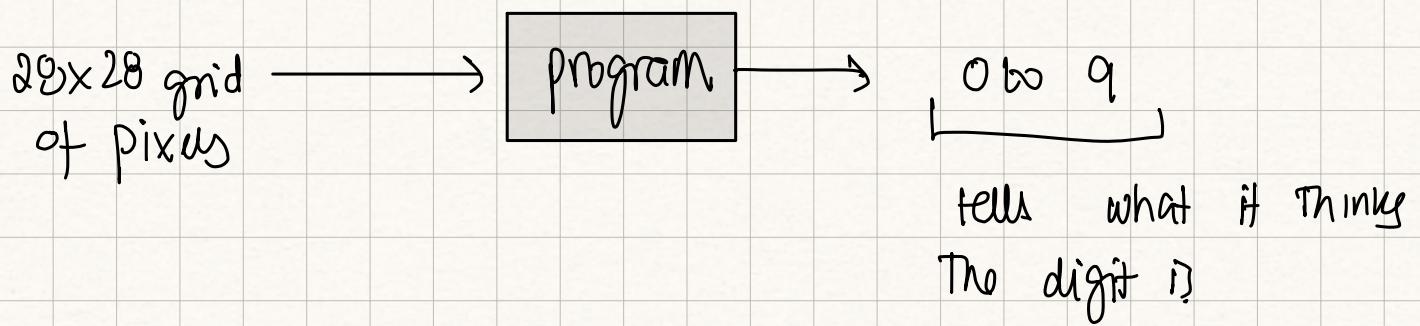


# DL 01 Notes : But what is a NN ?

- Brain has no problem recognising  $28 \times 28$  pixel imgs of 3.
- light sensitive cells that are firing when seeing each img are diff. but visual cortex resolves as it representing the exact same idea!
- Impossible to write a program that



## Goal for 1st 2 Videos :

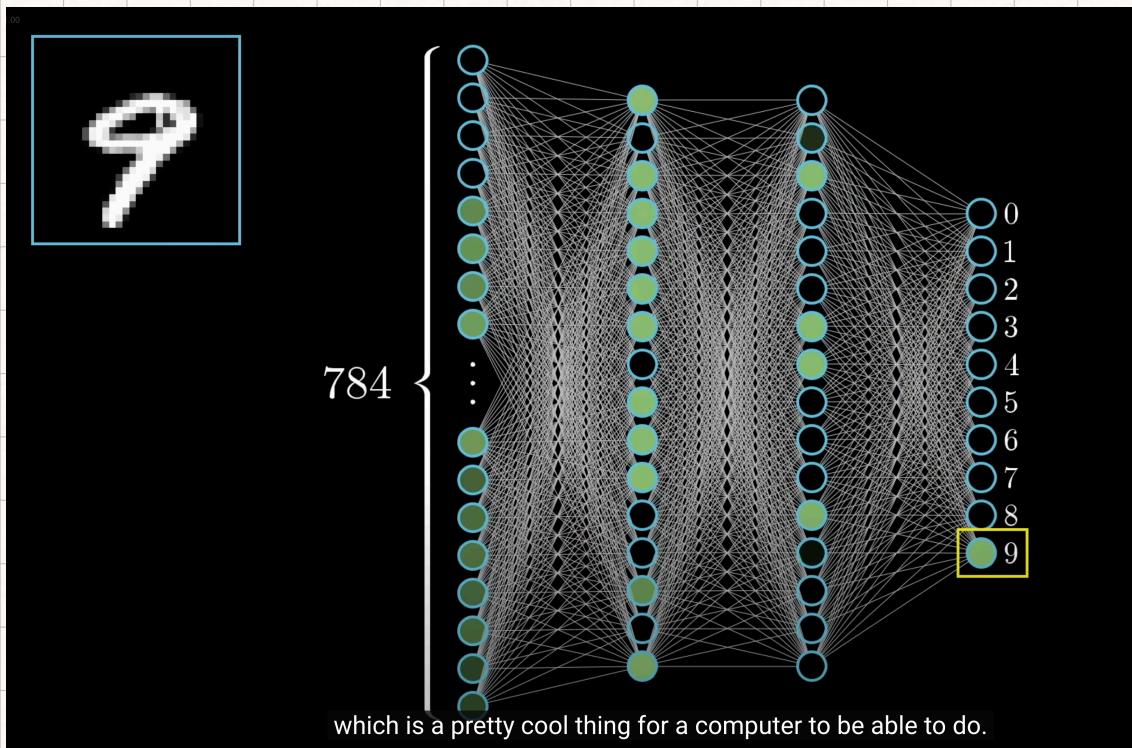
- What a NN actually is
- Visualise what it does -
- The structure of a NN is motivated

# o intuition for a NN "learning"

- NN to recognise handwritten digits (MNIST)

Many variants

- CNN - img recognition
- LSTM - speech recognition.

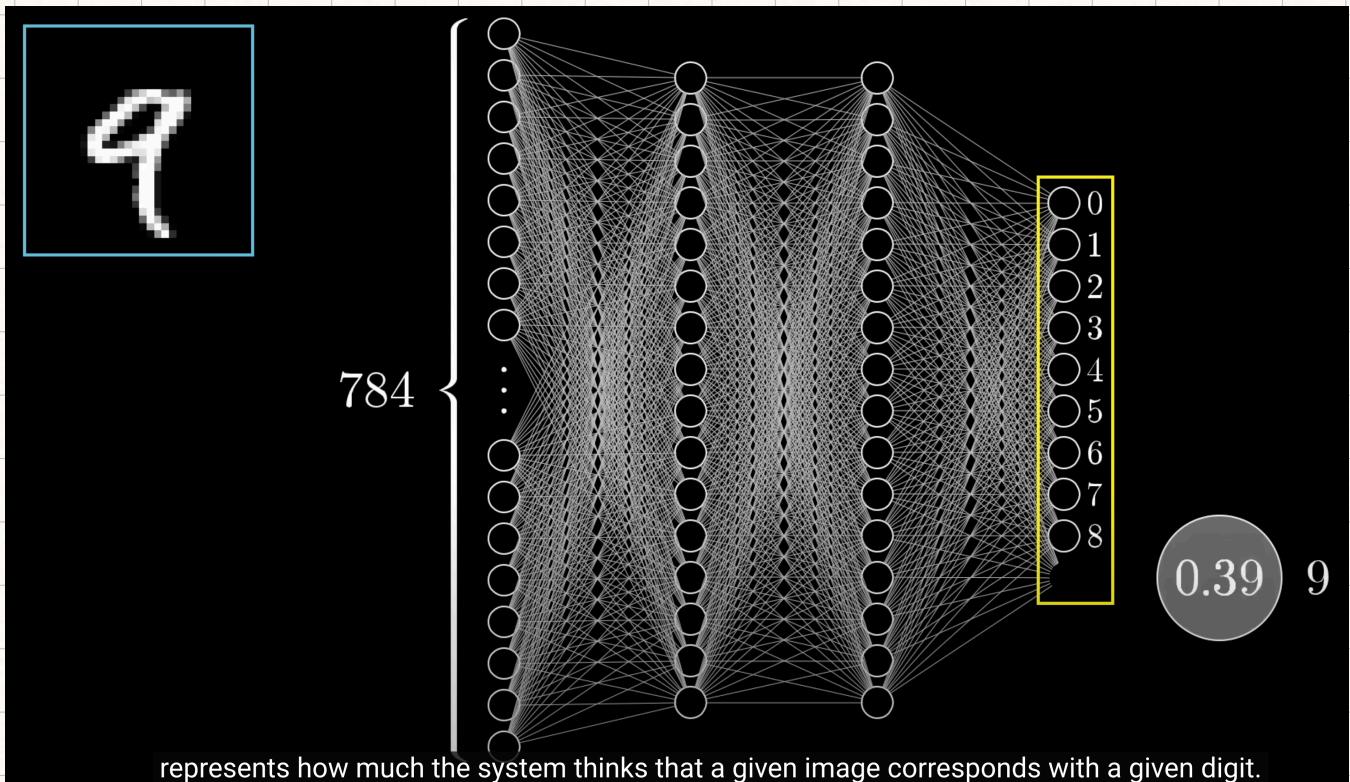
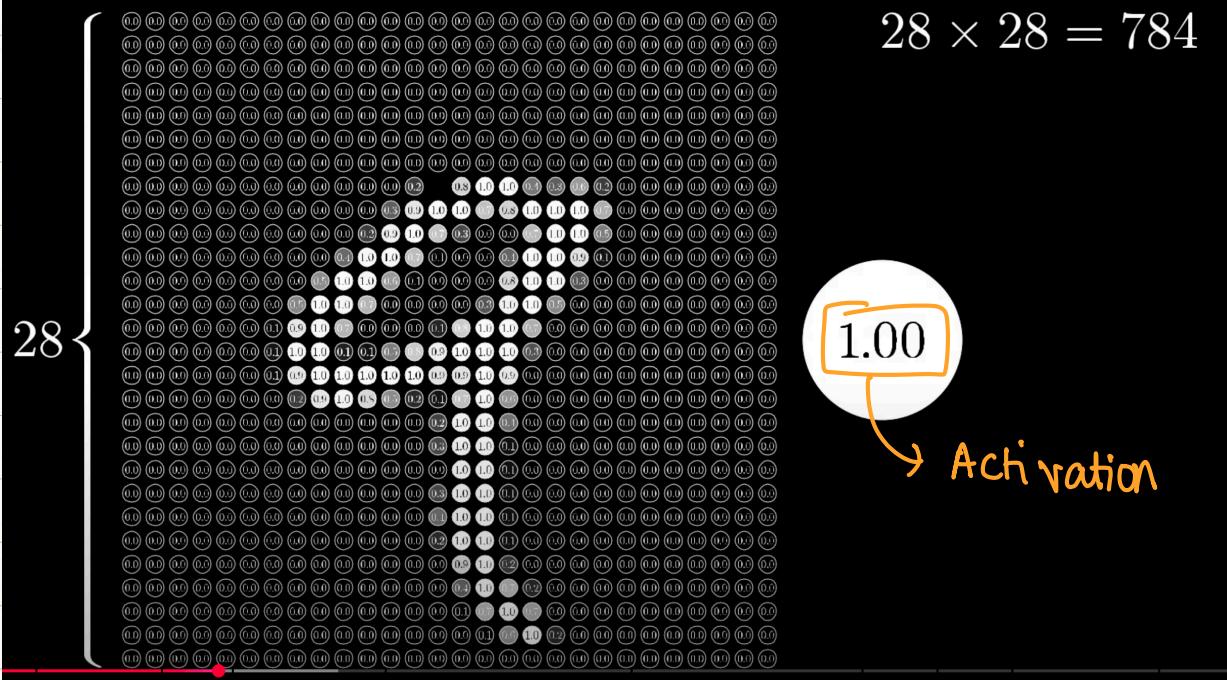


Neuron

0.2

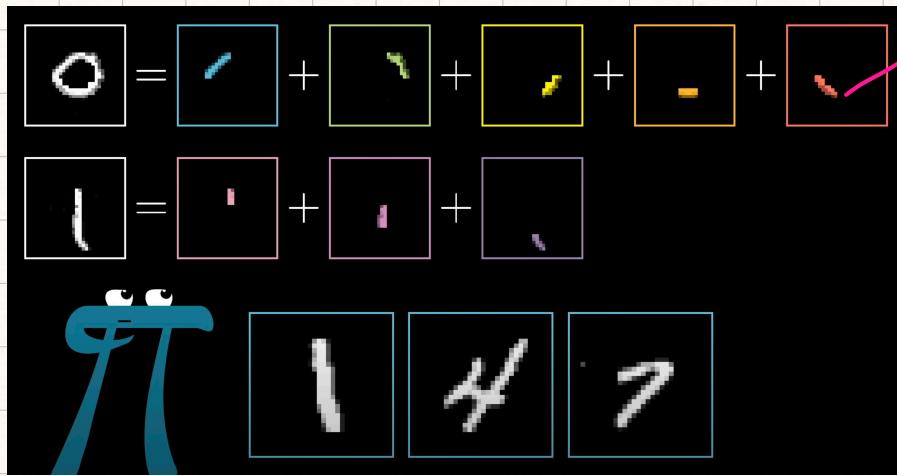
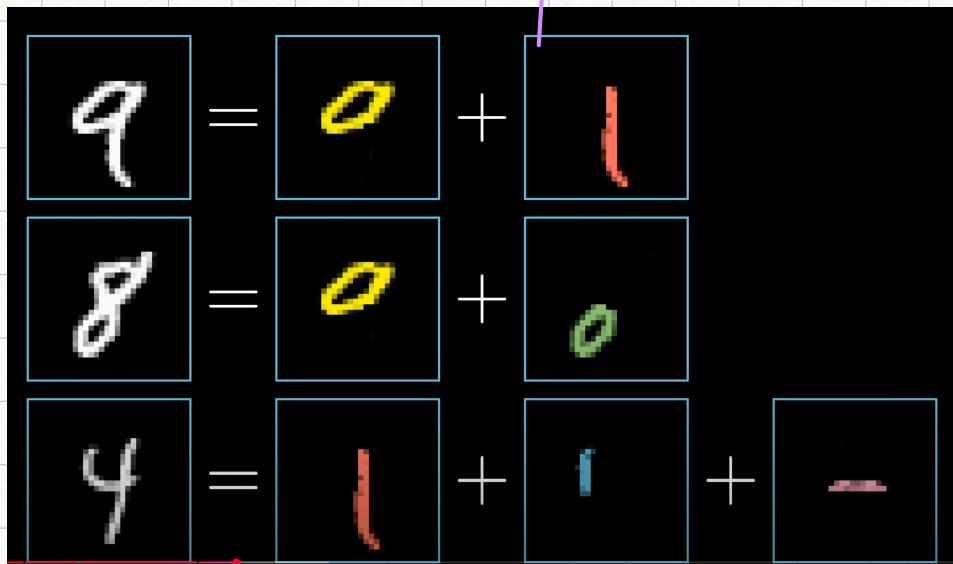
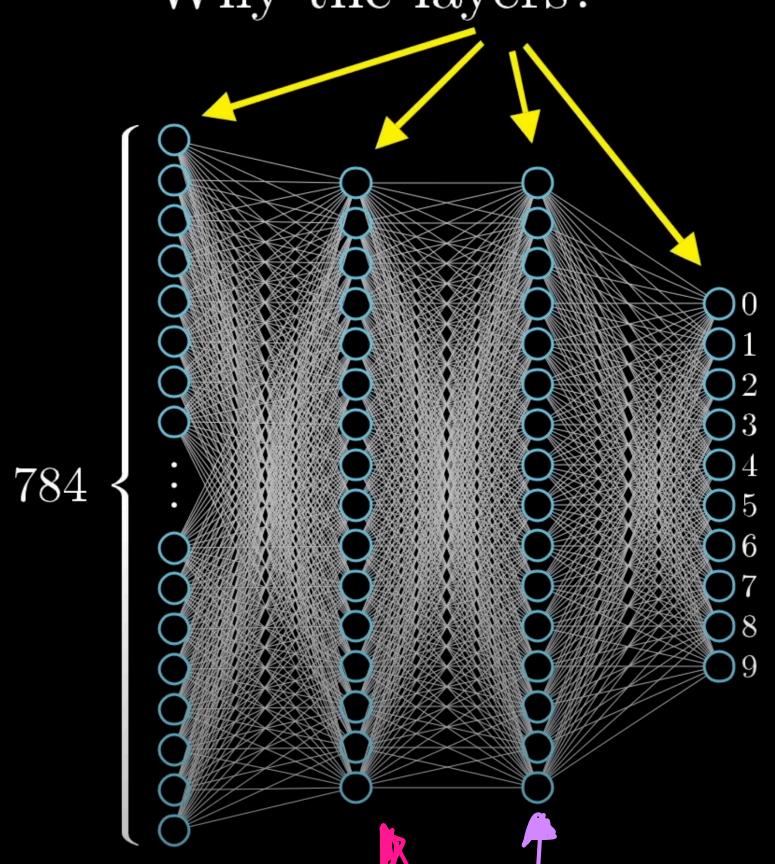
Neuron → Thing that holds a number

$$28 \times 28 = 784$$



- each input img causes some neurons in the 1st layer to light up, which lights up other neurons & so on...

