## Deep Neural Network (L layered)



$x_1$
$x_2$ → O → $\hat{y}$
$x_3$

Logistic → 1 Layer
Regression NN
[SHALLOW]



$x_1$
$x_2$
$x_3$

1 hidden layer
(NN) → 2 layer
[DEEP]

\* 4 layer NN = 3 hidden layers

No. of nodes:   Layer 1 = 5

| | |
|---|---|
| Layers = L (4) | |
| No. of nodes in layer $l$. | $= n^{[l]}$ |

Layer 2 = 5
Layer 3 = 3
Layer 4 = 1

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

$n^{[L]} = n^{[4]} = 1$

| | |
|---|---|
| $a^{[l]}$ = activations of layer $l$. | → $a^{[l]} = g^{[l]}\left(z^{[l]}\right)$ |
| $w^{[l]}$ = weights for $z^{[l]}$ | |
| $b^{[l]}$ = biases. | |

## ○ FORWARD PROPAGATION

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

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

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

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

$\vdots$

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

$a^{[4]} = \hat{y} = g^{[4]}\left(z^{[4]}\right)$

**vectorized**

$Z^{[1]} = w^{[1]}X + b^{[1]}$   for $l=1$ to $L$

$A^{[1]} = g^{[1]}\left(z^{[1]}\right)$

$Z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$   Here

$A^{[2]} = g^{[2]}\left(z^{[2]}\right)$   $L = 4$

$\vdots$

$\hat{y} = A^{[4]} = g^{[4]}\left(z^{[4]}\right)$

## • MATRIX DIMENSIONS

| | |
|---|---|
| $L = 5$ | $z^{[1]} = w^{[1]} x + b^{[1]}$ |
| $n^{[1]} = 3$ | $(3,1) = (3,2)(2,1) + (3,1)$ |
| $n^{[2]} = 5$ | $z^{[1]} = (n^{[1]}, 1)$ |
| $n^{[3]} = 4$ | $w^{[1]} = (n^{[1]}, n^{[0]})$ |
| $n^{[4]} = 2$ | $x = (n^{[0]}, 1)$ |
| $n^{[5]} = 1$ | $b = (n^{[1]}, 1)$ |

$x_1$
$x_2$

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

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

$0 \to \hat{y}$

$\emptyset$

## *vectorized

$$\underset{\underset{(n^{[1]}, m)}{\uparrow}}{Z^{[1]}} = \underset{\underset{(n^{[1]}, n^{[0]})}{\uparrow}}{w^{[1]}} \underset{\underset{(n^{[0]}, m)}{\uparrow}}{x} + \underset{(n^{[1]}, 1)}{b^{[1]}} \xrightarrow[\text{cast}]{\text{broad}} (n^{[1]}, m)$$

---

$$Z^{[l]}, A^{[l]} : (n^{[l]}, m) = dZ^{[l]}, dA^{[l]}$$

when $l = 0 \qquad A^{[0]} = X = (n^{[0]}, m)$

---

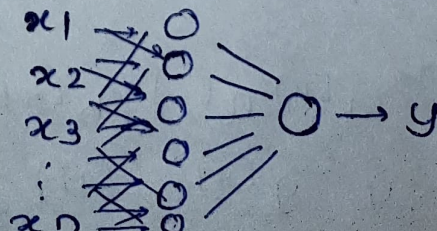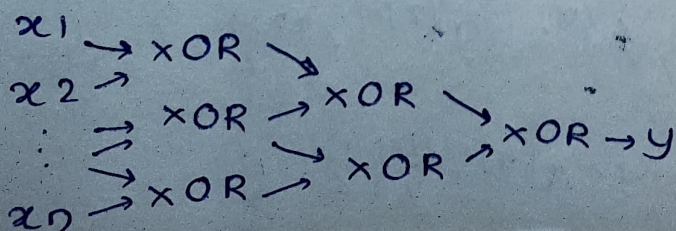## • DEEP NETWORK Importance

