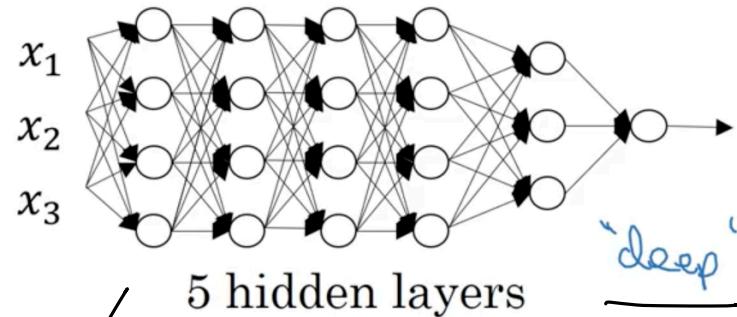
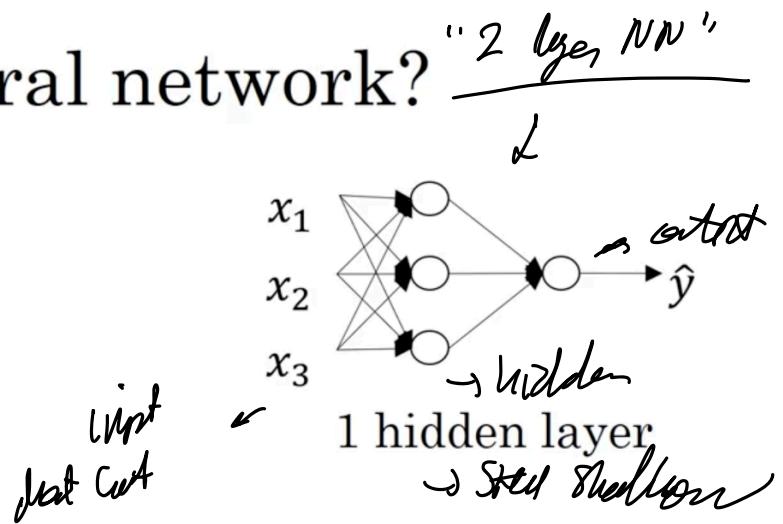
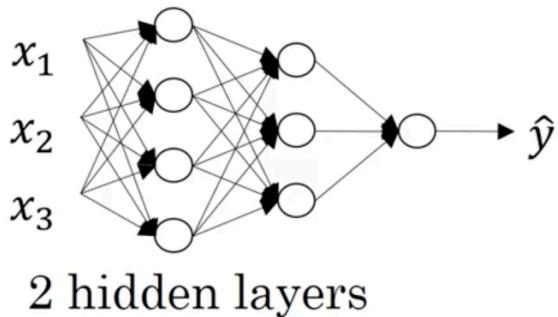
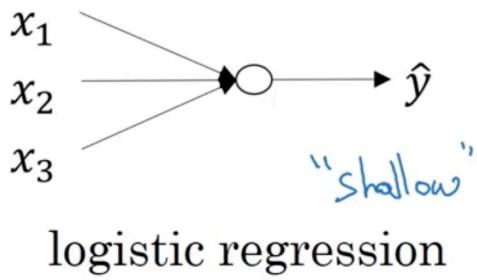


Deep L. Layer Neural Network

What is a deep neural network? "2 layer NN"