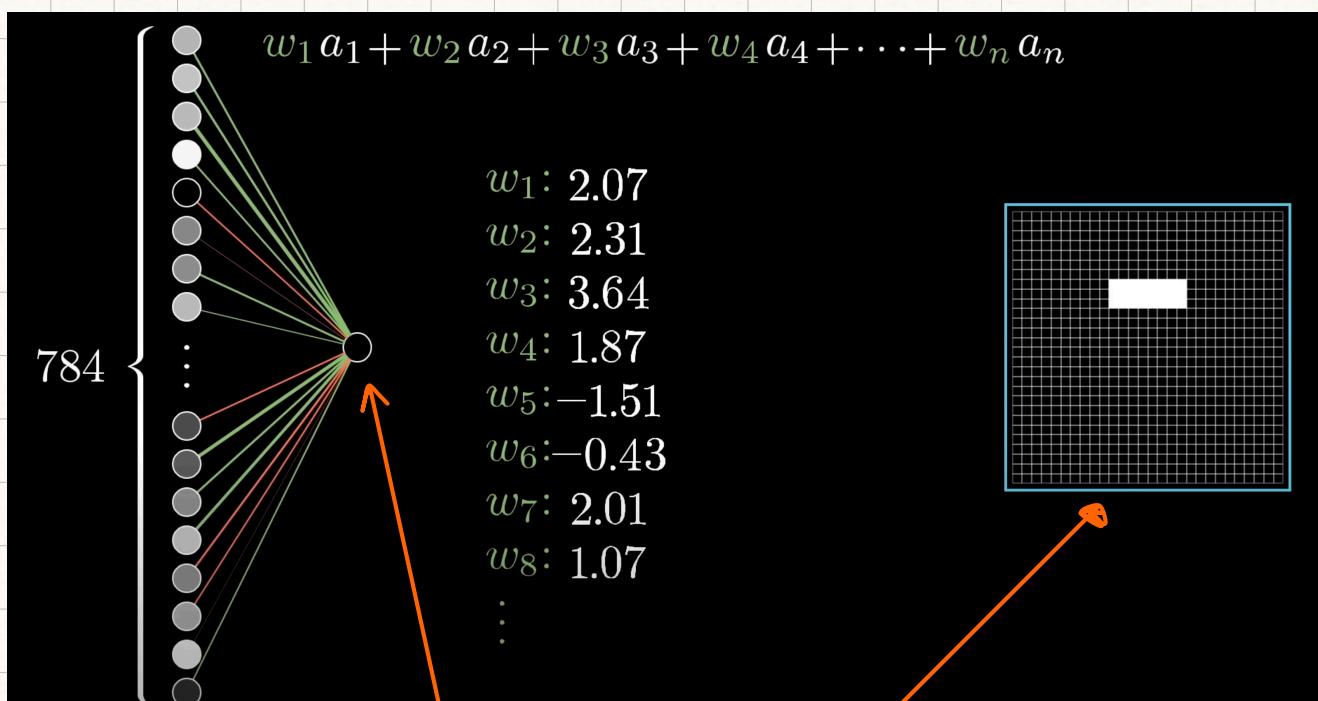
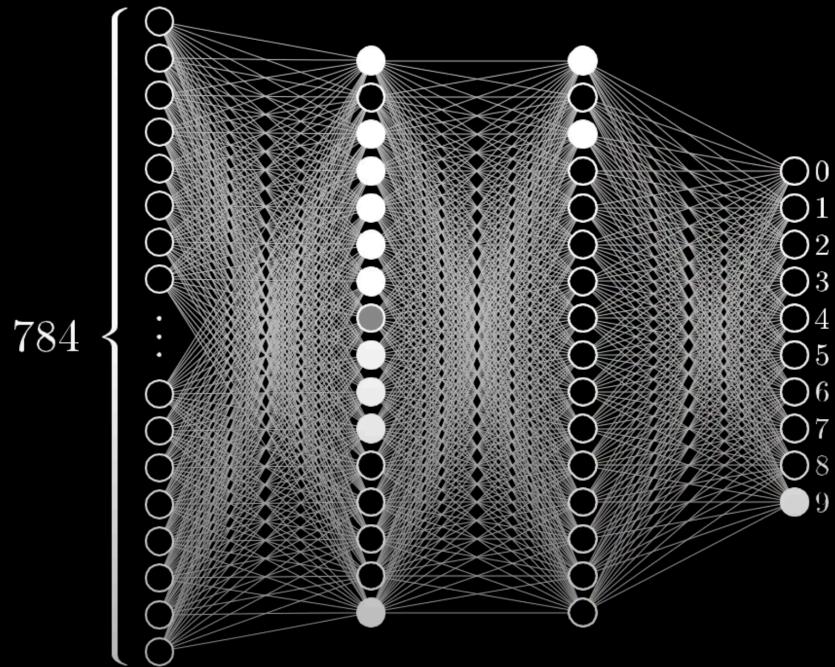
# Why the layers?



now do you recognize

a loop?

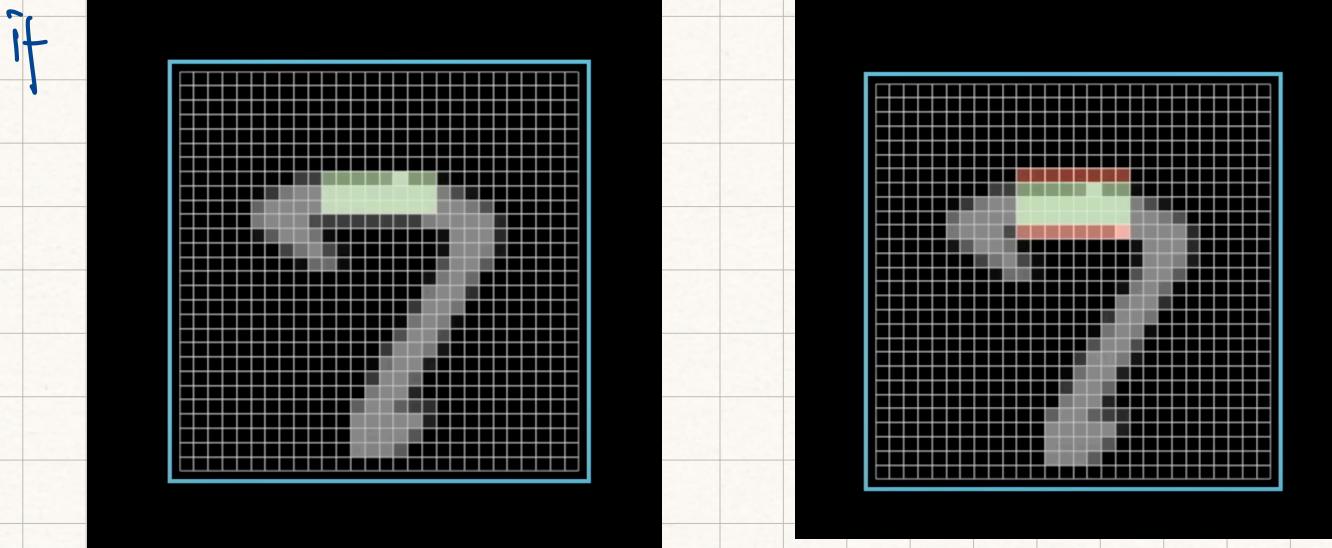
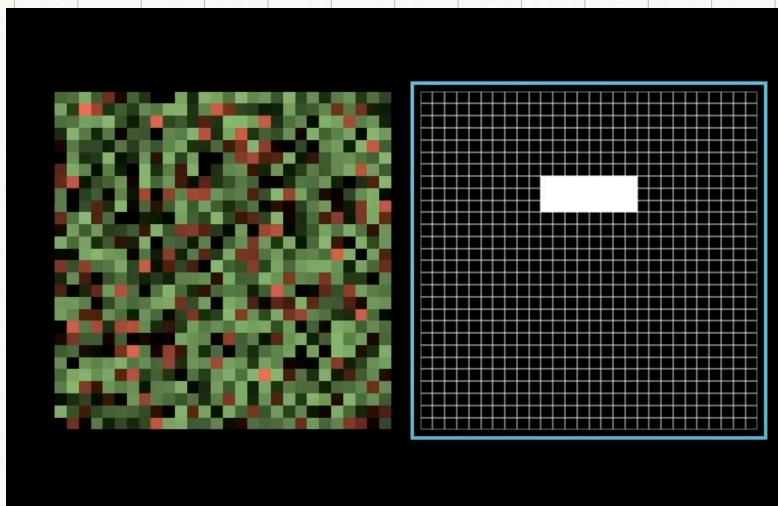
A line?



We want this neuron to light up when there's an edge here

so, we think of the params as "weights"

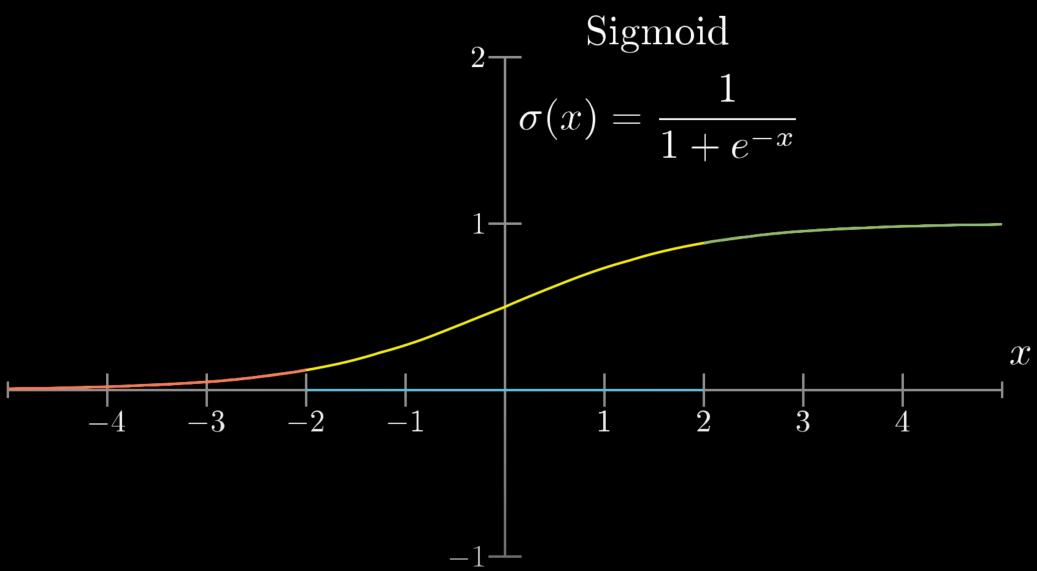
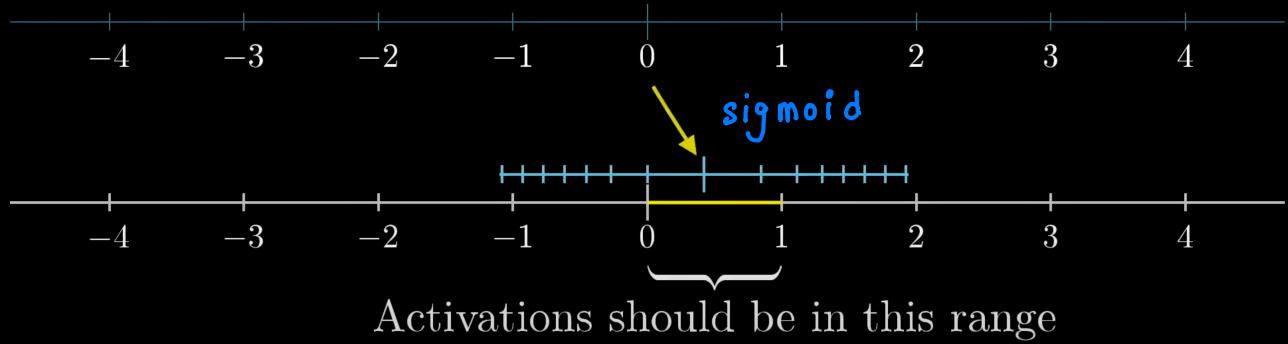
Useful to think of weights as a grid of their own:



only light up when  
that patch has white

sum is ↑↑ when  
middle pixel - white  
& surrounding → dark.

$$w_1a_1 + w_2a_2 + w_3a_3 + w_4a_4 + \cdots + w_na_n$$



$$\sigma(w_1a_1 + w_2a_2 + \cdots + w_na_n)$$

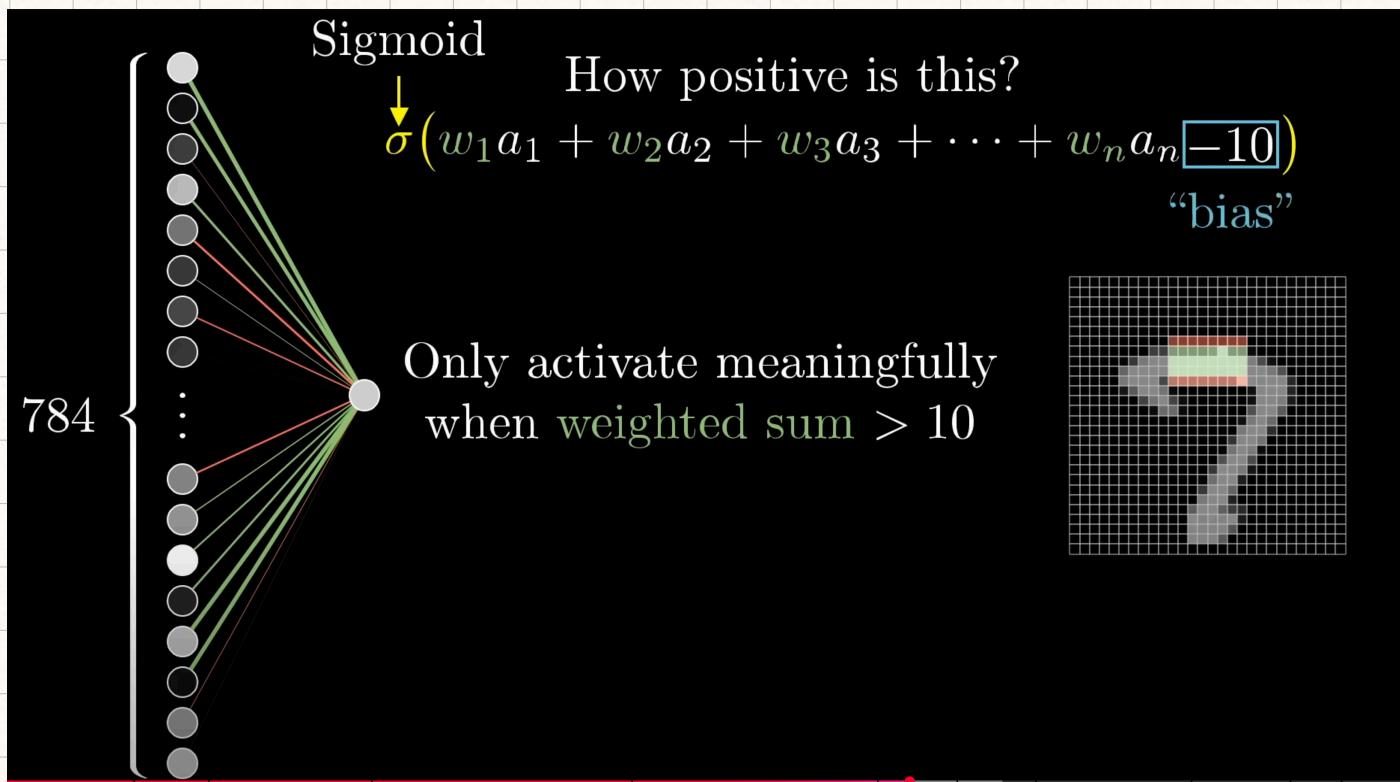
Say you only want neuron to be active when  
The weighted sum is more than 10