As we progress through layers of networks,
the features detected go from  simple  to
complex.

eg: $y = x_1$ XOR $x_2$ XOR .... XOR $x_n$ (LOGICAL)

Deep network $[\log n]$          one hidden layer
layers                              $[2^{n-1}]$ nodes

$x_1 \to$ XOR
$x_2 \to$ XOR $\to$ XOR
$\vdots \to$ XOR $\to$ XOR $\to$ XOR $\to y$
$x_n \to$ XOR

$x_1$
$x_2$
$x_3$
$\vdots$
$x_n$
$\to y$

# • FORWARD & BACKWARD PROPAGATION

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

→ <u>Forward</u> : input $a^{[l-1]}$, output $a^{[l]}$
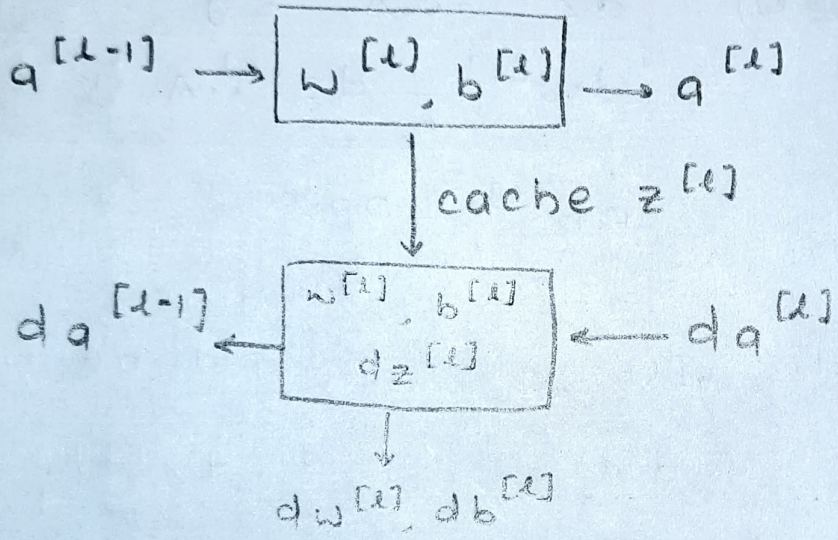
$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$, cache $z^{[l]}$

