# Deep - < layer N·N

→ ~~show~~ shallow Vs Deep is a matter of depth ie no of hidden layers

→ Logistic reg. is a 1 layer N·N

→

# Forward propagation in Deep Network

For a single training eg:-
for 1st hidden layer

→ We $z^{[1]} = W^{[1]} X + b^{[1]}$
→ $a^{[1]} = g^{[1]}(z^{[1]})$

For 2nd hidden layer
$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
$a^{[2]} = g^{[2]}(z^{[2]})$

..... and so on

→ In O/P layer
$z^{[last]} = W^{[last]} a^{[last-1]} + b^{[last]}$
$a^{[last]} = g(z^{[last]}) = \hat{y}$

general $z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$
rule: $a^{[l]} = g^{[l]}(z^{[l]})$

## If vectorized

$$Z^{[1]} = W^{[1]} X + b^{[1]} \qquad (X = A^{[0]})$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
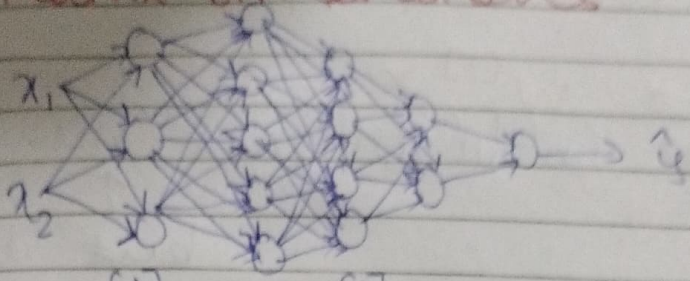$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
$$A^{[2]} = g^{[2]}(Z^{[2]})$$

..... and so on

for
$l=1$ to
↑ last

no
way
to eq
do
without
for loop

$$Z^{[2]} = \begin{bmatrix} & \vdots & & \vdots & \\ & Z^{[2](1)} & \cdots & Z^{[2](m)} & \\ & \vdots & & \vdots & \end{bmatrix}$$

# Matrix Dimensions



$n^{[0]} = 2$
$n^{[1]} = 3$
$n^{[2]} = 5$
$n^{[3]} = 4$
$n^{[4]} = 2$
$n^{[5]} = 1$

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

$$z^{[1]} = (3,1) = (n^{[1]}, 1) \quad b = (n^{[1]}, 1)$$
$$X = (2,1) = (n^{[0]}, 1) \quad = (3,1)$$
$$W = (3,2) = (n^{[1]}, n^{[0]})$$

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

$$\begin{bmatrix} . \\ . \\ . \end{bmatrix} = \begin{bmatrix} . . \\ . . \\ . . \end{bmatrix} \begin{bmatrix} . \\ . \end{bmatrix} \quad \begin{bmatrix} . . \\ . . \\ . . \end{bmatrix} \begin{bmatrix} . \\ . \\ . \end{bmatrix}$$

$$\underset{\underset{(5 \times 1)}{\uparrow}}{z^{[2]}} = \underset{\underset{(5 \times 3)}{\uparrow}}{W^{[2]}} \underset{\underset{(3,1)}{\uparrow}}{a^{[1]}} + \underset{\underset{(5,1)}{\uparrow}}{b^{[2]}}$$

→ $W^{[l]} = (n^{[l]}, n^{[l-1]})$

→ dw shld be same dimension as W

ie $dw^{[l]} = (n^{[l]}, n^{[l-1]})$

→ $db^{[l]} = [n^{[l]} - 1]$ — same as b

→ $a^{[l]} = g^{[l]}(z^{[l]})$

∴ z & a shld hv same dimension $= (n^{[l]}, 1)$

for m examples

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$
$$(n^{[1]}, m) \quad (n^{[1]}, n^{[0]})(n^{[0]}, m) \quad (n^{[1]}, 1)$$

↑

but by
python boroadcasting
it is duplicated to
$(n^{[1]}, m)$

$\therefore Z^{[\ell]}, A^{[\ell]} : (n^{[\ell]}, m)$

$dZ^{[\ell]}, dA^{[\ell]} : (n^{[\ell]}, m)$ (same
as $Z$ & $A$)

# Why deep Networks work well
→ Detects simple things like edges
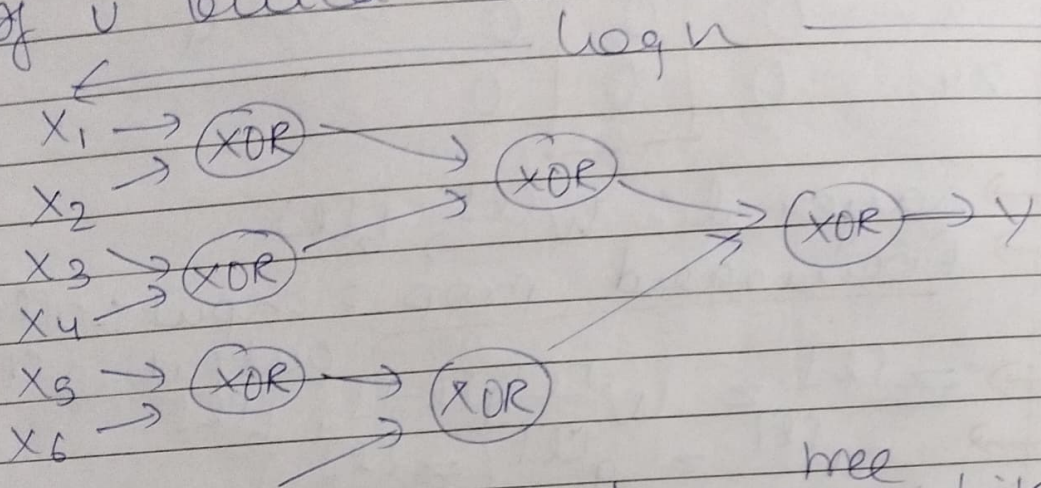first and then builds them
into more complex notes
things.