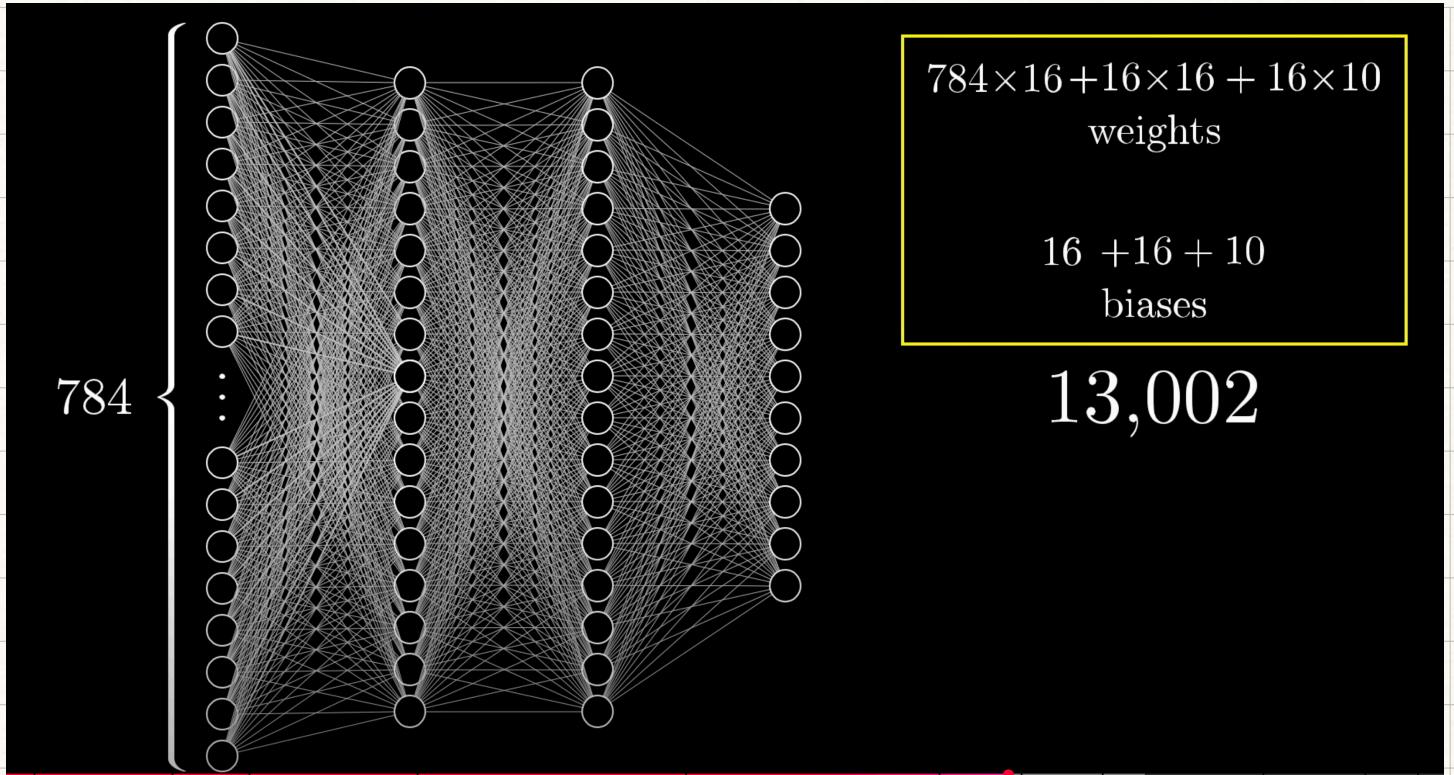
⇒ bias for inactive.

so,  $\sigma(w_1a_1 + w_2a_2 + \cdots + w_na_n - 10)$

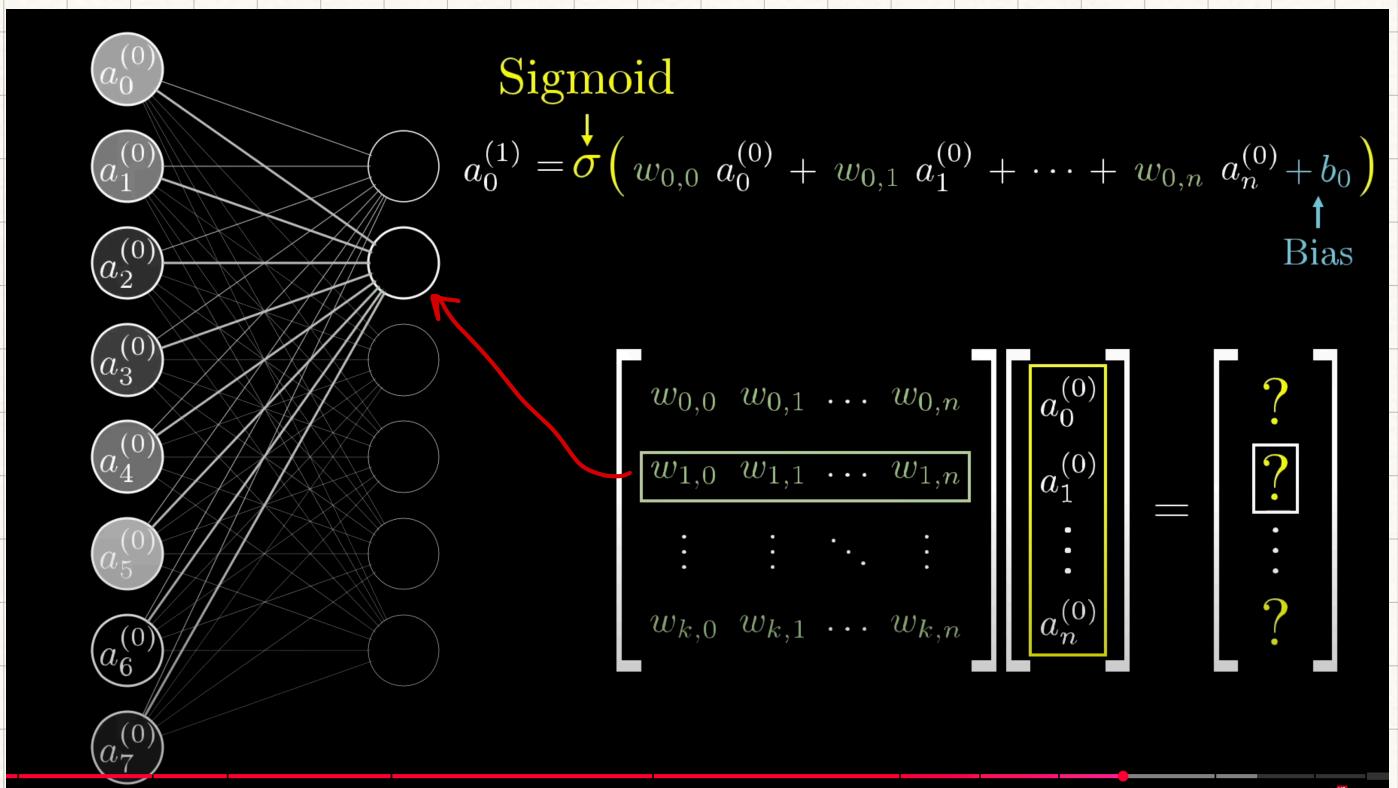
weights  $\rightarrow$  what pixel pattern the neuron is picking up on

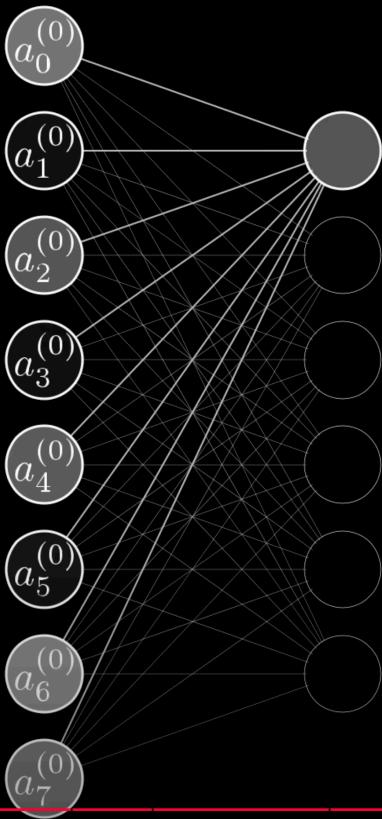
bias  $\rightarrow$  how high, weighted sum needs to be b4 the neuron starts to get meaningfully active





"training"  $\rightarrow$  finding  
 numbers  
 problem      valid setting for these  
 so that it can solve the  
 @ hand.



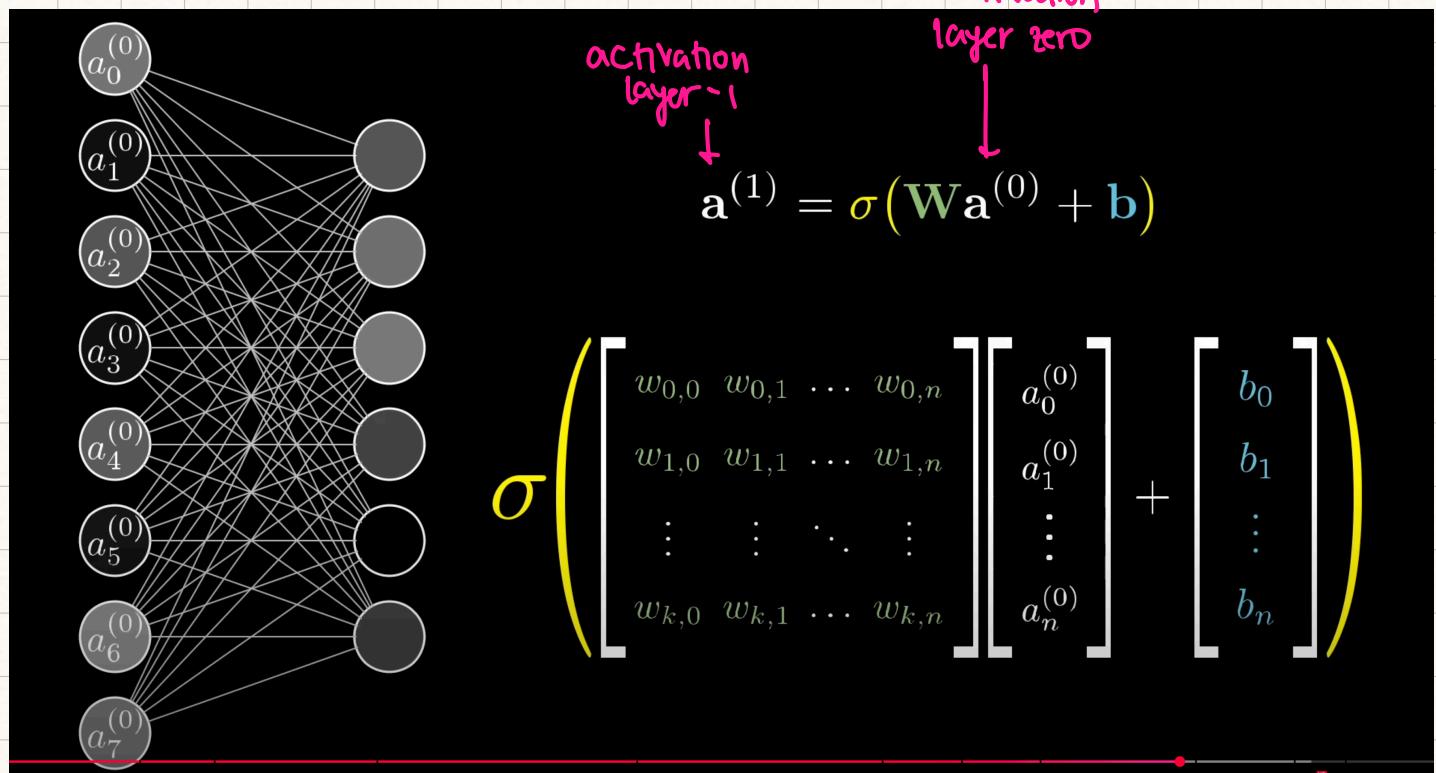


Sigmoid

$$a_0^{(1)} = \sigma(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0)$$

Bias

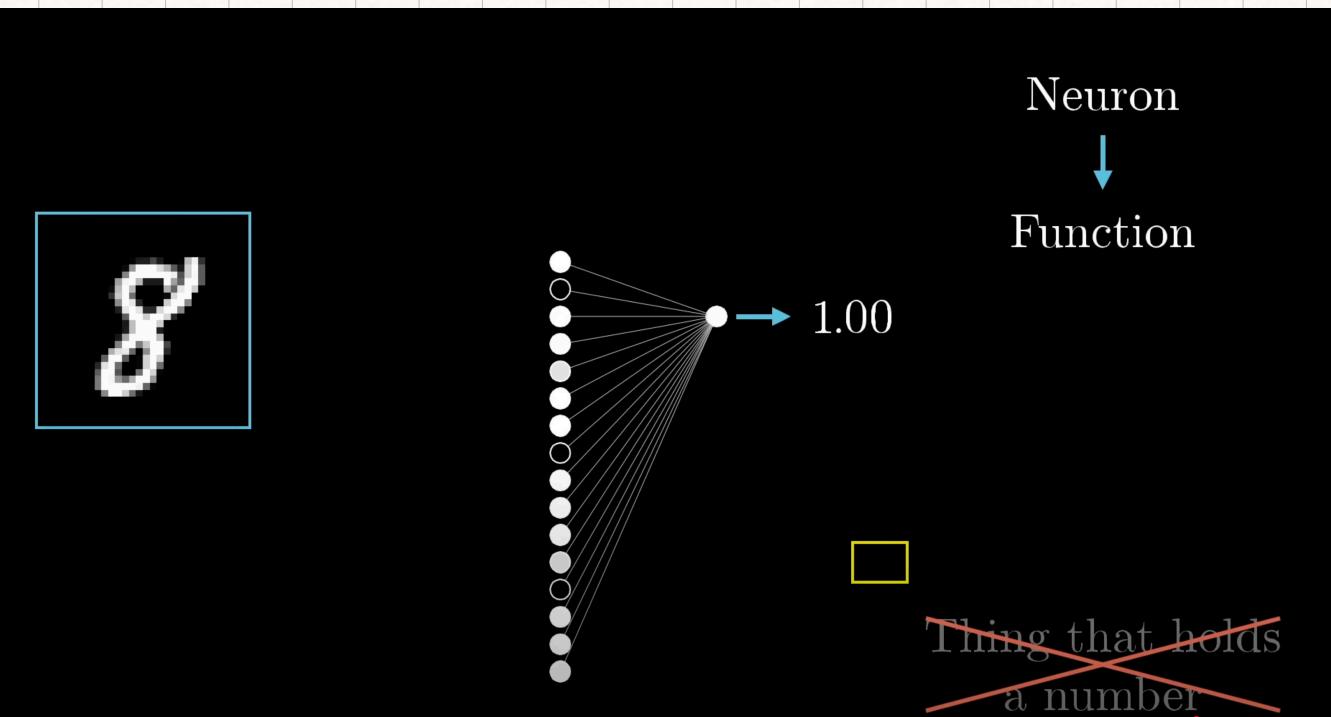
$$\sigma \left( \begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1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