↗
deep layer nn

can recognize

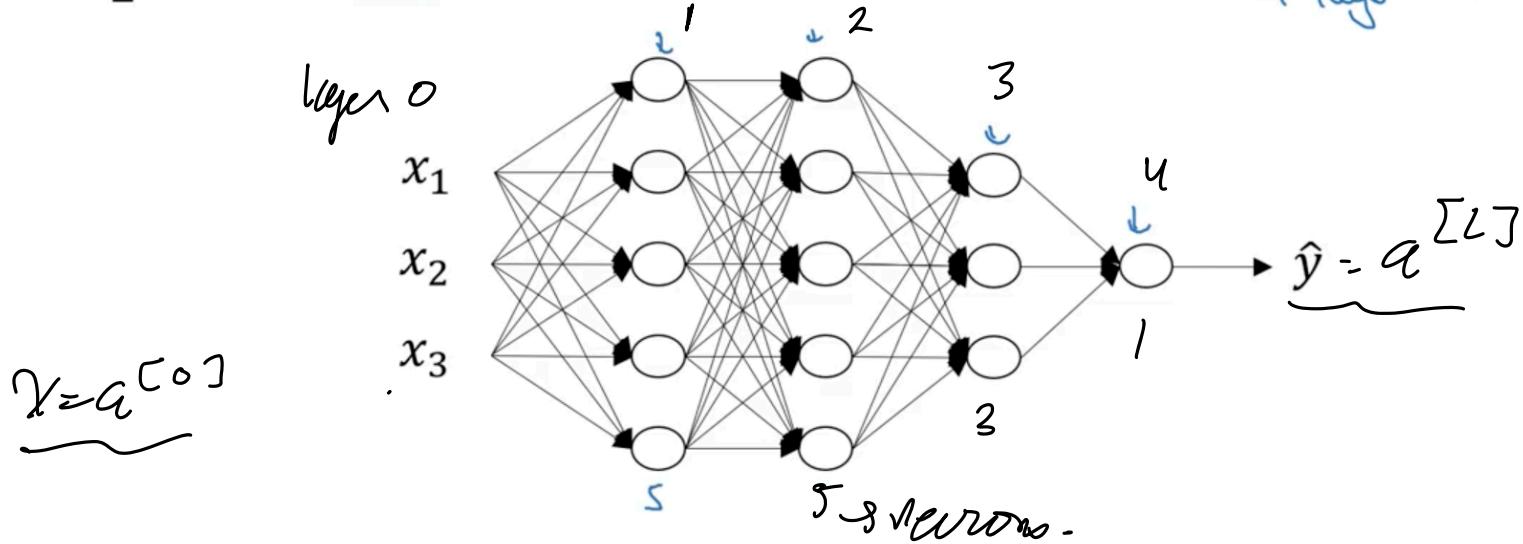
intensive feature

ones can't

Andrew

Deep neural network notation

4 layer NN



$L = 4$ (# of layers)

$n^{[l]}$ = # units in layer l (neurons in layer l)

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = 1$$

$$n^{[0]} = n_x = 3 \rightarrow \text{inputs}$$

$$n^{[L]} = 1$$

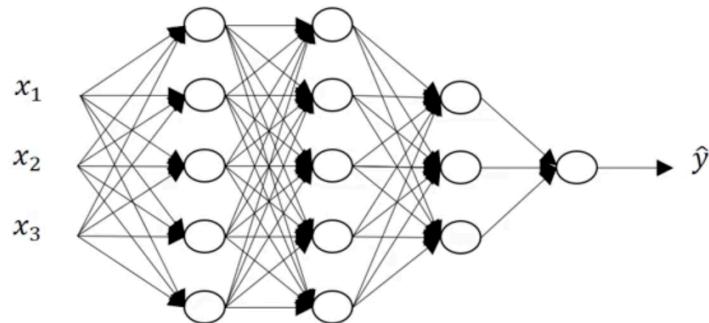
$a^{[l]}$ = Activations in layer l

$$a^{[l]} = g^{[l]}(z^{[l]}), \quad w^{[l]} = \text{weights for } z^{[l]}$$

$$b^{[l]} = \text{biases for } z^{[l]}$$

FORWARD PROPAGATION IN A DEEP NETWORK

Forward propagation in a deep network



Giving fixed training example x :

$$x: z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

⋮

$$z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]})$$

$$= \hat{y}$$

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

$$a^{[L]} = g^{[L]}(z^{[L]})$$

$\longrightarrow m$

VECTORIZED

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

:

$$\hat{y} = g(Z^{[4]}) = A^{[4]}$$

$$A^{[4]}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$\longrightarrow m$

$$A^{[1]} = \begin{bmatrix} & \\ & \\ n^{[1]} & \end{bmatrix}$$

comes from

for loop

for the different

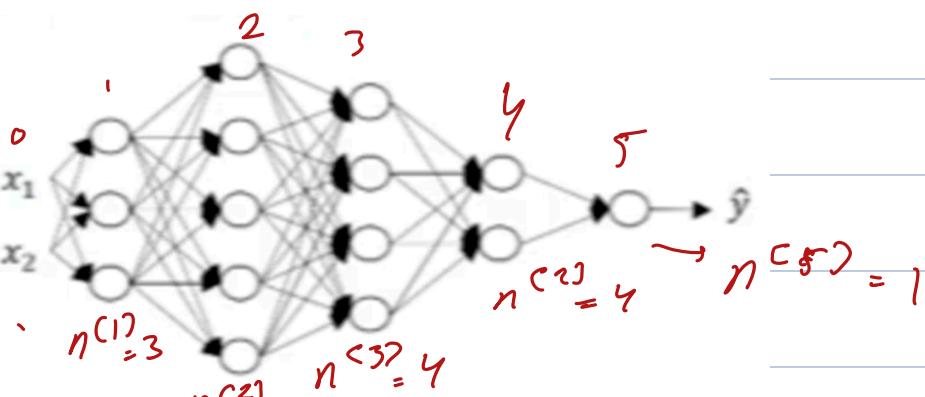
types -> So one need to use
after loop here.