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

→ <u>Backward</u> : input $da^{[l]}$, output $da^{[l-1]}$

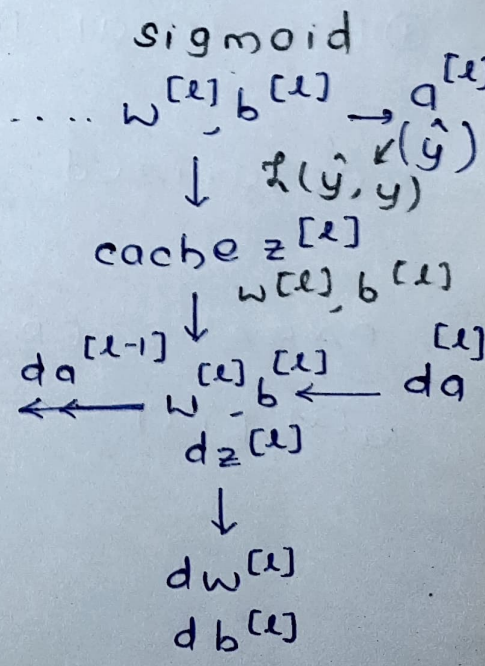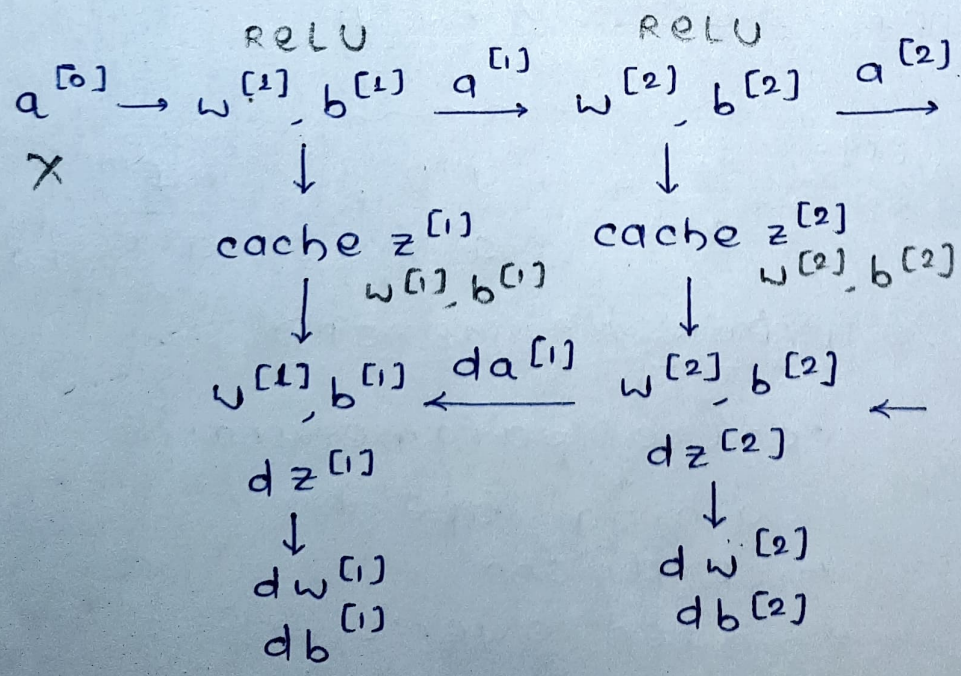~~cache~~     cache $(z^{[l]})$     $dw^{[l]}$

$db^{[l]}$

layer $l$

$a^{[l-1]} \longrightarrow \boxed{w^{[l]}, b^{[l]}} \longrightarrow a^{[l]}$

$\downarrow$ cache $z^{[l]}$

$da^{[l-1]} \longleftarrow \boxed{\begin{array}{c} w^{[l]}, b^{[l]} \\ dz^{[l]} \end{array}} \longleftarrow da^{[l]}$

$\downarrow$

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

(cache)
⇕
passes info
from forward
to back
propagation
step.

RELU

$a^{[0]} \longrightarrow w^{[1]}, b^{[1]} \xrightarrow{a^{[1]}}$

$x \qquad \downarrow$

cache $z^{[1]}$

$\downarrow w^{[1]}, b^{[1]}$

$w^{[1]}, b^{[1]} \xleftarrow{da^{[1]}}$

$dz^{[1]}$

$\downarrow$

$dw^{[1]}$

$db^{[1]}$

RELU

$w^{[2]}, b^{[2]} \xrightarrow{a^{[2]}} \dots \dots$

$\downarrow$

cache $z^{[2]}$

$\downarrow w^{[2]}, b^{[2]}$

$w^{[2]}, b^{[2]} \longleftarrow$

$dz^{[2]}$

$\downarrow$

$dw^{[2]}$

$db^{[2]}$

sigmoid

$w^{[l]}, b^{[l]} \longrightarrow a^{[l]}$

$\downarrow \mathcal{L}(\hat{y}, y)$

cache $z^{[l]}$

$\downarrow w^{[l]}, b^{[l]}$

$da^{[l-1]} \downarrow$

$\longleftarrow w^{[l]}, b^{[l]} \longleftarrow da$

$dz^{[l]}$

$\downarrow$

$dw^{[l]}$

$db^{[l]}$

$\Rightarrow w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$

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

## • FORWARD PROPAGATION FOR LAYER $l$

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

output: $a^{[l]}$, cache $(z^{[l]})$

| VECTORIZED | $z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$ |

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

## • BACKWARD PROPAGATION FOR LAYER $l$

Input: $da^{[l]}$

output: $da^{[l-1]}, dw^{[l]}, db^{[l]}$     | VECTORIZED |

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

$\checkmark dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$    $dw^{[l]} = \dfrac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$