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

```
class Network(object):
    def __init__(self, *args, **kwargs):
        #...yada yada, initialize weights and biases...

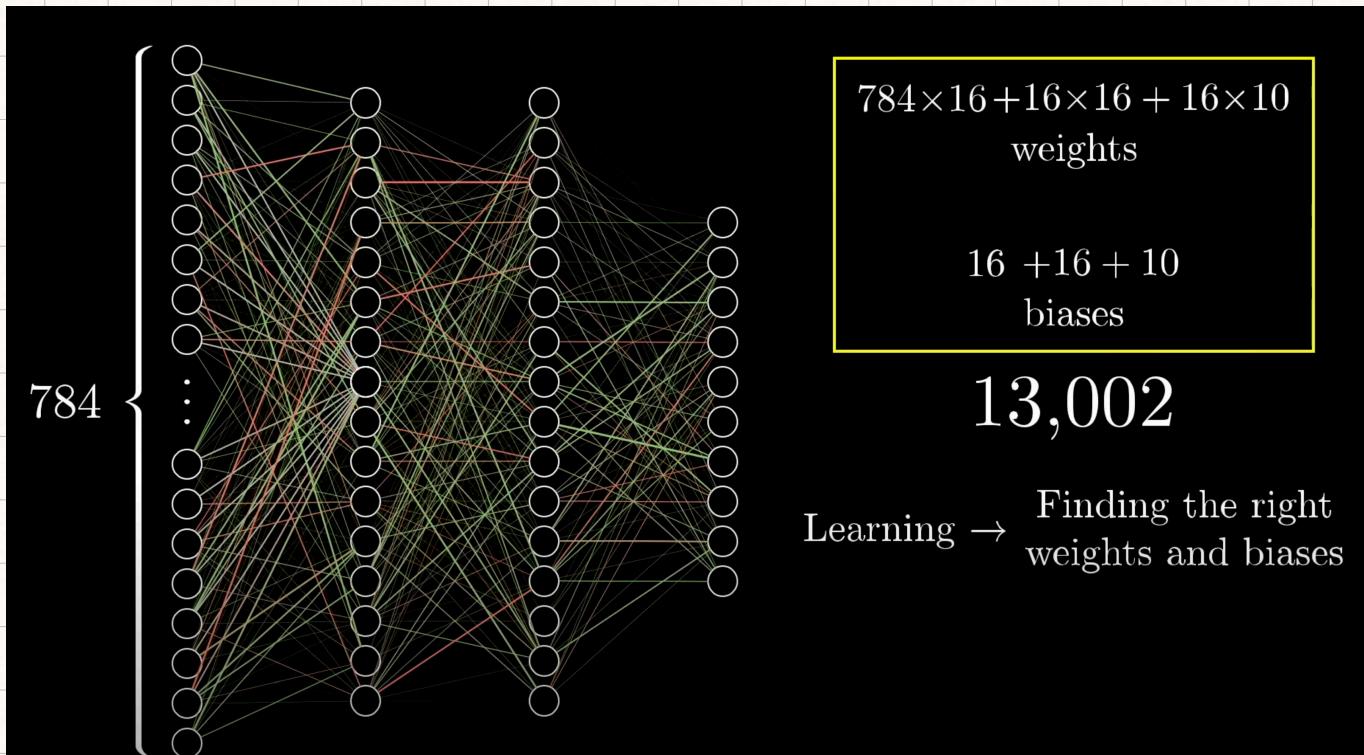
    def feedforward(self, a):
        """Return the output of the network for an input vector a"""
        for b, w in zip(self.biases, self.weights):
            a = sigmoid(np.dot(w, a) + b)
        return a
```



ReLU      ↗      sigmoid ↗

# 2<sup>nd</sup> Video || DL 2 || Gradient Descent.

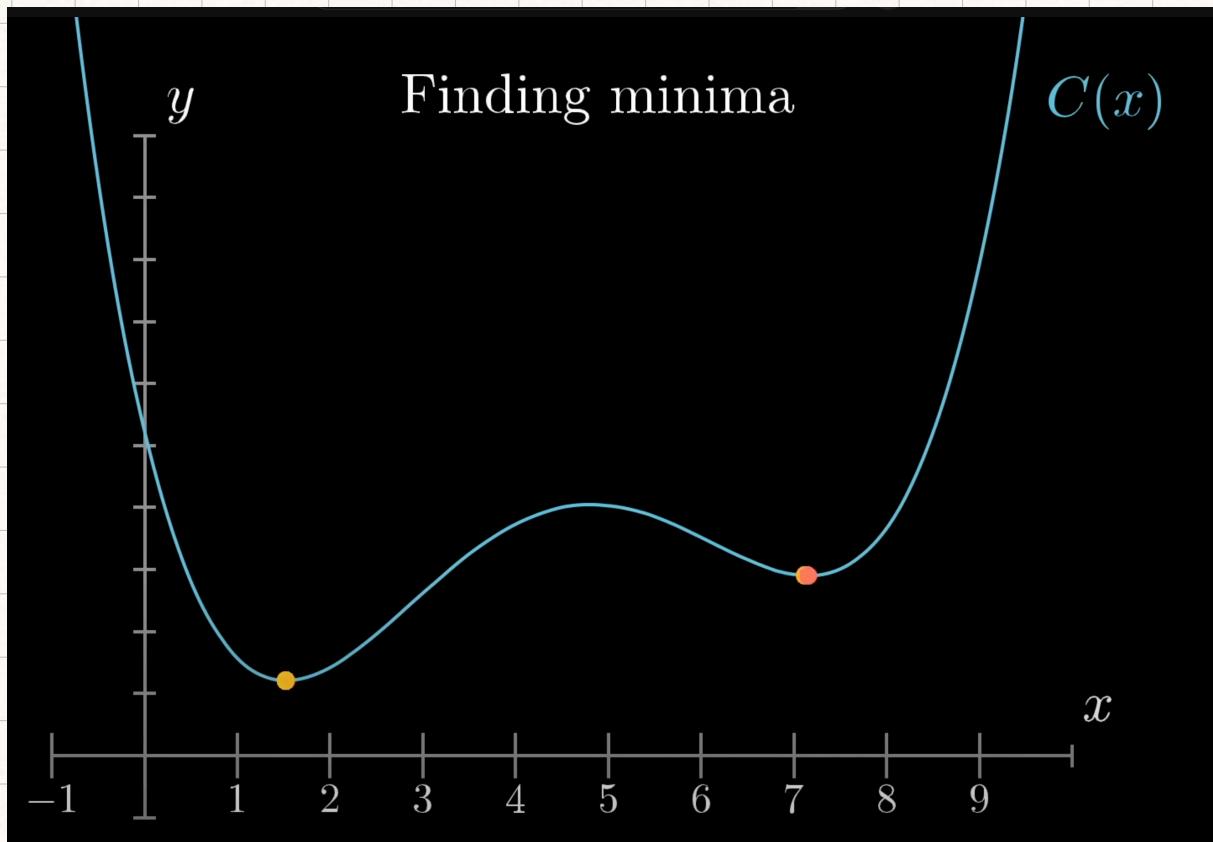
- <Recap>



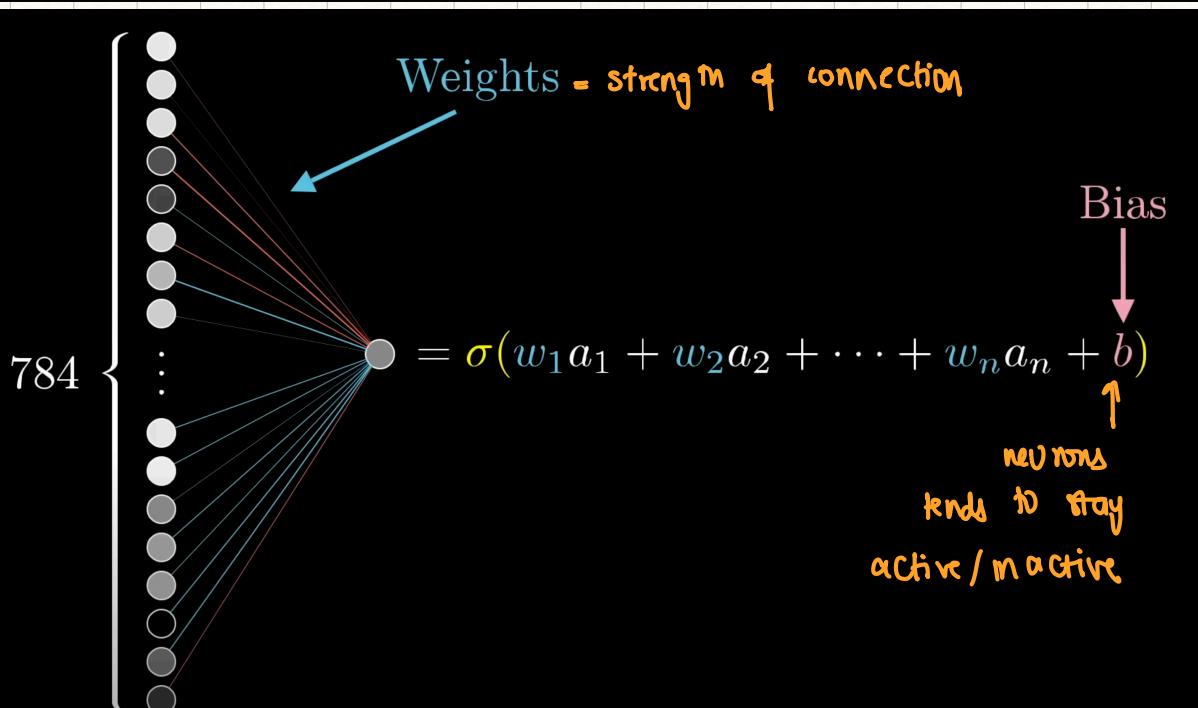
</Recap>

want an algo that can be shown a lot of  
training data  $\rightarrow$  images + labels  
will adjust its w & b to improve its  
performance on training data

Hopefully it'll generalise to images beyond training data.

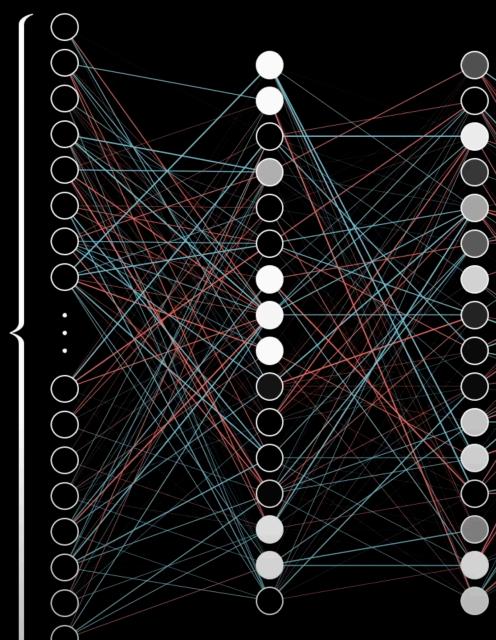


basically a calculus exercise, finding min. of a certain function.





784



What's the “cost”  
of this difference?



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

Utter trash

define “cost”

Cost of

3.37

$$\begin{cases}
 0.1863 \leftarrow (0.43 - 0.00)^2 + \\
 0.0809 \leftarrow (0.28 - 0.00)^2 + \\
 0.0357 \leftarrow (0.19 - 0.00)^2 + \\
 0.0138 \leftarrow (0.88 - 1.00)^2 + \\
 0.5242 \leftarrow (0.72 - 0.00)^2 + \\
 0.0001 \leftarrow (0.01 - 0.00)^2 + \\
 0.4079 \leftarrow (0.64 - 0.00)^2 + \\
 0.7388 \leftarrow (0.86 - 0.00)^2 + \\
 0.9817 \leftarrow (0.99 - 0.00)^2 + \\
 0.3998 \leftarrow (0.63 - 0.00)^2
 \end{cases}$$

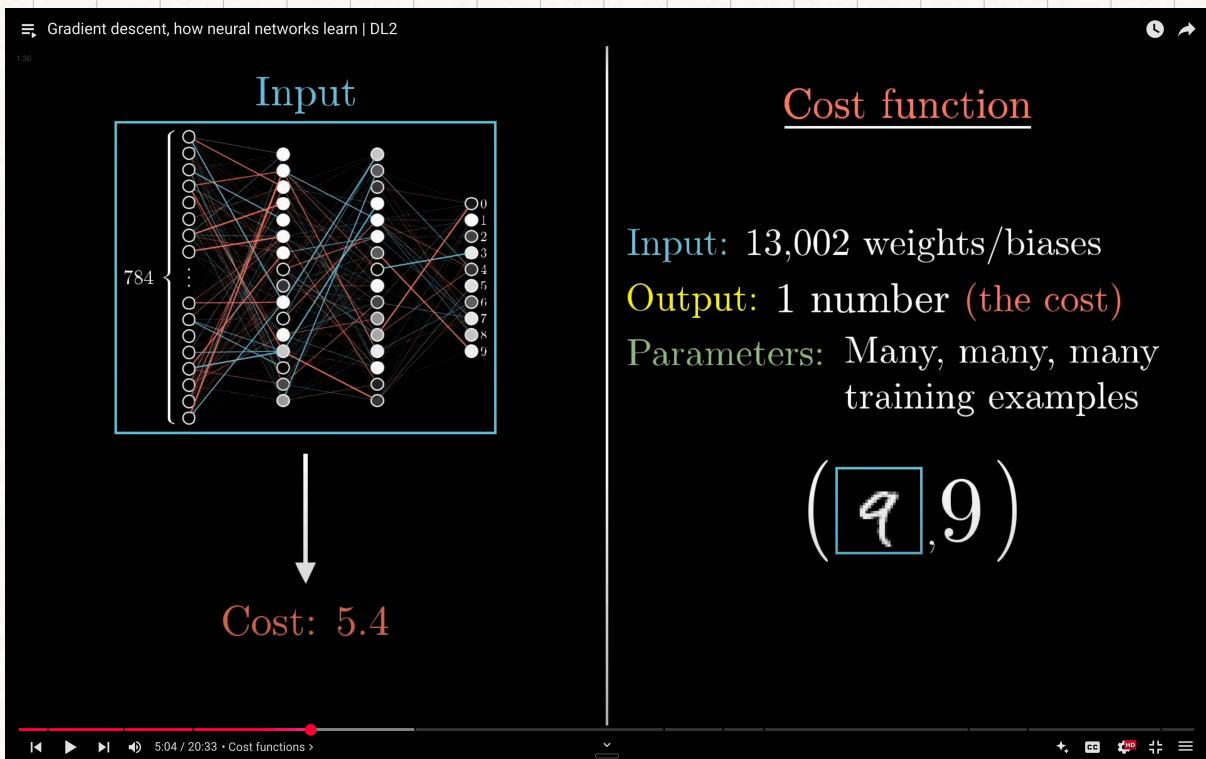
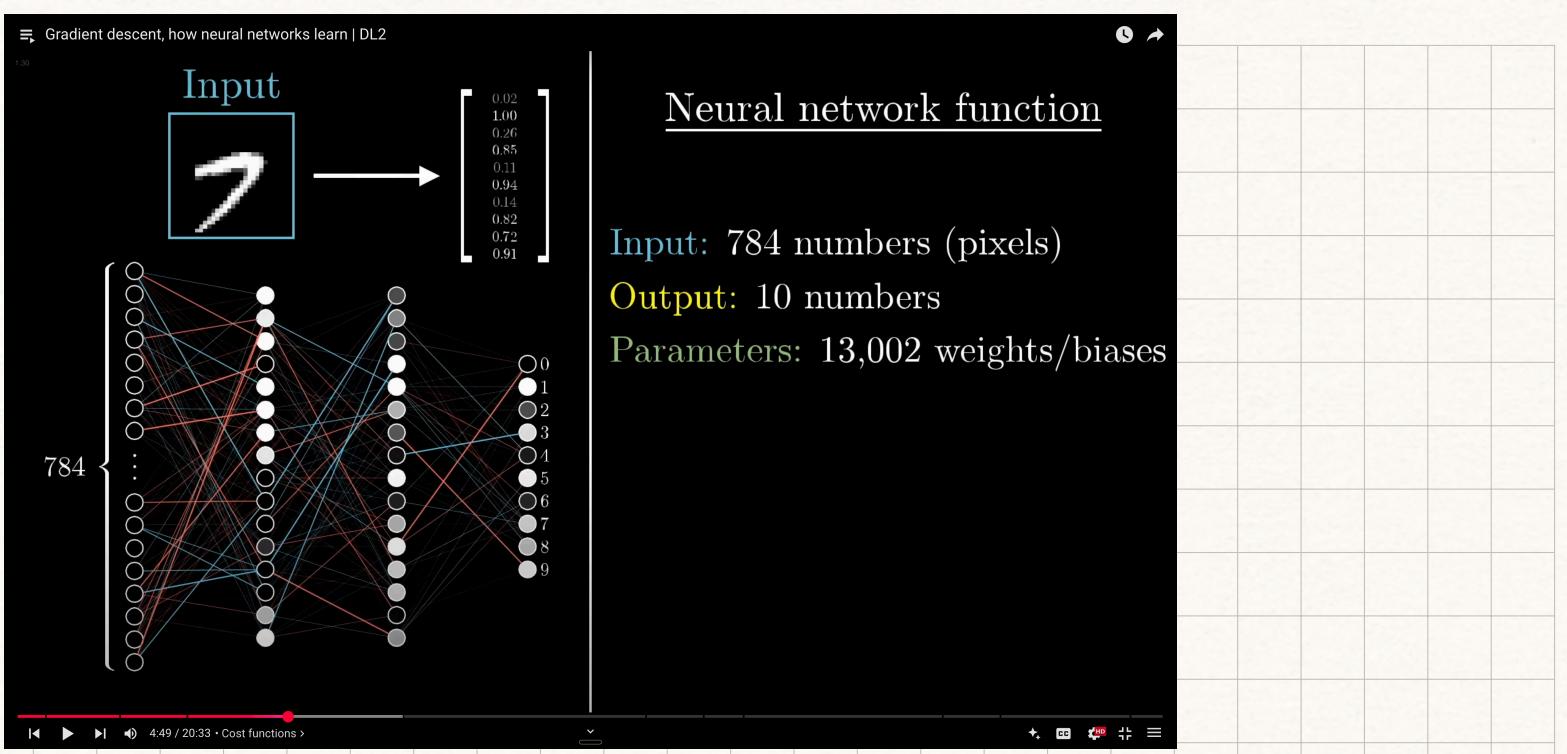
What's the “cost”  
of this difference?