Getting MATRIX DIMENSIONS RIGHT

Parameters $W^{[L]}, b^{[L]}$

$$L=5$$

$$n^{[0]} = n_x = 2$$



$$n^{[0]} = n_x = 2, n^{[1]} = 3, n^{[2]} = 4, n^{[3]} = 4, n^{[4]} = 1$$

~~$n^{[l-1]}$~~

$$z^{[1]} = w^{[1]} x + b^{[1]} \quad \xrightarrow{\text{grave for now.}}$$

↓

$$\begin{array}{c} (3,1) \\ + \\ (2,1) \\ + \\ (n^{[l]}, 1) \\ \hline (m) \end{array}$$

↓

$$(3,2)$$

$$(n^{[1]}, n^{[2]})$$

$$w^{[1]} = (n^{[1]}, n^{[0]})$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}. \quad w^{[l]} = (n^{[l]}, n^{[l-1]})$$

$$b^{[l]} = (n^{[l]}, 1)$$

Stay the same.

$$\begin{aligned} w^{[l]} : (n^{[l]}, n^{[l-1]}) &= dw^{[l]} \\ b^{[l]} : (n^{[l]}, 1) &= db^{[l]}. \end{aligned}$$

$$a^{[l]} = f^{[l]}(z^{[l]}) : a^{[l]} f^{[l]} = (n^{[l]}, 1)$$

$$\Rightarrow (n^{[l]}, m)$$

der Verteilung

$$Z^{[L]} = A^{[L]} = (n^{[L]}, m) \xrightarrow{\text{using examples}}$$

↳ first with size $(n^{[L]}, 1)$ and get broadcasted.

