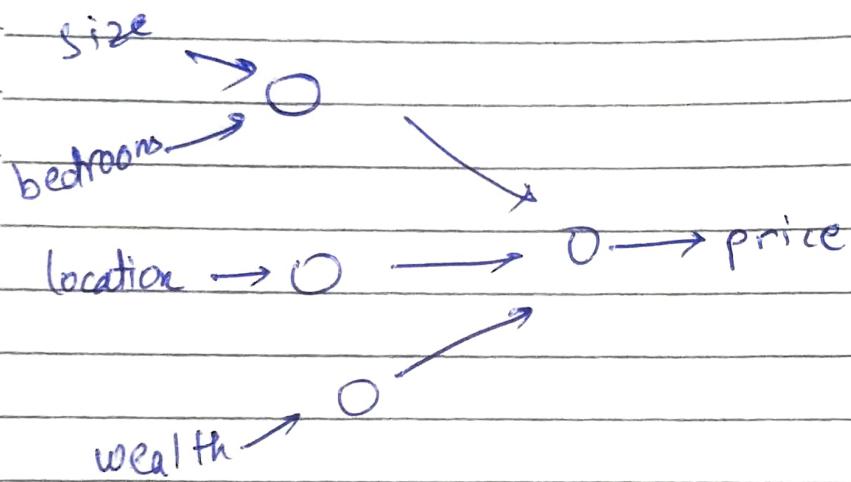
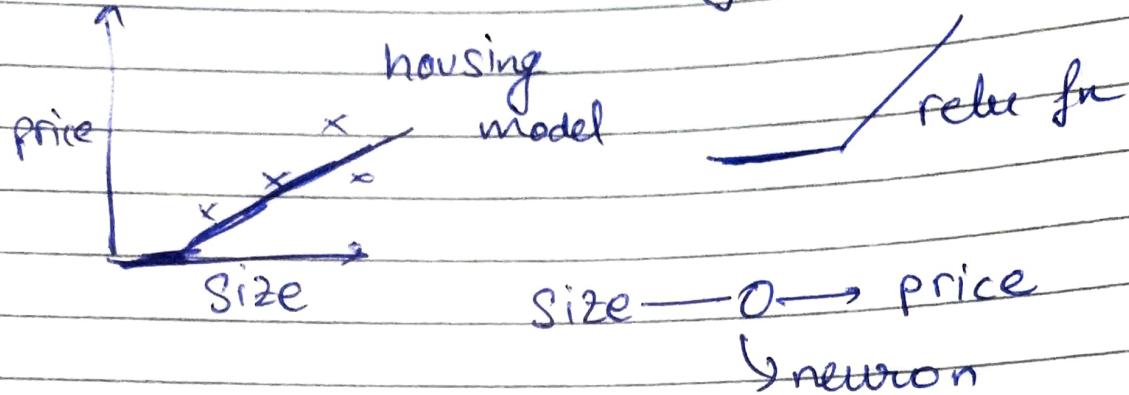
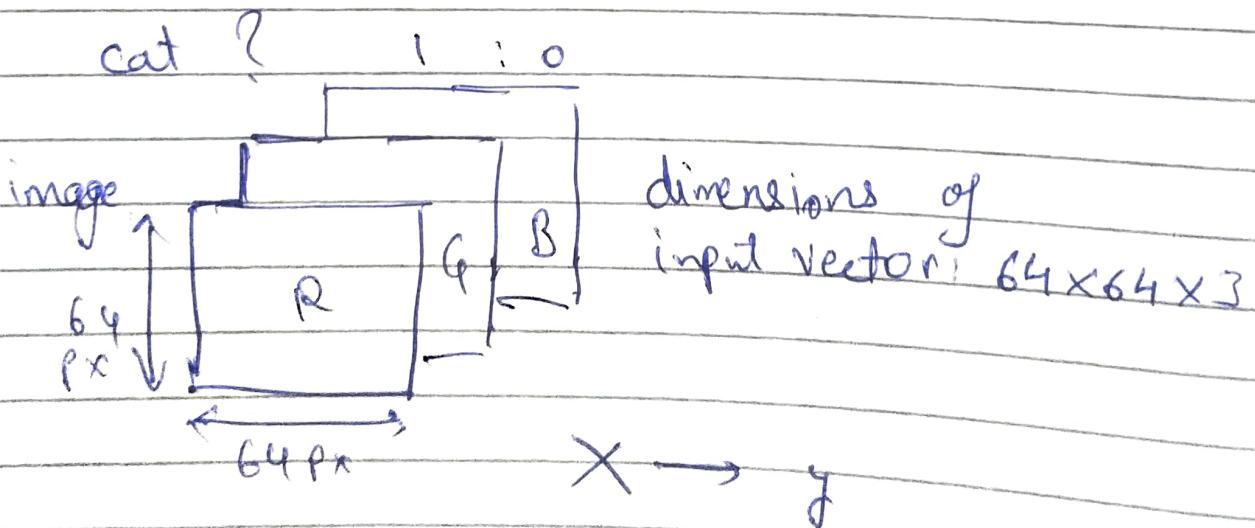