★ Circuit theory and deep learning
There are functions you can
compute with a "small"
$L$-layer deep neural network
that shallower networks require
exponentially more hidden
units to compute

→ Circuit theory noctures the
thinking about what types
of functions you can
compute with different

## AND, OR and NOT gates.

→ If we r trying to compute $x_1$ XOR, $x_2$ XOR ..... $x_n$ XOR if we have n features

If u build XOR tree like this:

← ———————— log n ————————→

$x_1 \rightarrow$ (XOR)
$x_2 \rightarrow$
           → (XOR)
$x_3 \rightarrow$ (XOR)
$x_4 \rightarrow$
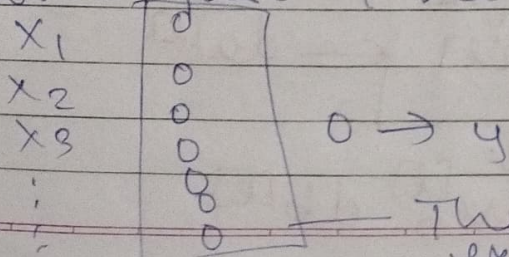                    → (XOR) → y
$x_5 \rightarrow$ (XOR) → (XOR)
$x_6 \rightarrow$

we can build a ckt (tree) like this

→ To compute XOR The depth of network will be on the order log n

→ no of nodes/ckt components/gates in the network is not that large. We dont need that many gates in order to compute exclusive OR

→ If we are forced to compute with just 1 hidden layer

$x_1$   0
$x_2$   0
$x_3$   0    0 → y
⋮     0
⋮     0   This layer will be ($2^{n-1}$) exponentially large

# Building blocks of deep neural Networks

$x_1$    ○ | ○ | ○

$x_2$    ○ | ○ | ○

$x_3$    ○ | ○ | ○    ○ → $\hat{y}$

$x_4$    ○ | ○ | ○

→ Layer $l$: $W^{[l]}$, $b^{[l]}$
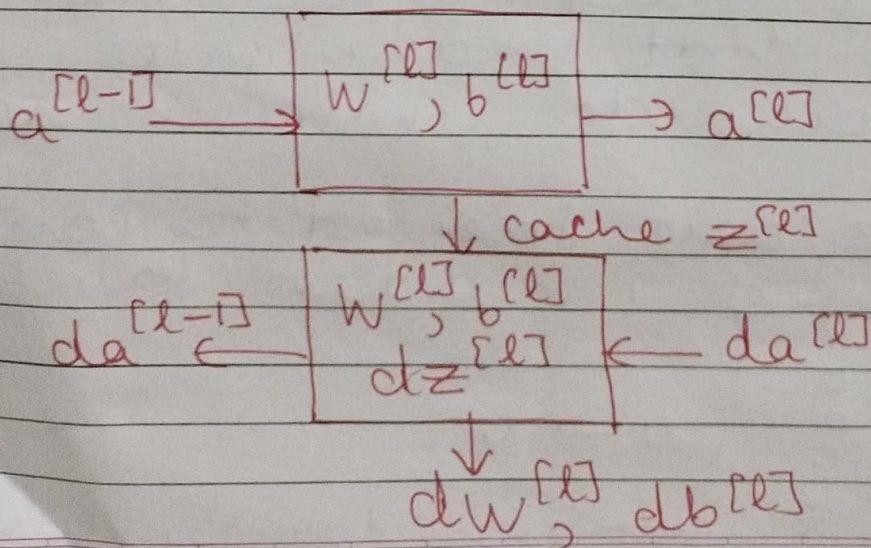
★ → <u>Forward prop</u>: Input: $a^{[l-1]}$
           output $a^{[l]}$