$$z^{[L]}, q^{[L]}: (n^{[L]}, 1)$$

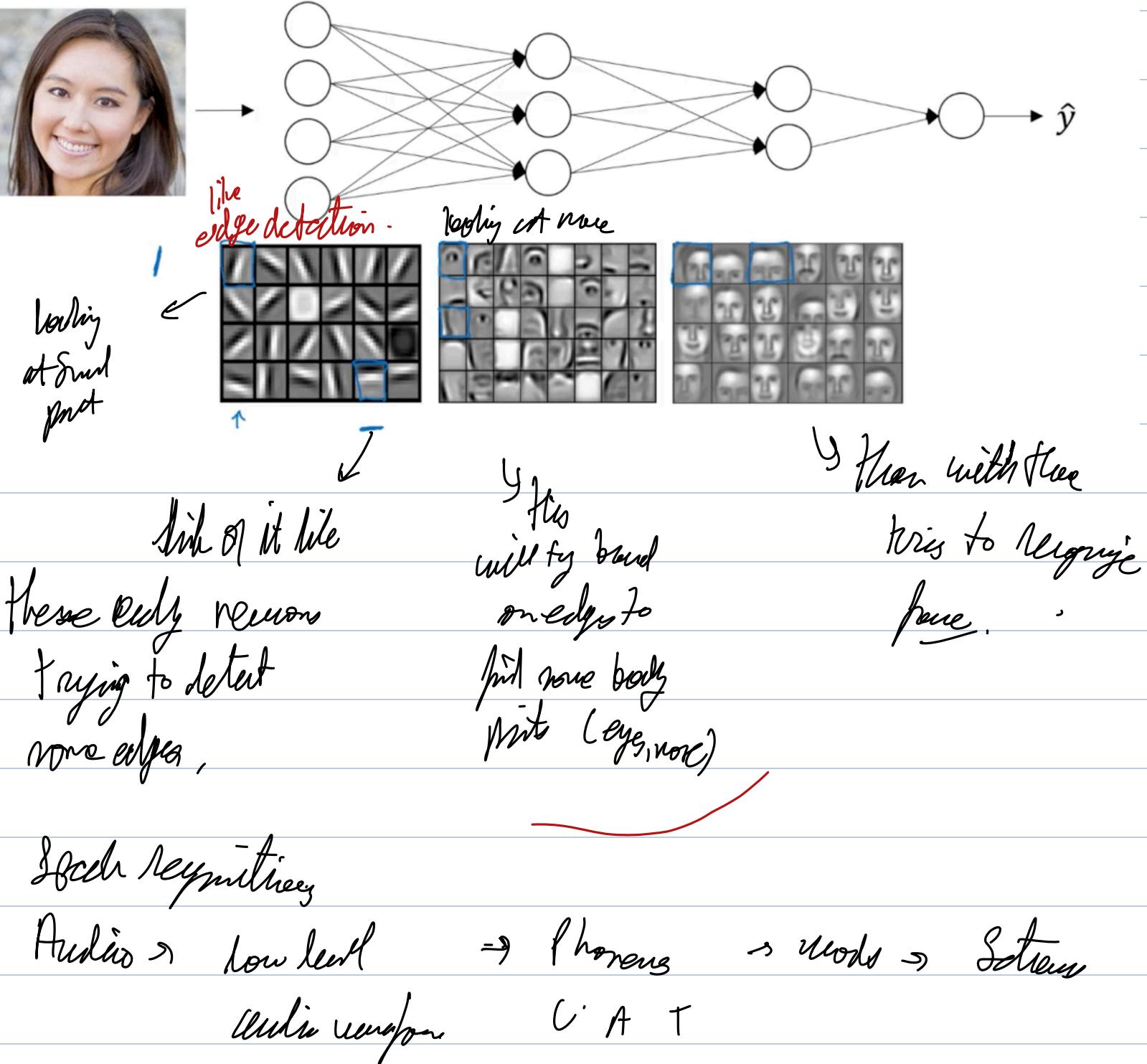
$$Z^{[L]}, A^{[L]}: (n^{[L]}, m)$$

$$l=0 \Rightarrow A^{(0)} = x = (n_x, m)$$

$$dZ^{[L]}, dA^{[L]}: (n^{[L]}, m)$$

WITH DEEP REPRESENTATIONS

Intuition about deep representation

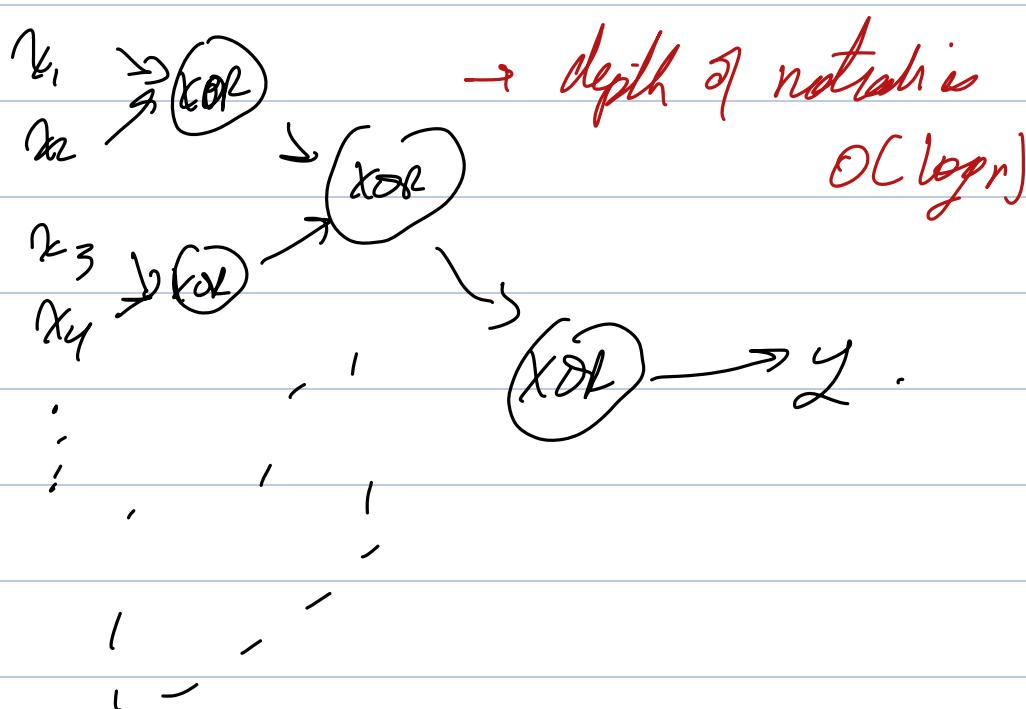


Circuit theory and deep learning

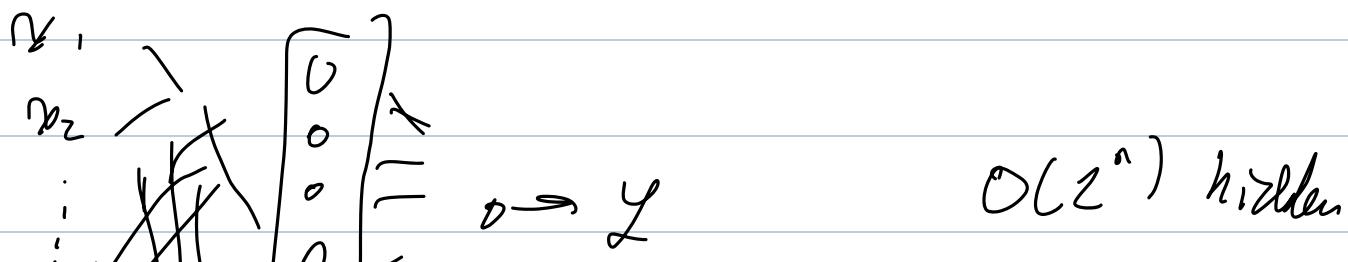
Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

↑
the middle
units / neurons
is small.