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

Utter trash

now, take avg over all training samples.



Cost function

$$C(w_1, w_2, \dots, w_{13,002})$$

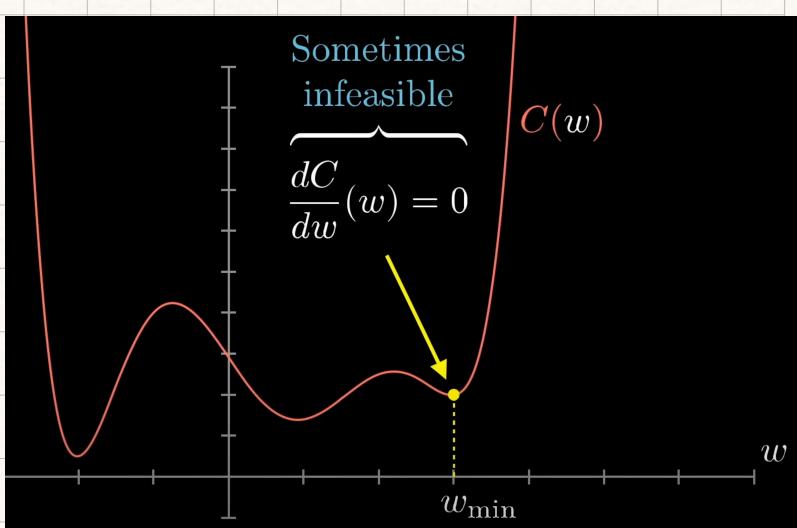
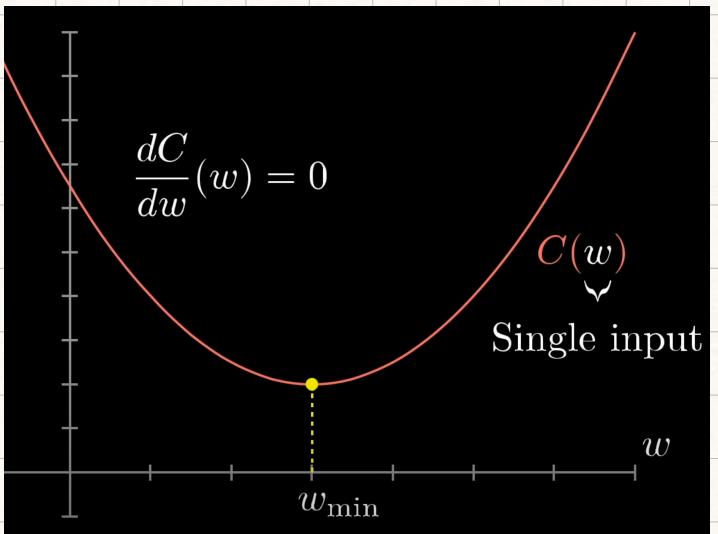
Weights and biases

Cost function

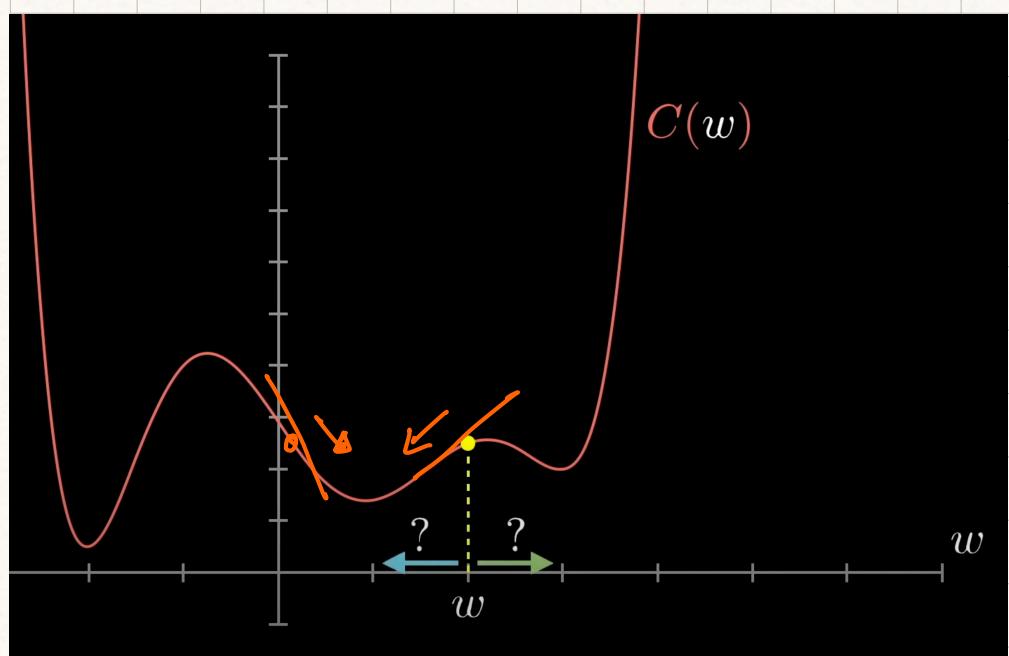
$$C(w)$$

Single input

easier to imagine



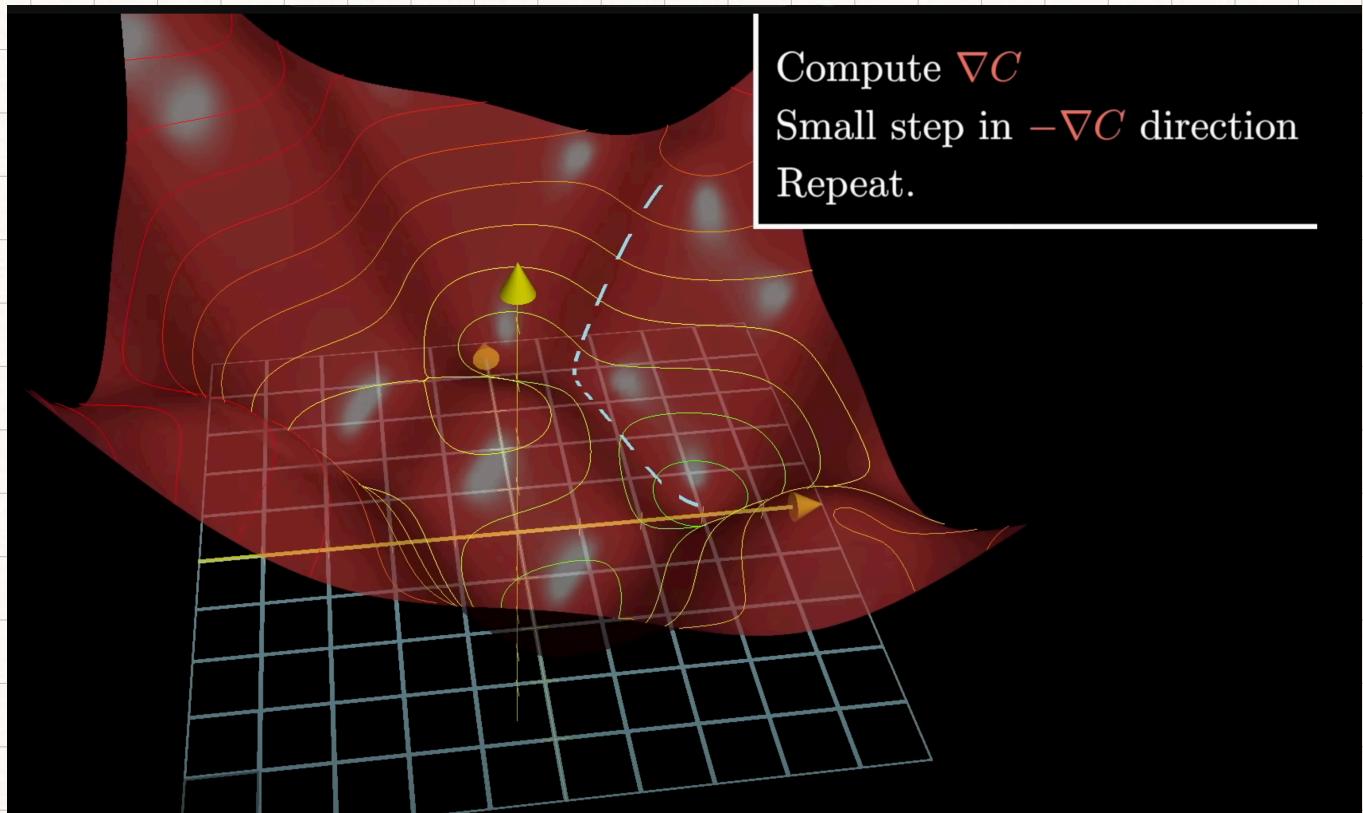
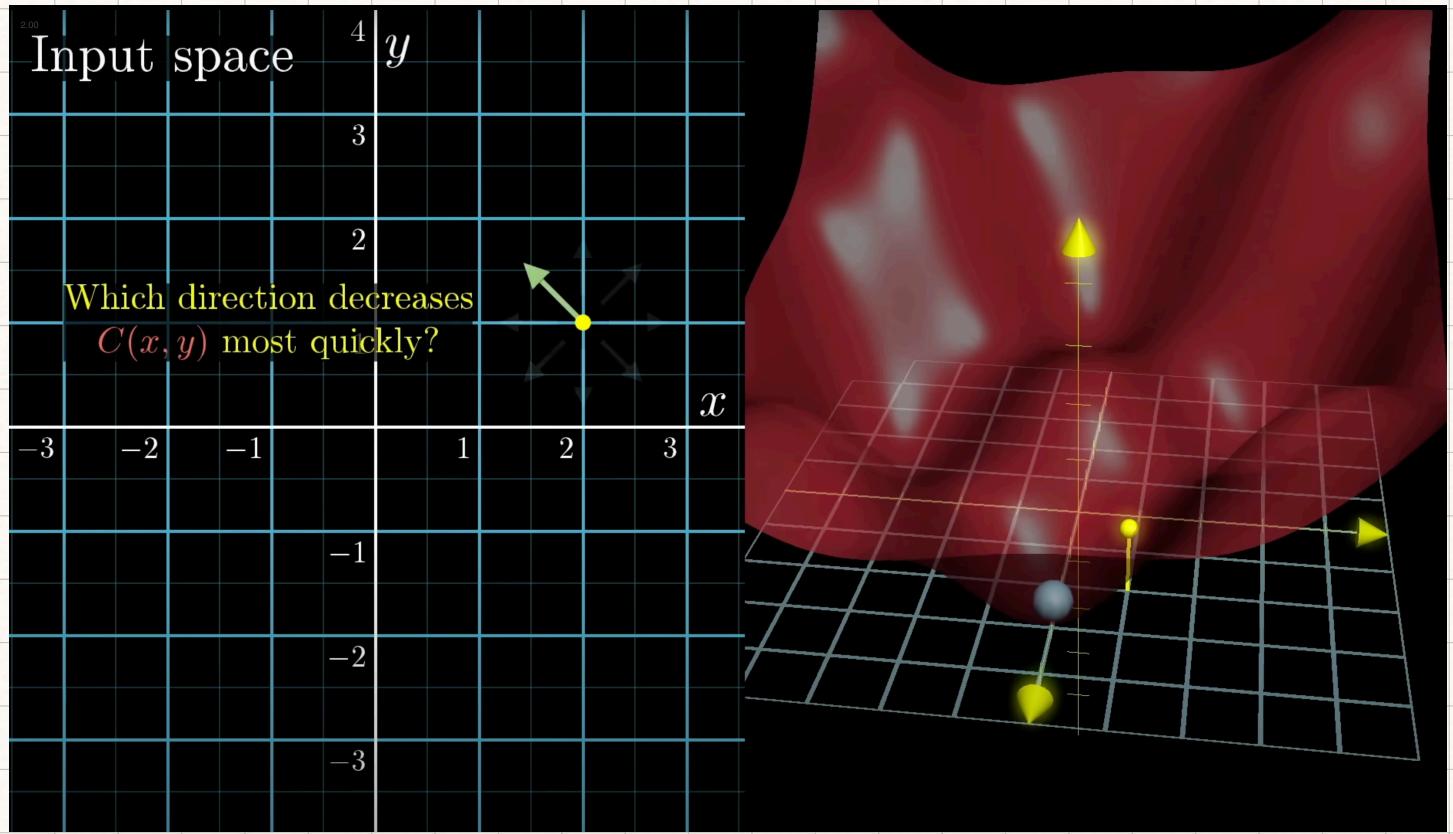
better way :



- "ball rolling down hill".
- many valleys That one might land in depending on which random input you start @ .

if step size  $\propto$  slope, Then steps get smaller

L smaller,



13,002 weights and biases

How to nudge all weights and biases

$$\vec{W} = \begin{bmatrix} 2.25 \\ -1.57 \\ 2.43 \\ -1.12 \\ 1.47 \\ \vdots \\ -0.76 \\ 3.50 \\ 1.21 + 0.82 \Rightarrow 2.03 \end{bmatrix}$$

$$-\nabla C(\vec{W}) =$$

$$\begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$

some direction which says which nudges to weights is going to cause most rapid decrease to cost func.  
z

Cost fn  $\rightarrow$  avg over all training data.

so minimising it  $\rightarrow$  better performance on all of those samples.

Back prop is how the gradient is calculated.

- Back prop next video.

What happens to each wt & each bias

↳ "intuitive feel" of what's happening.