## NN & DL

### # Intro to Deep Learning



### # Binary Classification - Logistic Reg as NN



## Notation

$(x, y)$   $x \in \mathbb{R}^{n \times m}$ ,  $y \in \{0, 1\}$

$m$  training examples:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$\swarrow$  split  
 $m_{\text{train}}$      $m_{\text{test}}$

$$X = \begin{bmatrix} | & | & | & & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ | & | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times m}$$

## # Logistic Reg

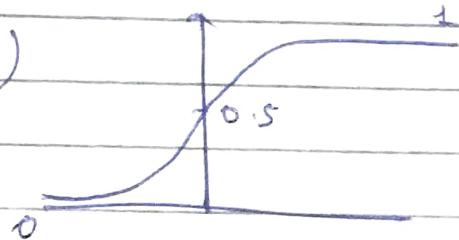
Given  $X$ , we want  $\hat{y} = P(y=1|x)$

$$x \in \mathbb{R}^{n \times n}$$

Parameters  $w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$

$$\text{Output } \hat{y} = \sigma(w^T x + b)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



## Alternative notation of assigning

$$x_0 = 1$$

$$\hat{y} = \sigma(\theta^T x)$$

does not work very well  
in case of NNs

\* Cost Function

$$\text{we want } \hat{y}^{(i)} \approx y^{(i)}$$

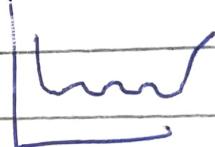
Logistic regression's cost fn -

if mse is considered

$$\text{i.e. } L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

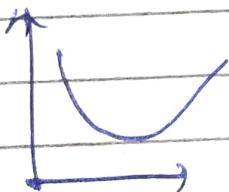
gives a curve with local minima  
which doesn't work well with  
gradient descent

$$\hat{y} = \sigma(w^T x + b)$$



Instead we use as Loss fn [Applied on single training ex]

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$



$y \in \{0, 1\}$   
 $\hat{y} \in (0, 1)$

$$\text{Cost fn} = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

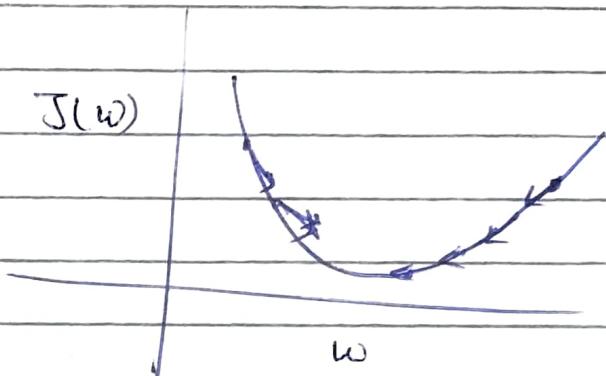
$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