Eg. $x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$.



Partial not allowed and can only use 1 layer.

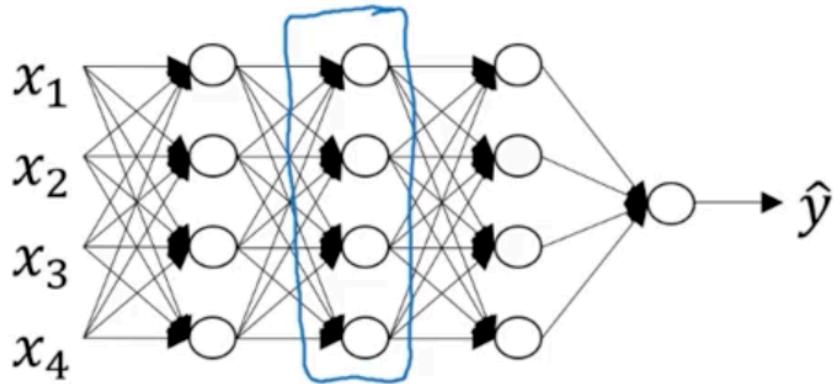


units

y exponentially large

BUILDING BLOCKS OF NEURAL NETWORKS

Forward and backward functions



Layer l : $w^{[l]}, b^{[l]}$

Forward: Input $a^{[l-1]}$, output $a^{[l]}$

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \Rightarrow \text{Cache } z^{[l]} \text{ for backprop}$$

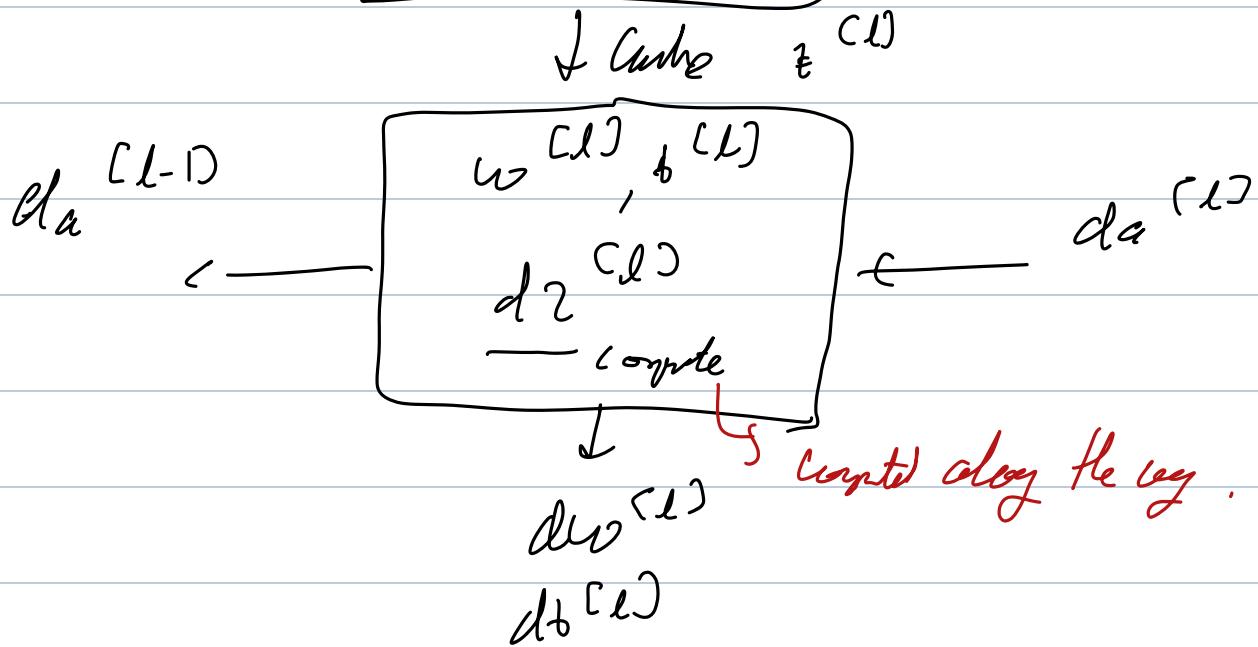
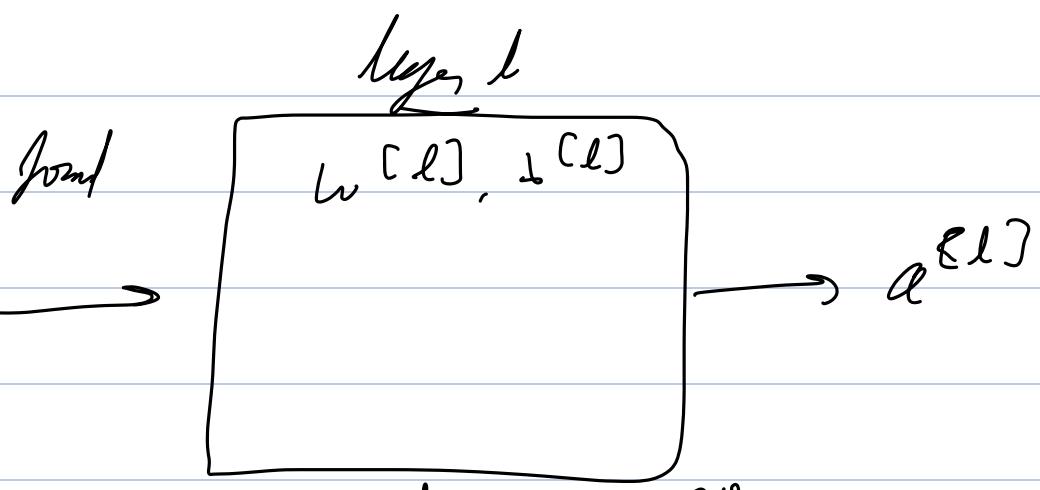
$$a^{[l]} = g^{[l]}(z^{[l]})$$