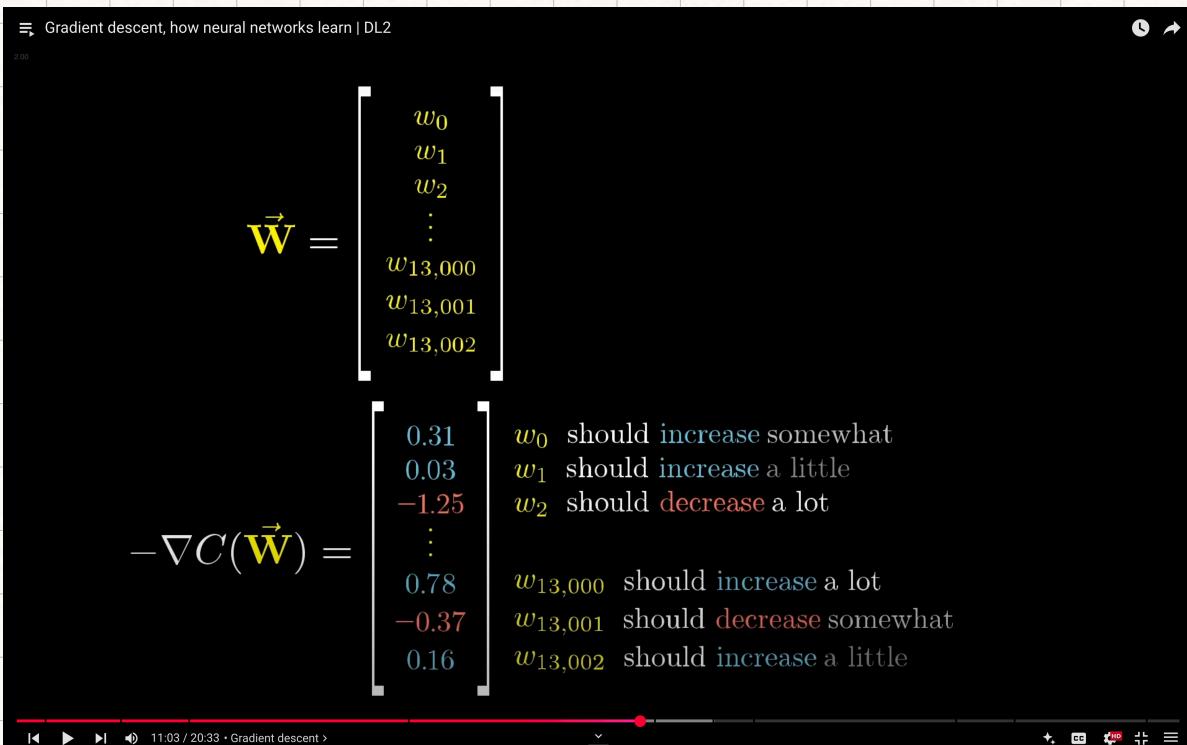
- Imp for cost function to be continuous & smooth

- process of repeatedly nudging the input of a M

by some multiple of the -ve gradient

↳ gradient descent.

- way to converge

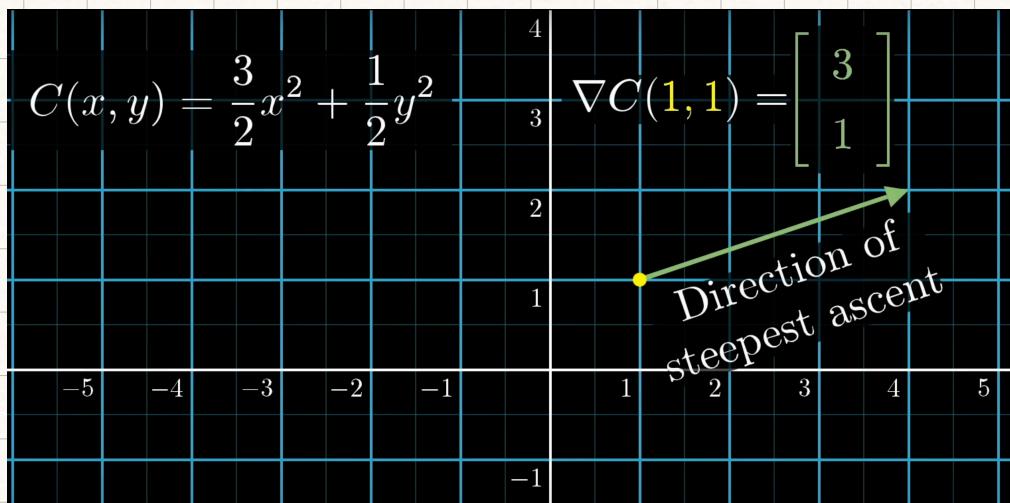
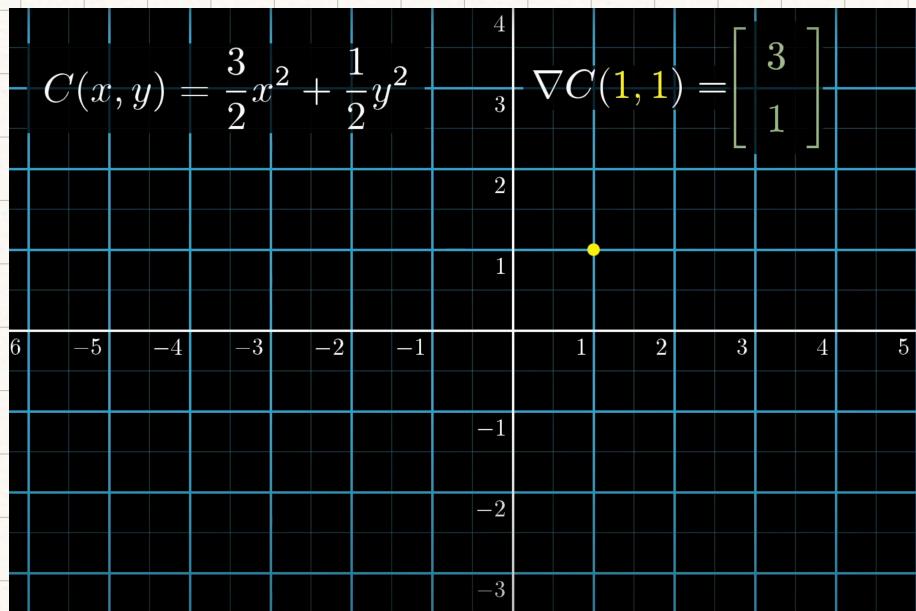


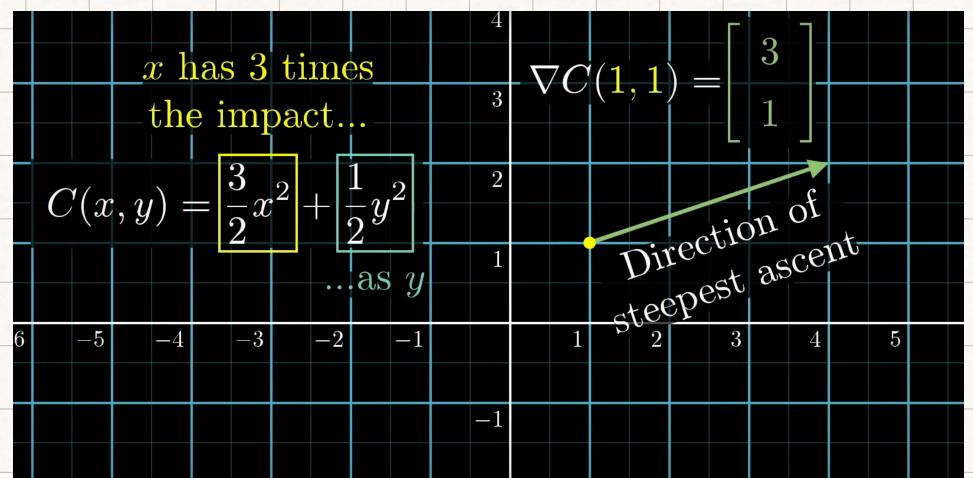
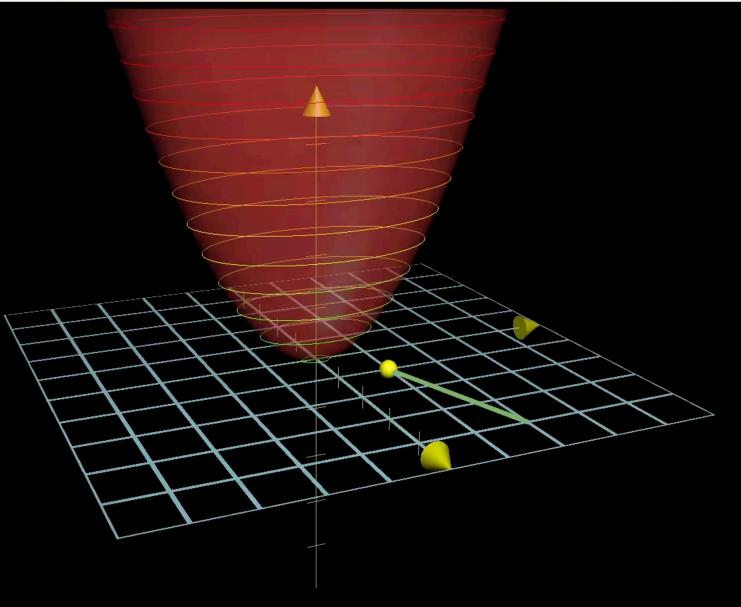
encodes relative imp. of each wt & bias.

→ which change  $\Rightarrow$  most bang for buck.

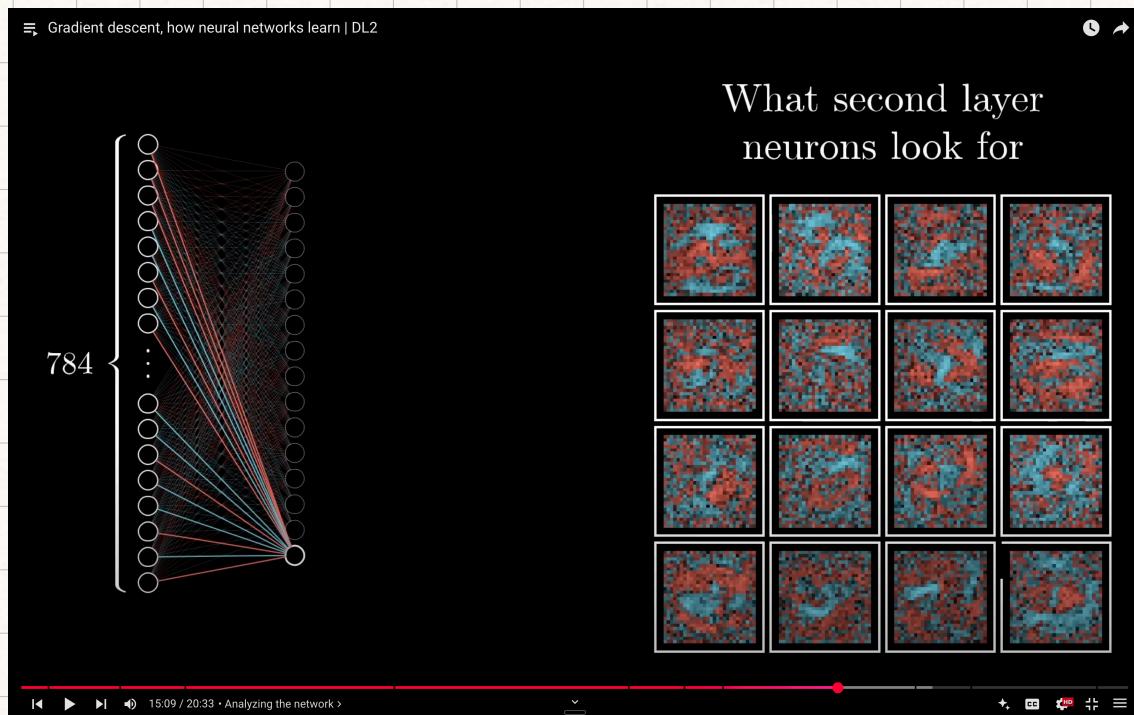
---

another way to think about direction.



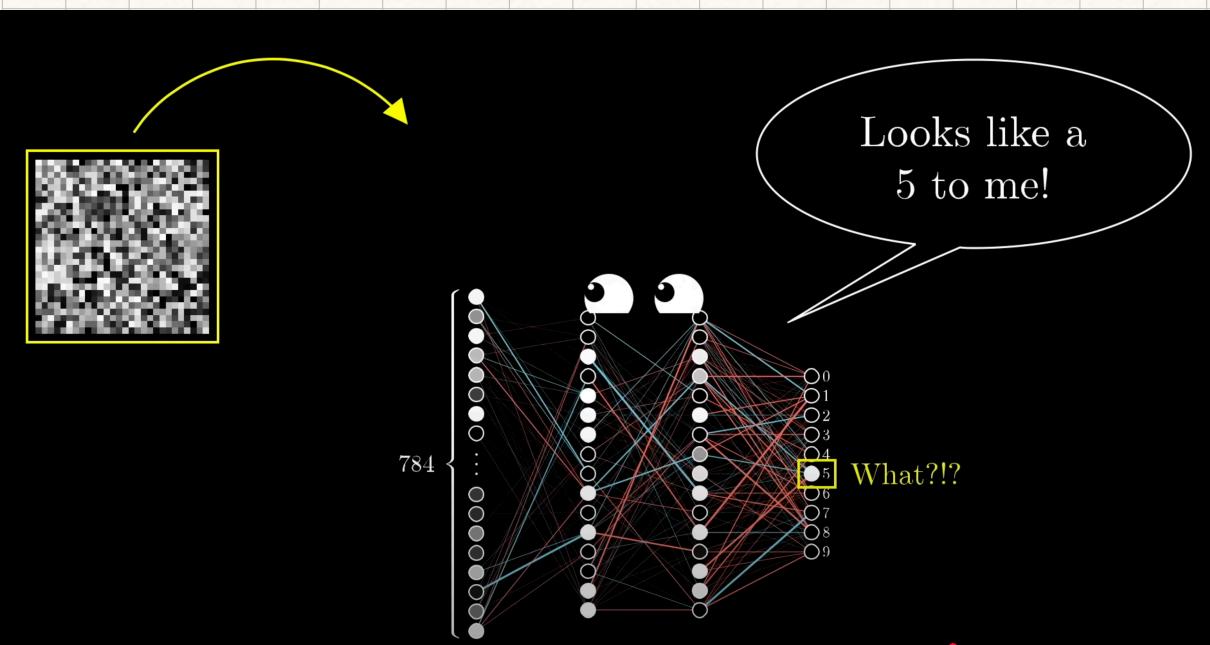


SOTA  $\rightarrow$  99.79 %



DOES NOT LOOK FOR THE "edges" WHICH  
WE HAD THOUGHT OF INITIALLY...  
...

- ALSO WHEN GIVEN RANDOM NOISE IT CONFIDENTLY  
PREDICTS "5"



- IT HAS NO IDEA HOW TO DRAW NUMBERS

## UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang\*  
Massachusetts Institute of Technology  
chiyuan@mit.edu

Samy Bengio  
Google Brain  
bengio@google.com

Moritz Hardt  
Google Brain  
mrtz@google.com

Benjamin Recht†  
University of California, Berkeley  
brecht@berkeley.edu

Oriol Vinyals  
Google DeepMind  
vinyals@google.com

→ shuffle labels of CNN  
• same training accuracy  
≡

## A Closer Look at Memorization in Deep Networks

→ little smarter

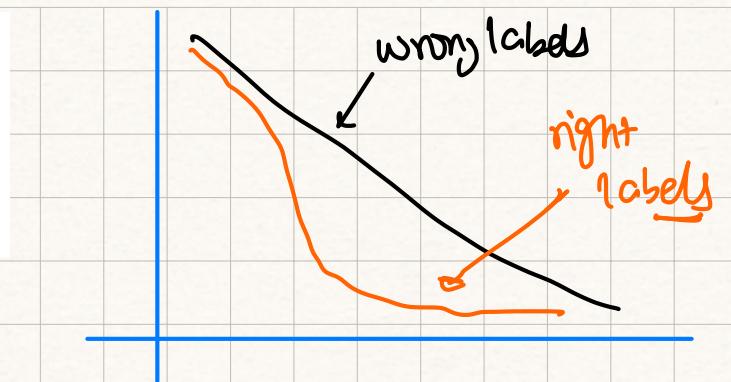
## The Loss Surfaces of Multilayer Networks

Anna Choromanska Mikael Henaff Michael Mathieu Gérard Ben Arous Yann LeCun  
achoroma@cims.nyu.edu mbh305@nyu.edu mathieu@cs.nyu.edu benarous@cims.nyu.edu yann@cs.nyu.edu

Courant Institute of Mathematical Sciences  
New York, NY, USA

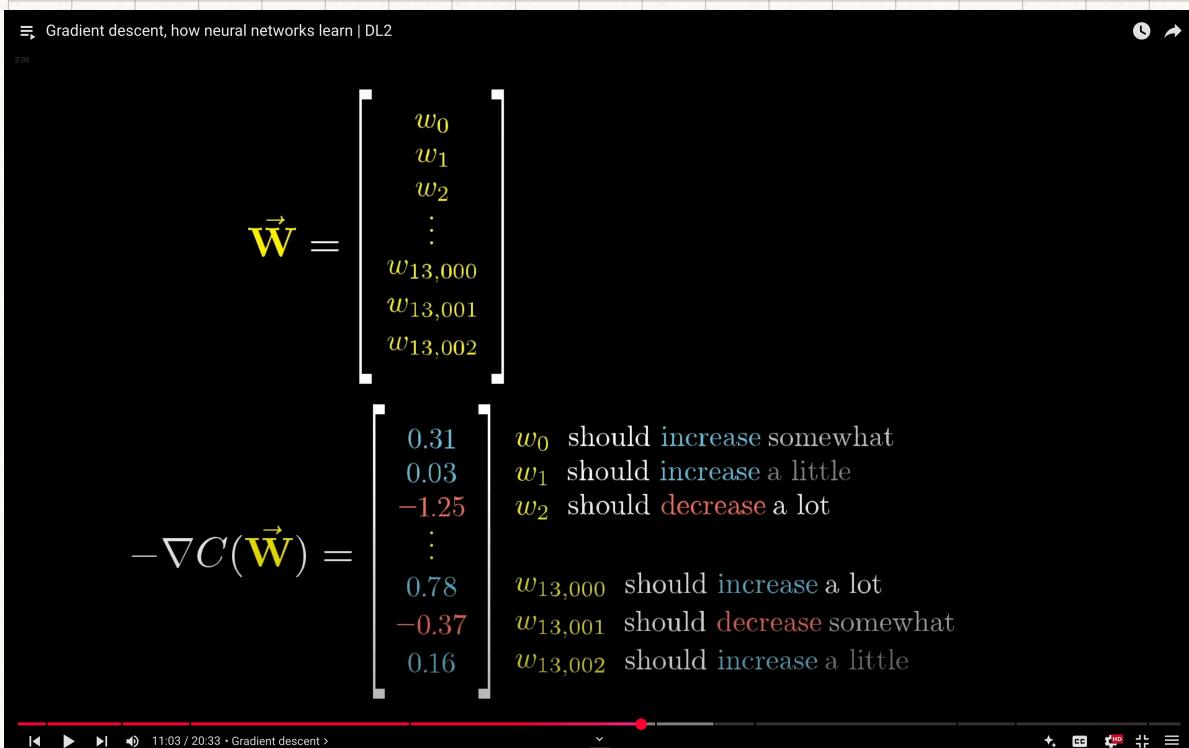


Roughly equal quality local minima.