→ $z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$

→ $a^{[l]} = g^{[l]}(z^{[l]})$

→ also cache $z^{[l]}$ as storing $z^{[l]}$ is useful for back propagation

★ → <u>Backward prop</u>: Input $da^{[l]}$, cache
           output $da^{[l-1]}$

Layer $l$

$a^{[l-1]}$ → | $W^{[l]}, b^{[l]}$ | → $a^{[l]}$

↓ cache $z^{[l]}$

$da^{[l-1]}$ ← | $W^{[l]}, b^{[l]}$ <br> $dz^{[l]}$ | ← $da^{[l]}$

↓

$dw^{[l]}, db^{[l]}$

$$a^{[0]} \rightarrow \boxed{W^{[1]}, b^{[1]}} \xrightarrow{a^{[1]}} \boxed{W^{[2]}, b^{[2]}} \xrightarrow{a^{[2]}}$$

$$\downarrow z^{[1]} \qquad \qquad \downarrow z^{[2]}$$

$$\boxed{\phantom{WWWW}} \rightarrow \cdots \cdots \boxed{\begin{array}{c} W^{[l]} \\ b^{[l]} \end{array}} \rightarrow \begin{array}{c} a^{[l]} \\ \hat{y} \end{array}$$

$$\downarrow z^{(3)}$$

$$\downarrow$$

$$\boxed{\begin{array}{c} W^{[1]}, b^{[1]} \\ dz^{[1]} \\ \downarrow \\ dw^{[1]} \\ db^{[1]} \end{array}} da^{[1]} \boxed{\begin{array}{c} W^{[2]} \\ b^{[2]} \\ dz^{[2]} \\ \downarrow \\ dw^{[2]}, \end{array}} \begin{array}{c} da^{[2]} \\ \\ \downarrow \\ dw^{[3]} \\ db^{[3]} \end{array} \boxed{\phantom{WW}} \leftarrow \cdots \cdots \leftarrow \boxed{\begin{array}{c} W^{[l]}, b^{[l]}, \\ dz^{[l]} \\ \downarrow \\ dw^{[l]} \\ db^{[l]} \end{array}} \begin{array}{c} da^{[l]} \\ \leftarrow \end{array}$$

$$W^{[l]} = W^{[l]} - \alpha dW^{[l]}$$

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$

$\rightarrow$ It is useful to store $z^{[1]}, W^{[1]}$ & $b^{[1]}$ .... in cache during forward prop to be used in backward prop.

# Forward and Backward Propagation

Forward prop. for layer $l$
- Input $a^{[l-1]}$ ____ $W^{[l]}, b^{[l]}$
- Output $a^{[l]}$, cache $(z^{[l]})$

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$
$$a^{[l]} = g^{[l]}(z^{[l]})$$

Vectorized

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$
$$A^{[l]} = g^{[l]}(z^{[l]})$$

Backward propagation for layer $l$
- Input $da^{[l]}$
- Output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$\rightarrow dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$
$$\uparrow$$
element wise multiplication

$$\rightarrow dW^{[l]} = dz^{[l]} a^{[l-1]T}$$
$$\rightarrow db^{[l]} = dz^{[l]}$$
$$\rightarrow da^{[l-1]} = W^{[l]T} dz^{[l]}$$

Vectorized :

$$dz^{[l]} = dA^{[l]} * g^{[l]'}(z^{[l]})$$
$$dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$
$$db^{[l]} = \frac{1}{m} np.sum(dz^{[l]}, axis=1, keepdims=True)$$

$$dA^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$$

→ In our forward prop for layer 0 we pass in the parameters x

→ Similarly for backward prop. we pass $da^{[l]}$ which is

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

# Parameters Vs Hyperparameters

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

Hyperparameters :
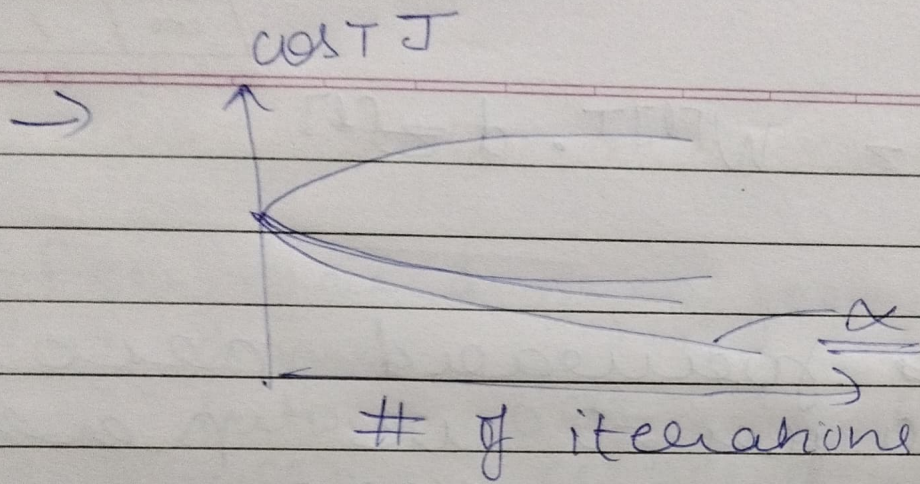→ learning rate $= \alpha$
→ Iterations
→ hidden layers. $l$
→ hidden units $n^{[1]}, n^{[2]} \ldots$,
    choice of activation functions

(Control the parameters)

These determine final value of parameters

* Applying deep learning is an emperical process

→ We keep trying diff values of cost function x and observe cost function

COST J



# of iterations

→ U Just have to try out diff values of hyperparameters and see what works best It is an emperical process