Backward: Input $da^{[l]}$, output $da^{[l-1]}$

Cache $(z^{[l]})$,

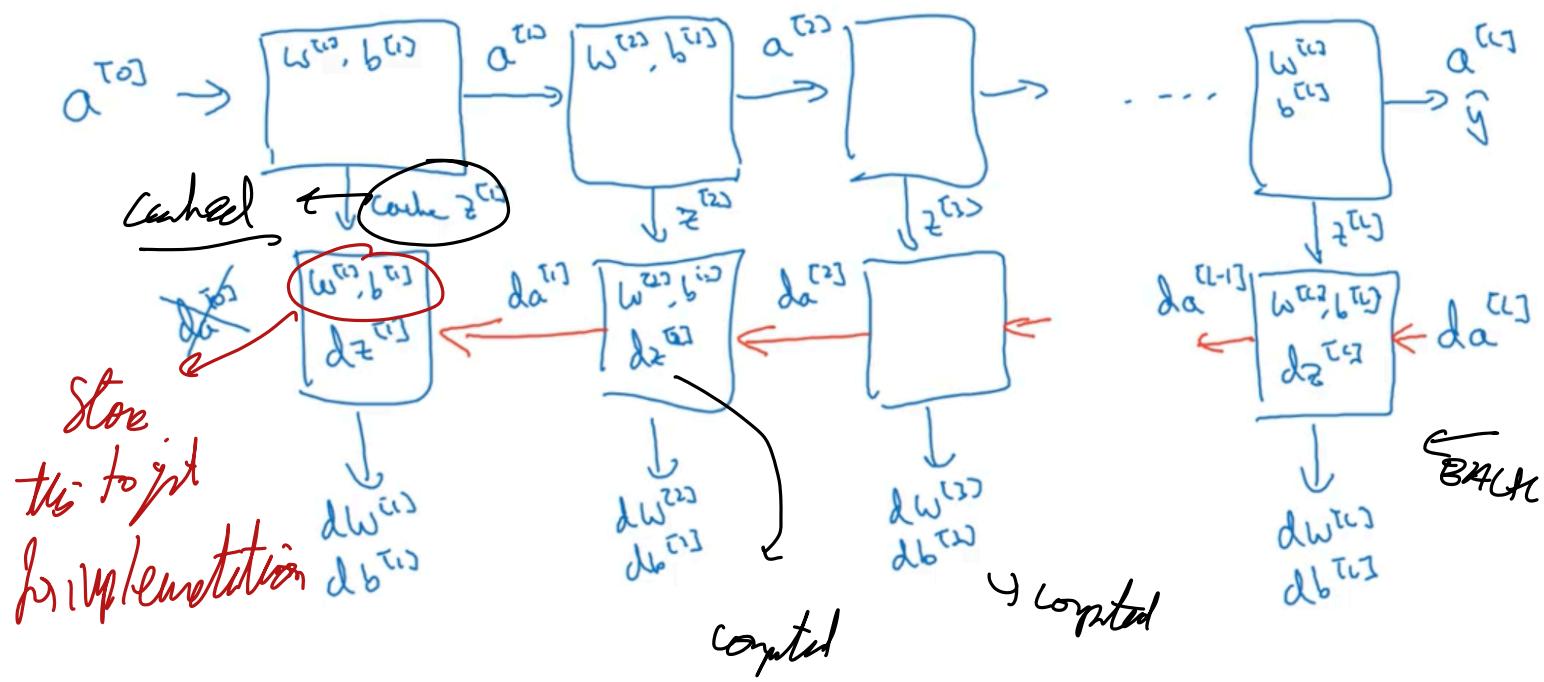
$$dw^{[l]}$$

$$db^{[l]}$$



Forward and backward functions

→ FORWARD



day the way.

Iteration roughly forward and back, and then:

$$w^{(l)} := w^{(l)} - \alpha d w^{(l)}$$

$$b^{(l)} := b^{(l)} - \alpha d b^{(l)}$$

Forward and Backward Propagation

Forward for layer l.

→ Input $a^{(l-1)}$ output: $a^{(l)}$, cache ($z^{(l)}$)
 ^{with activations}
 ^{using $w^{(l)}$, $b^{(l)}$}

$$\begin{aligned} z^{(l)} &= w^{(l)} a^{(l-1)} + b^{(l)} \\ a^{(l)} &= g^{(l)}(z^{(l)}) \end{aligned}$$

Vectorized

$$\begin{cases} z^{(l)} = w^{(l)} A^{(l-1)} + b^{(l)} \\ A^{(l)} = g^{(l)}(z^{(l)}) \end{cases}$$