# DL 03 - 3B1B Playlist // Back prop intuition.

- Back prop is the algo used to compute the gradient.



$$C(w) = \text{some num}$$

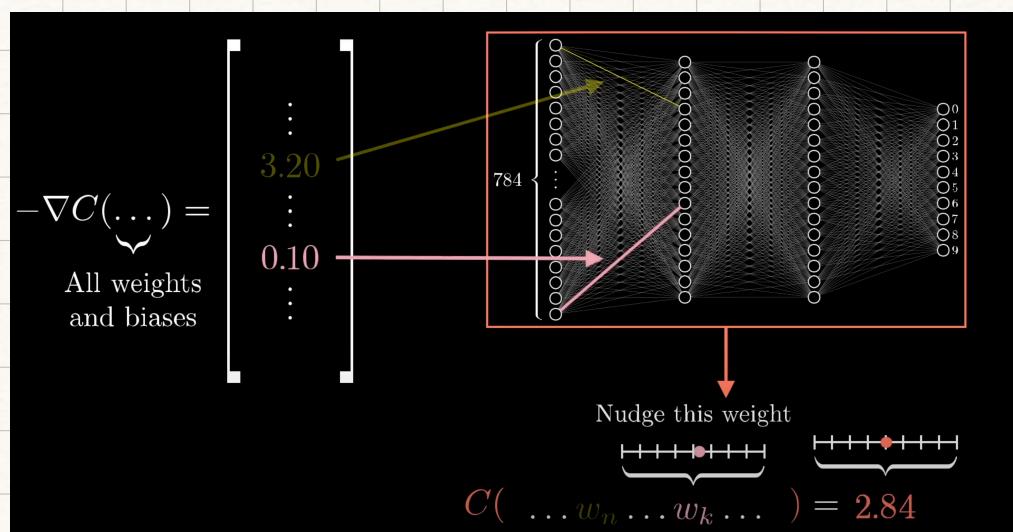
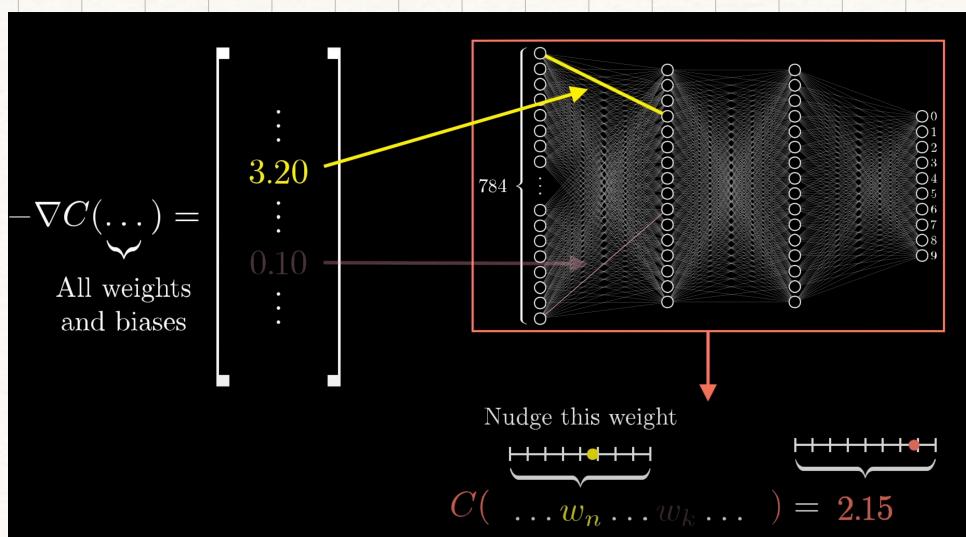
$-\nabla C(w)$  = tells how to change wt & bias to  
(gradient) most efficiently decrease the cost.

o direction in 13K dim  $\rightarrow$  xx {bad intuition}.

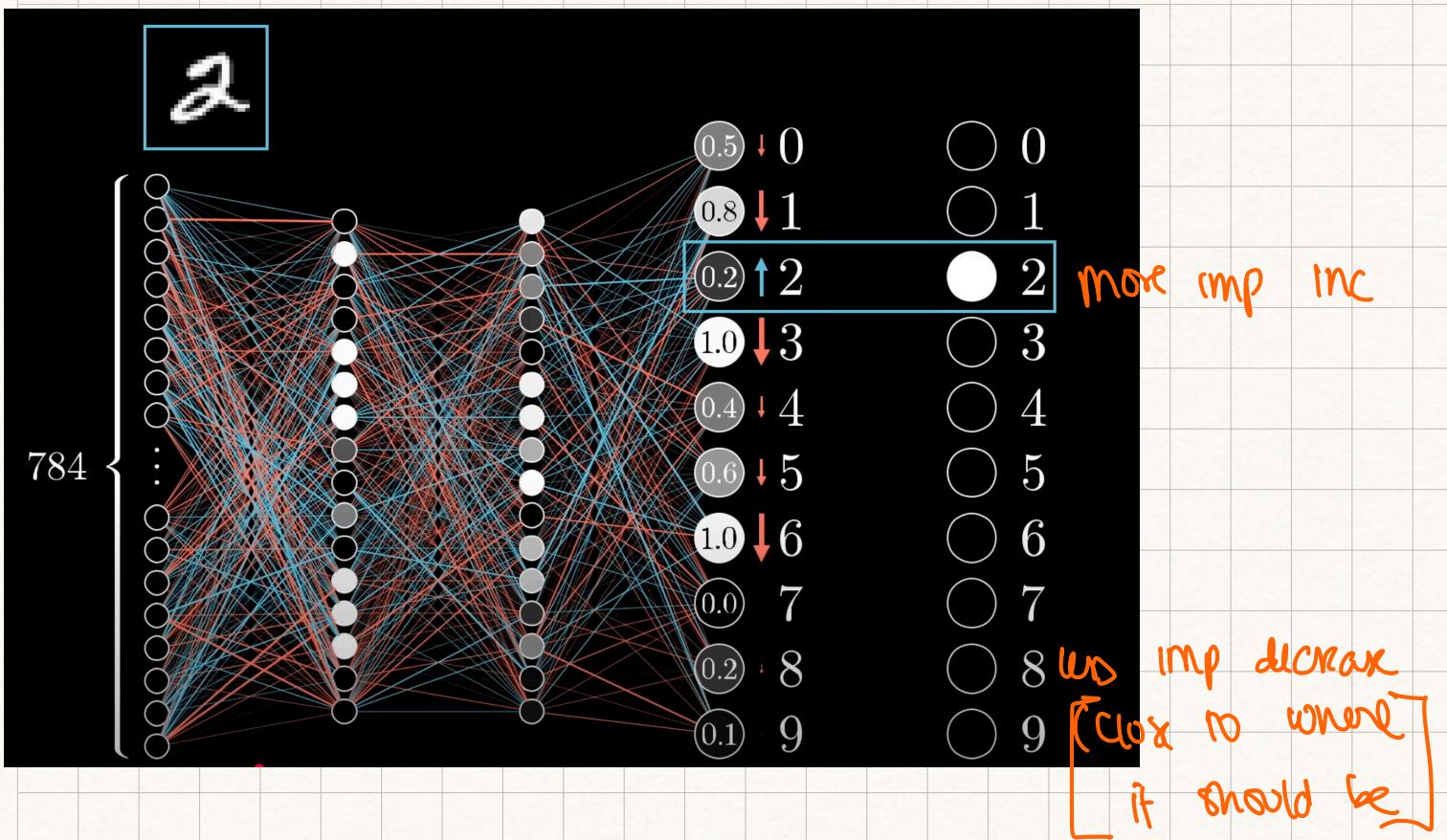
o magnitude tells us how sensitive the C is to  
each wt & bias

Eg :

Loss is 32x more sensitive to change in yellow  
as compared to a change in pink.



# Intuitive back prop for 1 example.

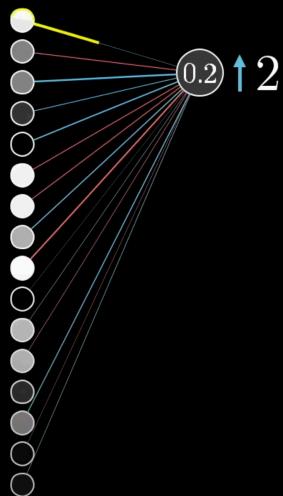


$$\textcircled{0.2} = \sigma(w_0a_0 + w_1a_1 + \dots + w_{n-1}a_{n-1} + b)$$

Increase  $b$

Increase  $w_i$

Change  $a_i$



- CONNECTION w/ brightest neurons from prev layer m. largest effect.

- Inc wt. of a unu w/  $\uparrow$  activation in prev layer, it has stronger infl. on C than inc the wt of connections w/ dummer neurons.

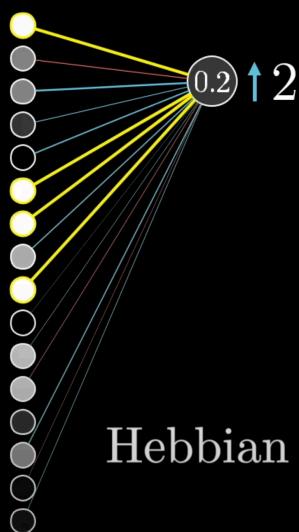
2

$$0.2 = \sigma(w_0a_0 + w_1a_1 + \dots + w_{n-1}a_{n-1} + b)$$

Increase  $b$

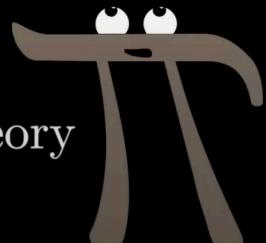
Increase  $w_i$   
in proportion to  $a_i$

Change  $a_i$



Hebbian theory

"Neurons that fire together wire together"



- strength connctn  $\uparrow$  b/w neurons that are most active & the ones we wish to become more active.

- Neurons firing while seeing a two, get max strongly linked to those "thinking about a two".

3rd way :

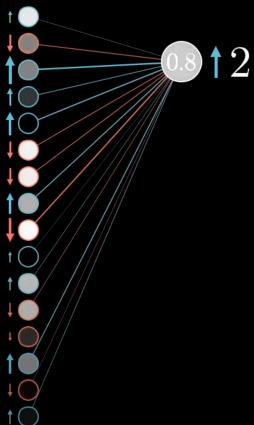
$$\textcircled{0.2} = \sigma(w_0 a_0 + w_1 a_1 + \dots + w_{n-1} a_{n-1} + b)$$

2

Increase  $b$

Increase  $w_i$   
in proportion to  $a_i$

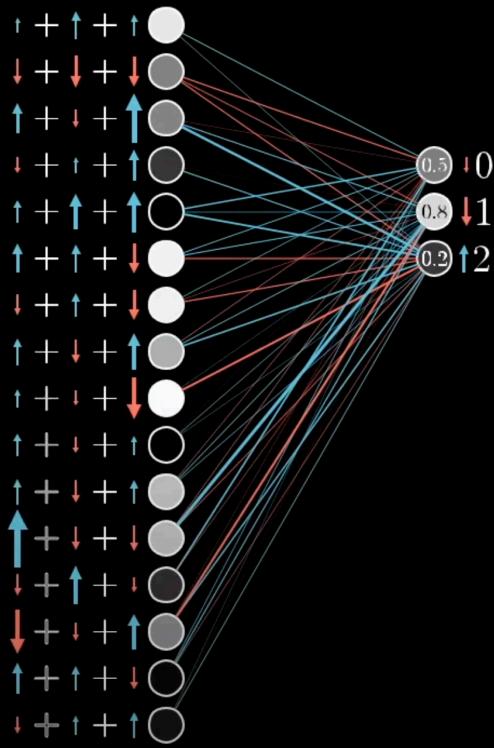
Change  $a_i$



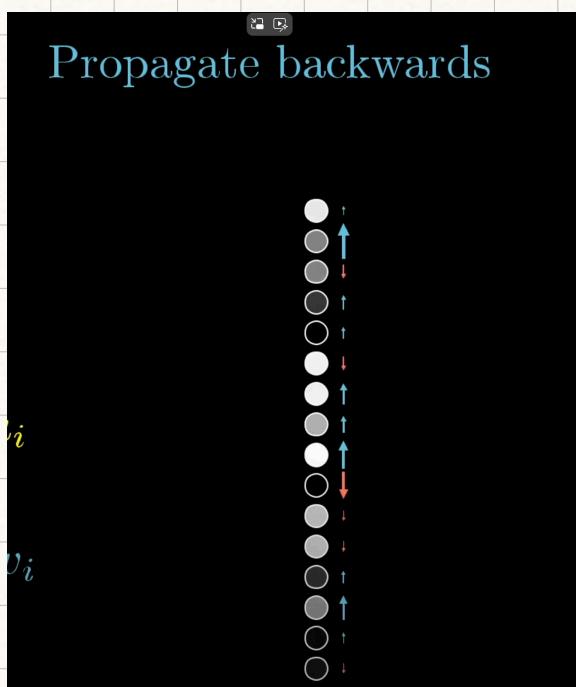
- changing activation
- everything connected w/ pos wt got brighter
- everything connected w/ neg wt got dimmer.
- most bang for buck by scaling changes  $\propto$  to ws.

This is only for the <sup>c</sup> <sup>-2'</sup> neuron in last layer

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



← Other neurons in the output layer w.r.t. other desires.  
 for what should happen to corresponding layer in proportion to their Wts.



← you get a list that you want to happen to 2nd to last layer.

Then recursively apply the same process to the prev layer.