## \* Gradient Descent

repeat of

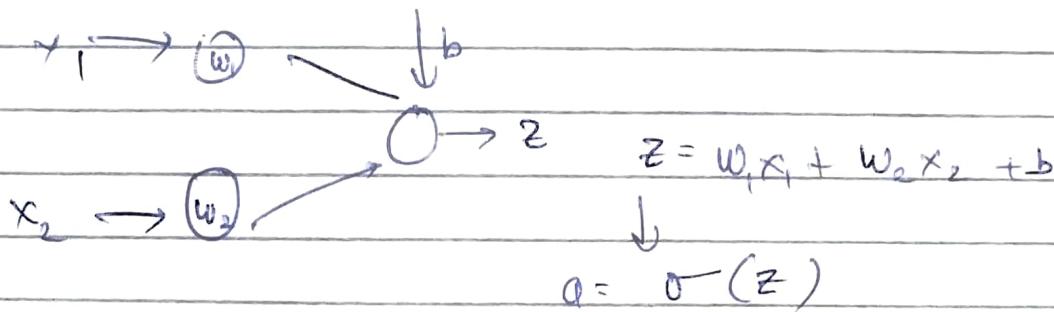
$$w := w - \alpha \frac{d J(w)}{dw}$$

}



$$J(w, b) \leftarrow : w := w - \alpha \frac{d J(w, b)}{dw}$$

$$b := b - \alpha \frac{d J(w, b)}{db}$$



$$\frac{d J(a, y)}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$$

$\nearrow J(a, y)$

$$\frac{dL}{dz} = \frac{dL}{da} \cdot \left( \frac{da}{dz} \right) \rightarrow a(1-a)$$

$$\therefore = a - y$$

$$\frac{\partial L}{\partial w_1} = dw_1 = x_1 dz$$

$$dw_2 = x_2 dz$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

$$db = dz$$

\* LR on

m examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m J(a^{(i)}, y^{(i)})$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} J(a^{(i)}, y^{(i)})$$

Initialize

$$J = 0 \quad dw_1 = 0 \quad dw_2 = 0 \quad db = 0$$

## 2 Vectorization

$$\cdot z = w^T x + b$$

avoid explicit for loops whenever possible

$$z = \text{np.dot}(w.T, x) + b$$

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

$$dz = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]$$

$$A = [a^{(1)} \ \dots \ a^{(n)}] \quad y = [y^{(1)} \ \dots \ y^{(m)}]$$

$$dz = A - y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]$$

Implementation!  
for i in range(epoch):  
 z = w^T x + b

$$= \text{np.dot}(w.T, x) + b$$

$$A = \sigma(z)$$

$$dz = A - y$$

$$dw = \frac{1}{m} x dz^T$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

$$\hat{y} : \text{If } y=1 \quad p(y|x) = \hat{y}$$

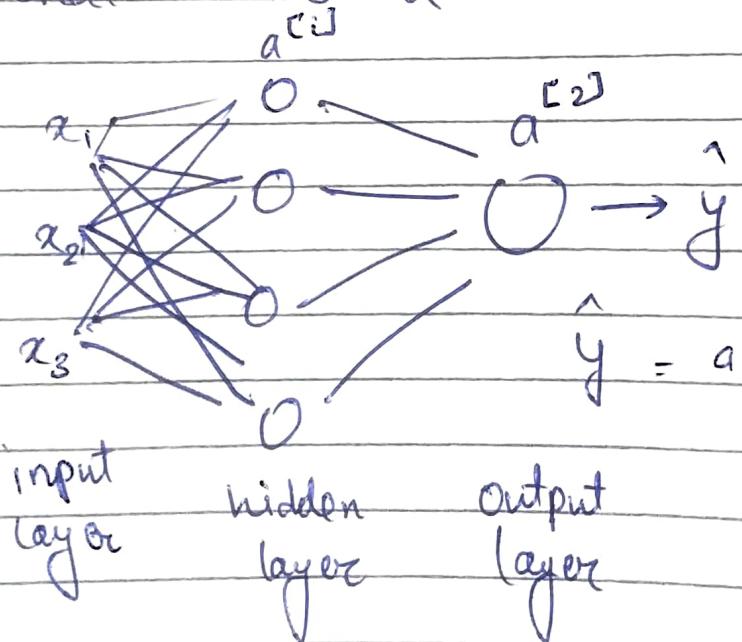
$$y=0 \quad p(y|x) = 1-\hat{y}$$

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{If } y=1: p(y|x) = \hat{y} \quad (1-\hat{y})^0$$

$$y=0: p(y|x) = \hat{y}$$

## # Shallow Neural Network



What a single node does -

$$z_1 = w_1^T x + b_1, \quad a_1 = \sigma(z_1)$$

Similarly

$$z_2 = w_2^T x + b_2, \quad a_2 = \sigma(z_2)$$

$$z_3 = w_3^T x + b_3, \quad a_3 = \sigma(z_3)$$

$$z_4 = w_4^T x + b_4, \quad a_4 = \sigma(z_4)$$

Vectorize this.