$a^{(0)}$ ← initial example

$$A^{(0)} \quad X = A^{(0)} \rightarrow [\rightarrow] \rightarrow [\rightarrow] \rightarrow [\rightarrow] \rightarrow$$

→

Backward at layer l

Input $da^{[l]}$ output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$

$$dW^{[l]} = dz^{[l]} * a^{[l-1]T}$$

$$db^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = W^{[l]} + dz^{[l]}$$

↓ Vectorized is the same idea

.

$$dz^{[l]} = \underbrace{dA^{[l-1]}}_{\text{Input}} * g^{[l]'}(z^{[l]})$$

$$\underbrace{dW^{[l]}}_{\text{sum}} = \frac{1}{m} \sum dz^{[l]} * A^{[l-1]T}$$

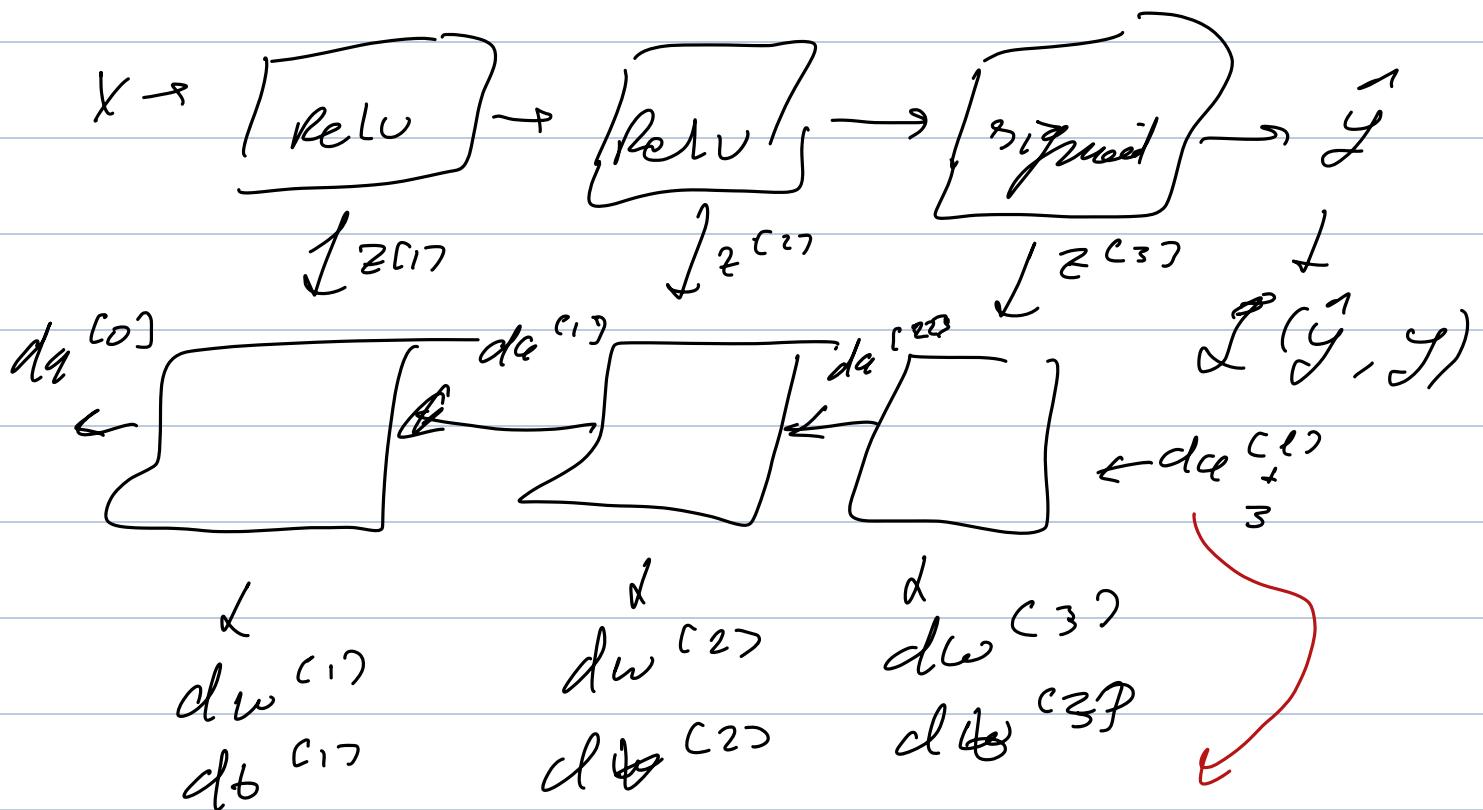
$$\underbrace{db^{[l]}}_{\text{sum}} = \frac{1}{m} \sum \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdim=True})$$