							...	Average over all training data
$w_0$	-0.08	+0.02	-0.02	+0.11	-0.05	-0.14	...	-0.08
$w_1$	-0.11	+0.11	+0.07	+0.02	+0.09	+0.05	...	+0.12
$w_2$	-0.07	-0.04	-0.01	+0.02	+0.13	-0.15	...	-0.06
:	:	:	:	:	:	:	..	:
$w_{13,001}$	+0.13	+0.08	-0.06	-0.09	-0.02	+0.04	...	+0.04

You do for all training ex. &  
avg the critics of each training example.



mini-batches

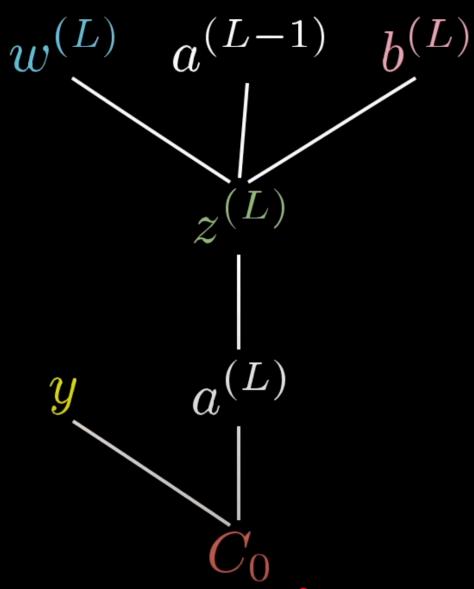
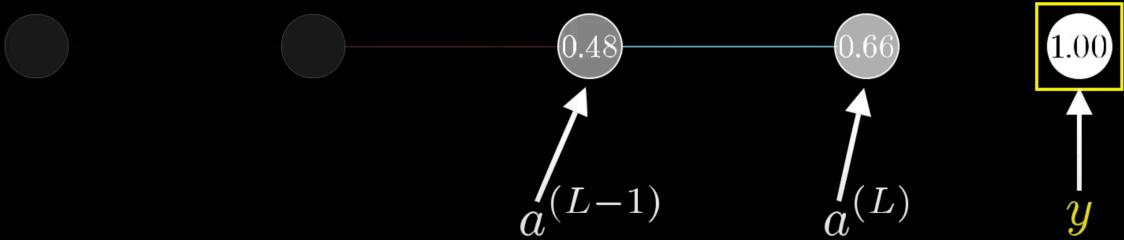
# DL 4 Backprop Calculus

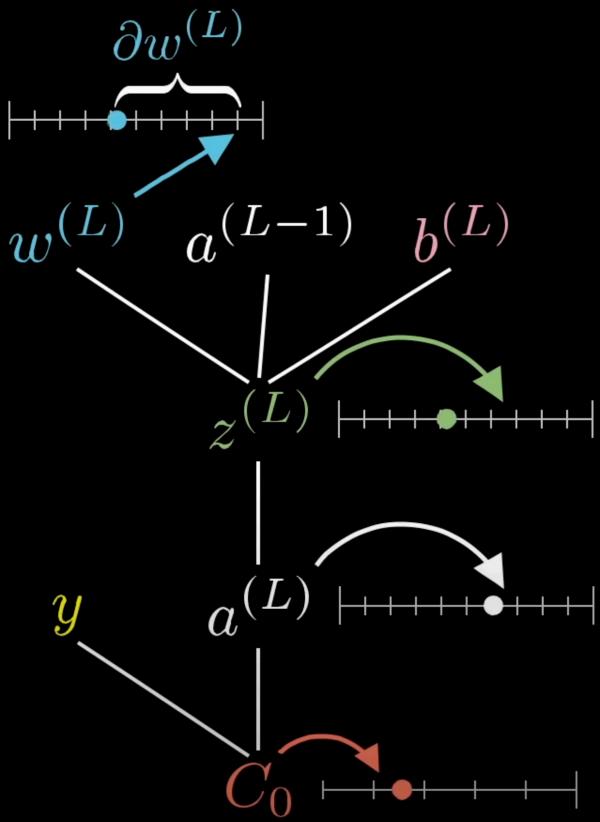
Cost  $\rightarrow C_0(\dots) = (a^{(L)} - y)^2$

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

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

Desired output





How sensitive is the cost fn. to small changes in w<sup>(L)</sup>?

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

whatever the resultant nudge to the cost is  
small nudge to w<sup>(L)</sup>

We break things up as  $\Delta w \rightarrow \Delta z \rightarrow \Delta a \rightarrow \Delta c$

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

← chain rule

↳ gives sensitivity of c to small changes in w<sup>(L)</sup>,

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

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

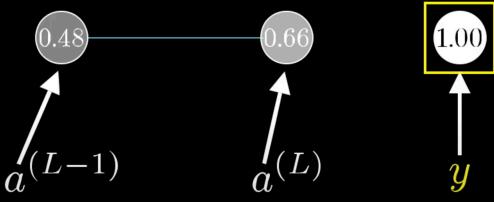
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

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

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

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

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



The amount that  $dW^{(L)}$  influenced the last layer depends on how strong the previous neuron is.

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

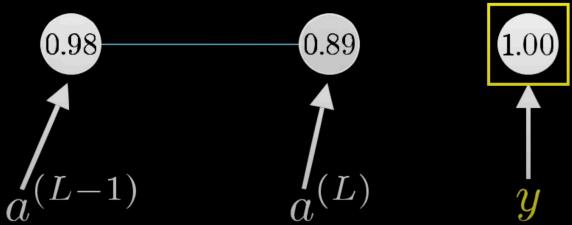
$$C_0 = (a^{(L)} - y)^2$$

Average of all training examples

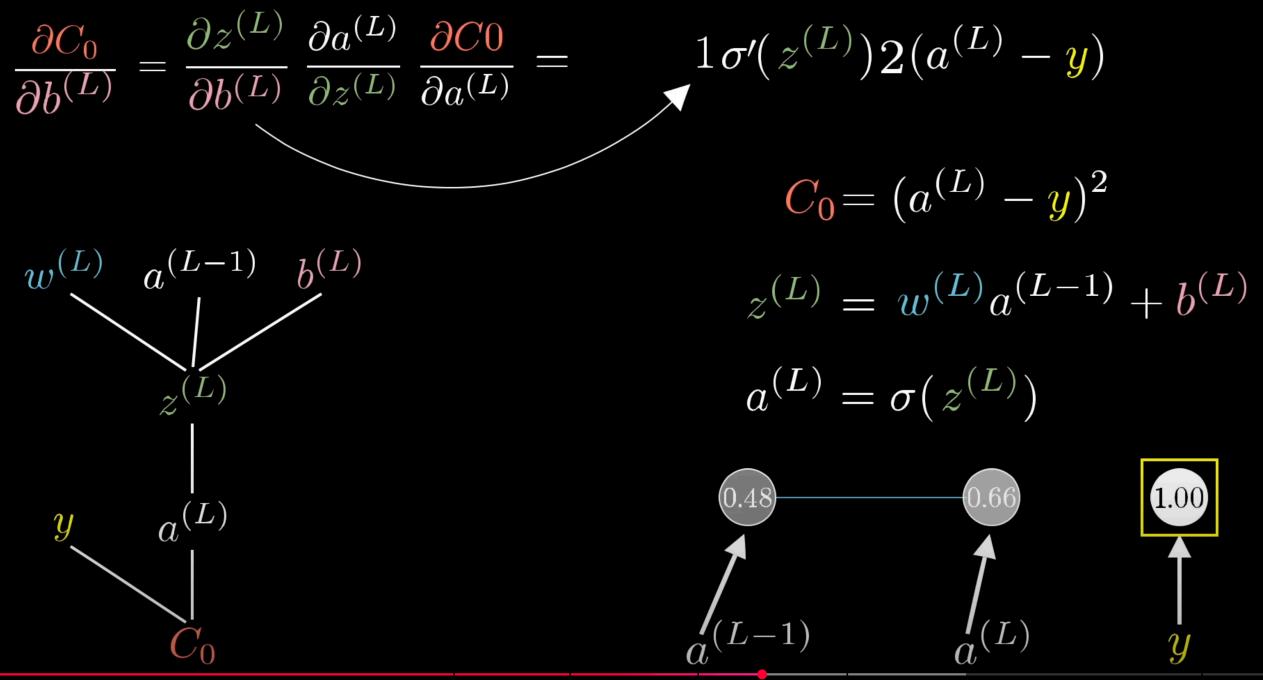
$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

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

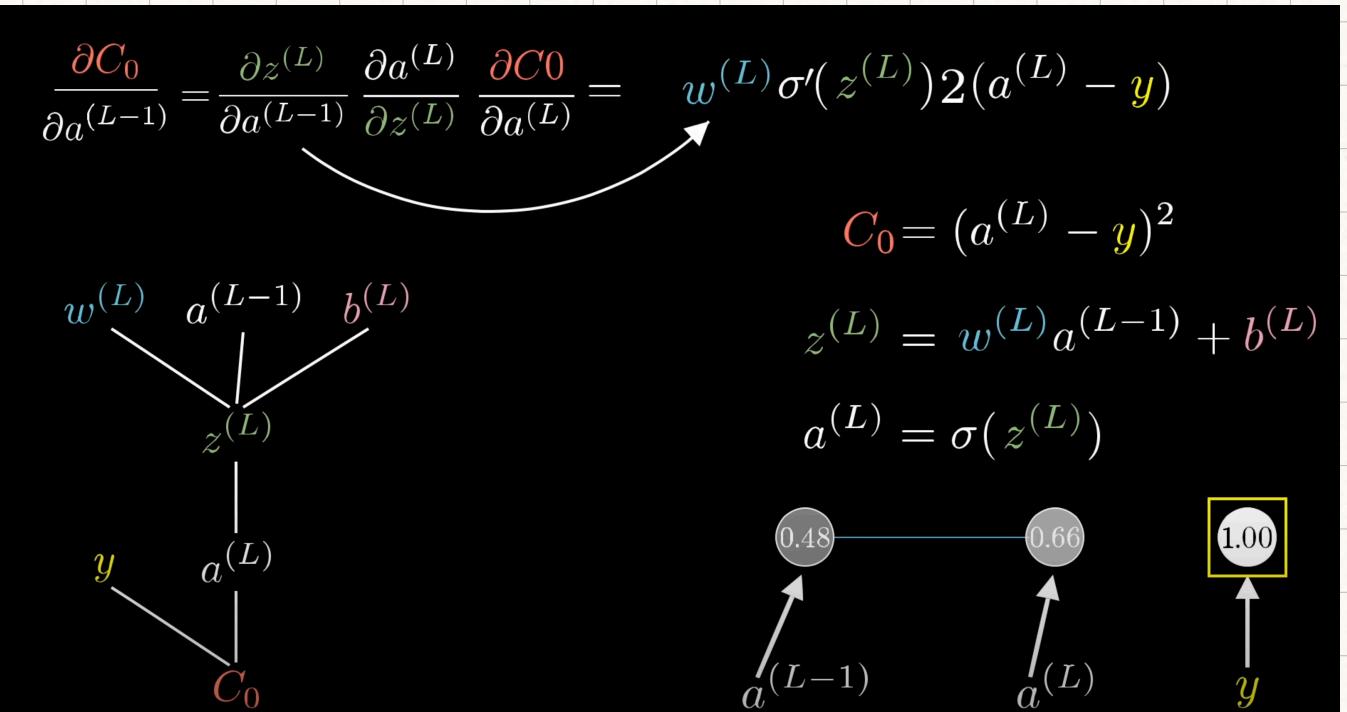
$$a^{(L)} = \sigma(z^{(L)})$$



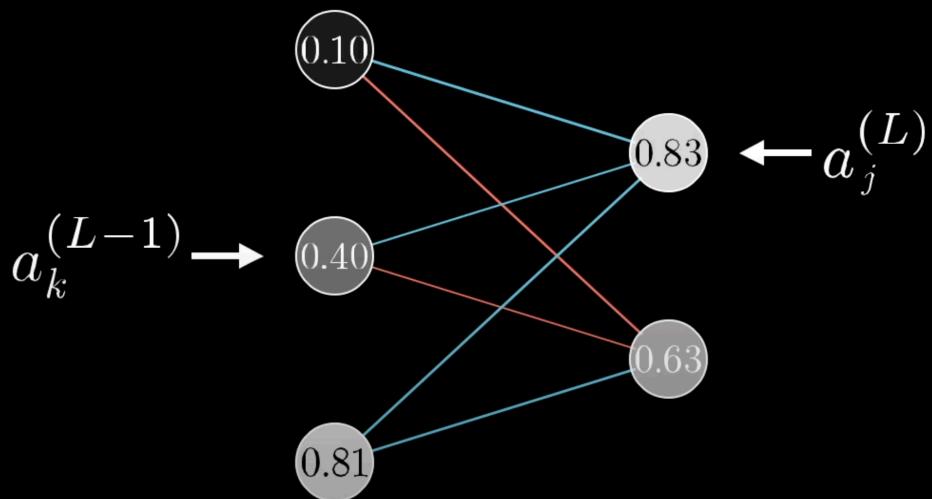
With bias  $\rightarrow$



With respect to back of previous layer ↗

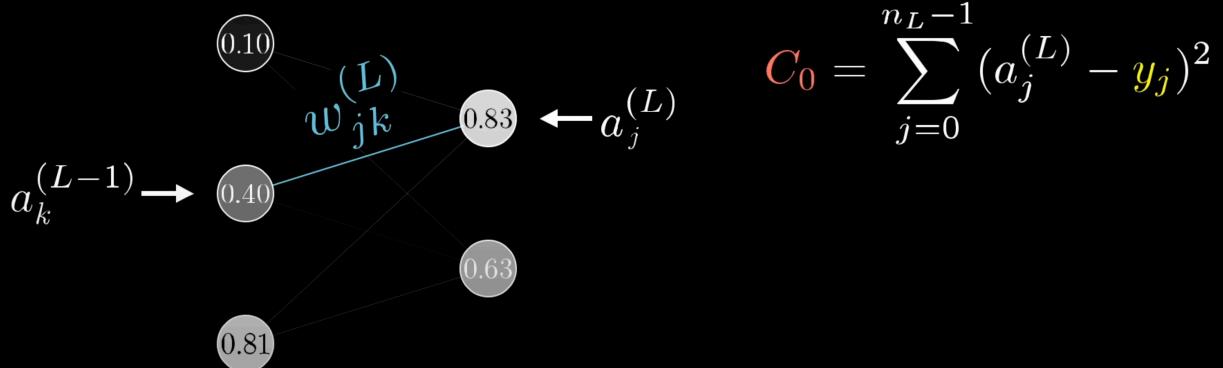


Now, what if multiple neurons?



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

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

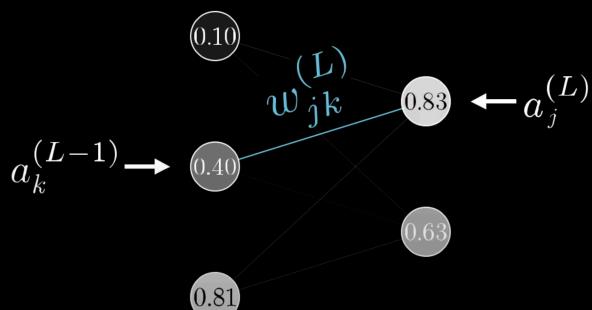


derivative is also very similar / same!

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

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

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



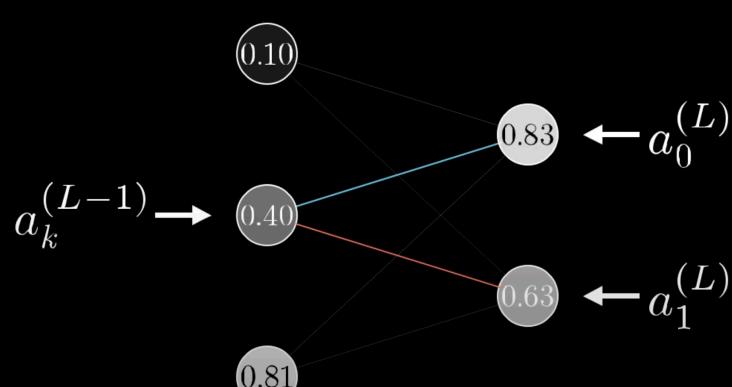
$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

However derivative for prior layer activation is slightly different:

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

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

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



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

NOK!

I plan to do another pass to beautify the notes.

Currently it's just raw.

# JUNK: ↴

-  → 28x28 pixels. ; low res
  - Brain can recognise 3's really easily ; crazy ;
- 

- Why the layers ; You understand "learning" in ML.
  - NN to recognise handwritten digits → 
- 

CNN → Image recognition

LSTM → Speech recognition.