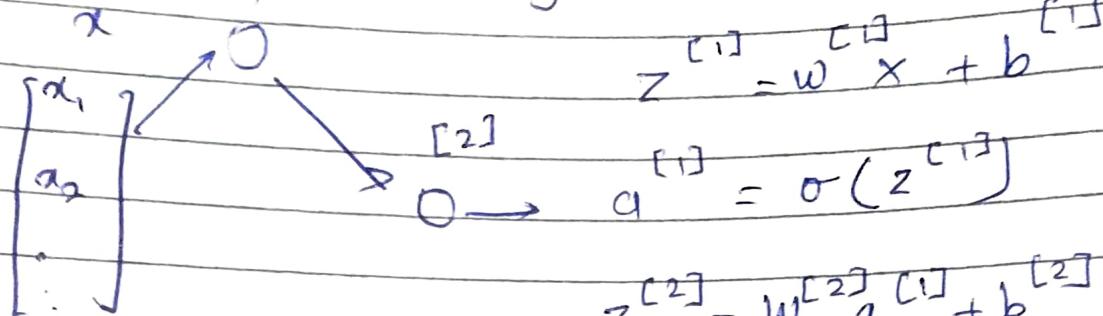
$$\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix}$$

Shape: [4, 3]      [3, 1]

$\sigma([ \cdot ])$

a

for  $i=1$  to  $m$ :



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

$\sigma$  is ~~not~~ not necessarily the sigmoid fn it can be any activation fn.

$$x = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ | & | & & | \end{bmatrix}_{n_x \times m}$$

$$z^{[1]} = \begin{bmatrix} | & | & & | \\ z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & & | \end{bmatrix}$$

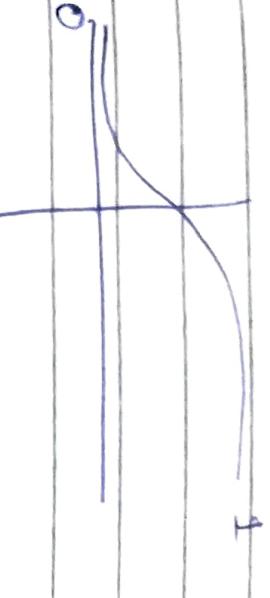
$$A^{[1]} = \begin{bmatrix} | & | & & | \\ a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

← training ex →

hidden units

## Activation funcs:

$$\text{sigmoid } a = \frac{1}{1+e^{-z}}$$



used in binary classification

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} - z$$

is much superior than sigmoid

those often pose problems for G.D.

- for large positive or negative values, as the slope is almost close to 0.  $\frac{dy}{dx} \rightarrow 0$

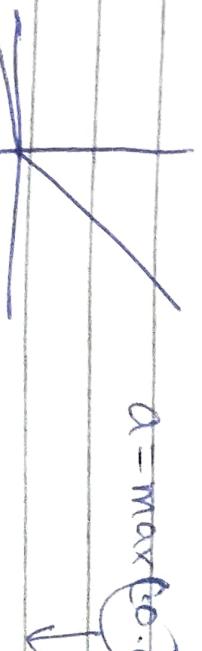
ReLU

$$a = \max(0, z)$$

Rectified Linear Unit

leaky ReLU

$$a = \max(0.01z, z)$$



can be any small number

Why use Non linear Activation fns?