$$dA^{[l-1]} = W^{[l]T} * dz^{[l]}$$

→ to next step.

→ what you need for gradient descent,

Forward ↗



$$dA^{[1]} = \left(-\frac{y^{(1)}}{a^{(1)}} + \frac{(1-y^{(1)})}{(1-a^{(1)})} \right) da^{[1]} = \frac{-y}{a} + \frac{(1-y)}{(1-a)}$$

$$\dots \frac{y^{(m)}}{a^{(m)}} + \frac{(1-y^{(m)})}{(1-a^{(m)})}$$

Really the derivative of $\frac{dy}{da}$ ↗ a

$$\frac{dy}{da} \rightarrow a$$

does this
if this
is Sigmoid
at the end

PARAMETERS VS HYPERPARAMETERS

~~W^(l)~~ is thought to optional ready or FFNN is depth ~~W^(l)~~ in CNTK 4.

PARAMETERS VS HYPERPARAMETERS

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \dots$

Hyperparameters: - learning rate α .

- # of iterations

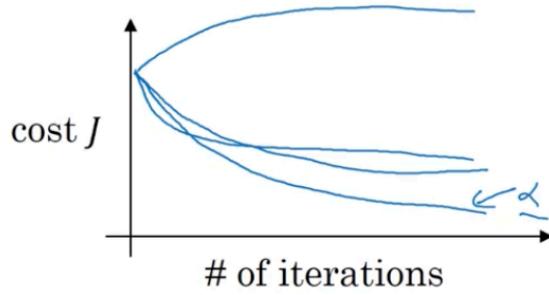
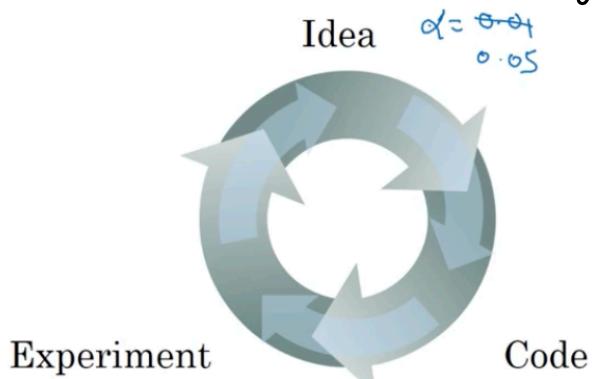
- # hidden layers L

- # hidden units $n^{(1)}, n^{(2)} \Rightarrow$ neurons per layer

- choice of activation fun.

Others: Momentum, mini-batch size etc.

Applied deep learning is a very empirical process \rightarrow try and see what works.



0 1 2 3 4,

[n_x, y, z, 2, 1]

w⁽¹⁾ = < n_x⁽¹⁾, n_y⁽¹⁾, n_z⁽¹⁾ >

i = 0 → [0, 5)

w₁ = degc(1).