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

$\checkmark da^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$    $db^{[l]} = \dfrac{1}{m}$ np.sum

$dz^{[l]} = w^{[l+1]T} dz^{[l+1]} * g^{[l]'}(z^{[l]})$    $(dz^{[l]}, \text{axis}=1,$

         keep dims = True)

         | $dA^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$ |

⊛ For final layer : $da^{[l]} = -\dfrac{y}{a} + \dfrac{(1-y)}{(1-a)}$

vectorized: $dA^{[l]} = \dfrac{-y^{(1)}}{a^{(1)}} + \dfrac{(1-y^{(1)})}{(1-a^{(1)})} + \dots + \dfrac{-y^{(m)}}{a^{(m)}} + \dfrac{(1-y^{(m)})}{(1-a^{(m)})}$

## • PARAMETERS & HYPERPARAMETERS

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

            ↓
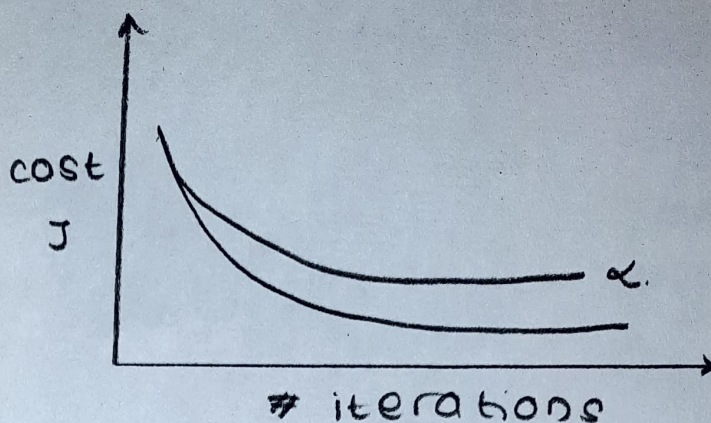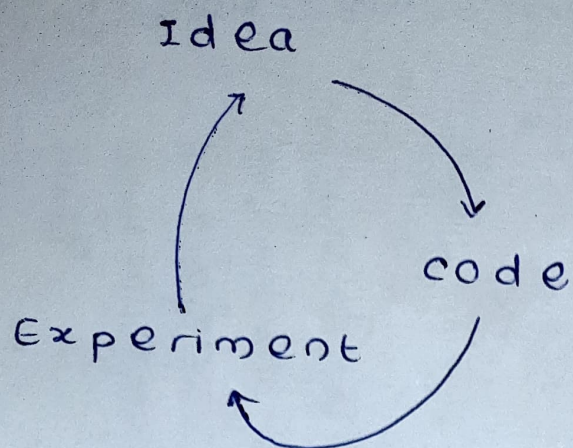
         control parameters

            ↓

         Learning rate $(\alpha)$

         # iterations

         # hidden layer $L$

         # hidden units $n^{[1]}, n^{[2]}, \dots$

         choice of activation function

Idea

↗ ↘

code

Experiment ↙



cost J

# iterations

$\alpha$

$\{IMP\}$ $dz^{[L]} = A^{[L]} - y$

$dw^{[L]} = \frac{1}{m} dz^{[L]} A^{[L-1]T}$

$db^{[L]} = \frac{1}{m} np.sum(dz^{[L]}, axis=1, keepdims=True)$

$dz^{[L-1]} = W^{[L]T} dz^{[L]} \circledast g'^{[L-1]}(z^{[L-1]})$

$\uparrow$ element wise multiplication

$\therefore dz^{[1]} = A^{[1]} - y$

$dw^{[1]} = \frac{1}{m} dz^{[1]} A^{[0]T}$   $\boxed{A^{[0]T} = X^T}$

$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$