Given  $x$ :

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

$$a^{[1]} = g^{[1]}(z^{[1]}) \rightarrow \text{make it } z^{[1]}$$

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

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

Let  $g(x) = x$

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

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

$$(w^{[2]} w^{[1]}) x + (w^{[2]} b^{[1]} + b^{[2]})$$

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

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

back to the ~~previous~~

layer form of eqn

Hence the hidden layer is useless

→ might just have these  $w^{[1]} b^{[1]}$  in the input layer

## Derivatives

$$g(z) = \frac{1}{1+e^{-z}} \quad g'(z) = g(z)(1-g(z))$$

$$g(z) = \tanh(z) \quad g'(z) = 1 - (\tanh(z))^2$$

# GD for NN

$$\text{Cost fn: } J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) \\ = \frac{1}{m} \sum_{i=1}^m f(\hat{y}_i, y_i)$$

$$\underbrace{\phantom{\sum_{i=1}^m f(\hat{y}_i, y_i)}}_a$$

repeat

forwards for computing derivatives

forward propagation:

$$z^{[1]} = w^{[1]} x + b^{[1]} \quad dz^{[2]} = A^{[2]} - Y$$

$$A^{[1]} = g^{[1]}(z^{[1]}) \quad \frac{d w^{[2]}}{d w^{[1]}} = \frac{1}{m} d z^{[2]} A^{[1]}$$

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

$$db^{[2]} = \frac{1}{m} \text{np.sum}(d z^{[2]}, \text{axis}=1)$$

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

Back Propagation

$$dz^{[1]} = w^{[1] \top} d z^{[2]} * g^{[1]}(z^{[1]})$$

$$d w^{[1]} = \frac{1}{m} d z^{[1]} x^\top$$

$$d b^{[1]} = \frac{1}{m} \text{np.sum}(d z^{[1]}, \text{axis}=1)$$

$$d w^{[2]} = \frac{1}{m} d z^{[2]} x^\top$$

duo [1] instead of dw[2]

8

Initialize  $w$  to small positive values

$w[1] = np.random([2, 2]) * 0.01$

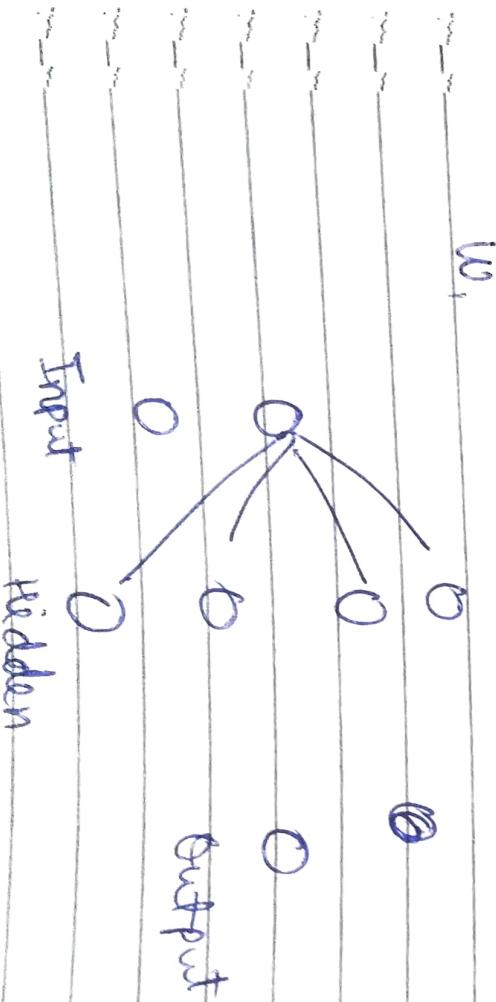
$b_1 = np.zeros([2, 1])$

If weight values are large  
you may end up on the  
flat sections of tanh & sigmoid  
fun. Hence gradient descent  
will be very slow.

## Programming Assignment

$x = 2 \times 400$

$y = 1 \times 400$



$$w_1 = (4, 2 \times 2) \quad b_1 = [1, 0, 0, 0]^T \quad 4 \times 1$$

$$w_2 = 1 \times 4$$

# # Deep Layer NN

## Notations

$$\alpha_1 \quad 0 \quad 0$$

$$\alpha_2 \quad 0 \quad 0 \quad 0$$

$$\alpha_3 \quad 0 \quad 0 \quad 0 \quad 0$$

$$V \quad 0 \quad 0 \quad 5 \quad 5 \quad 3 \quad 1$$

$$L = 4 \quad (\text{has } 4 \text{ layers})$$

$$n^{[l]} = 0 \# \text{ units in layer } l$$

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

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

$w^{[l]}, b^{[l]}$  = weights & bias for layer  $l$

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

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

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

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

b: 35

$$a^{[L]} = g^{[L]}(z^{[L]}) \quad \text{dim}(a) = (n^{[L]}, 1)$$

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

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

generalized to for  $L=1$  to  $L$

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

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

$$\dim(a^{[L]}) = (n^{[L]}, n^{[L-1]})$$

Why deep?

Ex speech recognition

Wave  $\rightarrow$  low level  $\rightarrow$  Phonetic sounds  $\rightarrow$  words from phonetic features

Pitch  
etc  
Sounds

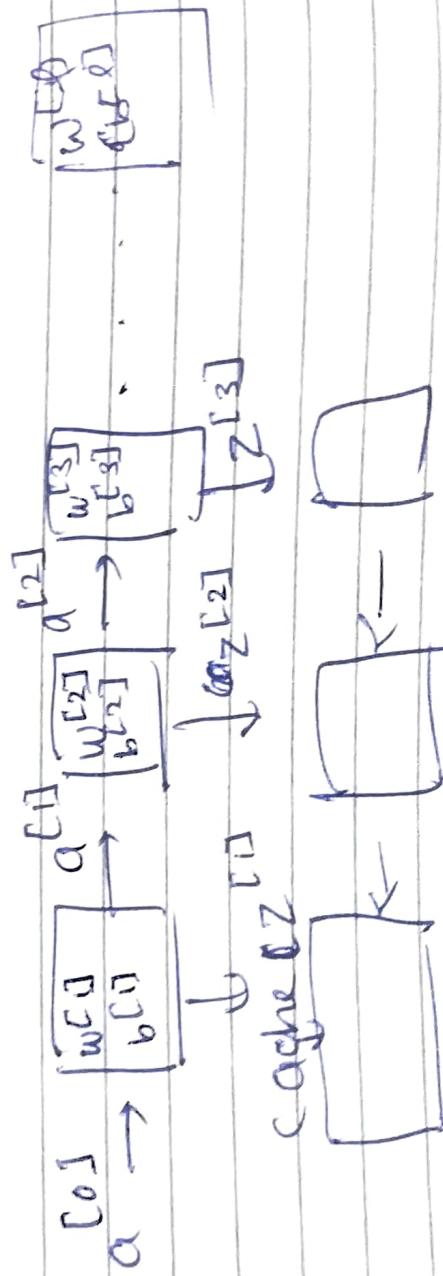
Analogous to many layer logic gate circuits.

## Forward & Backward pass

$$\begin{array}{c}
 \begin{array}{ccccc}
 x_1 & 0 & 0 & 0 & 0 \\
 x_2 & 0 & 0 & 0 & 0 \\
 x_3 & 0 & 0 & 0 & 0 \\
 x_4 & 0 & 0 & 0 & 0
 \end{array} \\
 \downarrow \\
 \text{Consider layer } l
 \end{array}$$

Forward : input:  $a^{[l]}$   
 output:  $a^{[l]}, \text{cache}(z^{[l]}, w^{[l]}, b^{[l]})$

Backward : input:  $da^{[l-1]}$   
 output:  $da^{[l]}$ ,  
 $dW^{[l]}$ ,  
 $db^{[l]}$



$$dw^{[l]} = dz^{[l]} * \alpha^{[l-1]T}$$

(2:30)

$$B.P. \quad d_{z^{[l]}}^{[l]} = d_{\alpha}^{[l]} * f^{[l]}(z^{[l]})$$

$$d_{w^{[l]}}^{[l]} = d_{z^{[l]}}^{[l]} * \underbrace{f^{[l-1]}(x)}_{a^{[l-1]}}$$

$$d_{b^{[l]}}^{[l]} = \frac{1}{m} (d_{z^{[l]}}^{[l]}) \rightarrow np.sum(), axis=1$$

$$d_{\alpha^{[l-1]}}^{[l]} = dw^{[l]} \cdot d_{z^{[l]}}^{[l]}$$

## # Parameters & Hyperparameters

$\downarrow$   $\alpha$ , iterations, hidden layer,  
 $w$  &  $b$   
hidden units, choice of  
activation func.