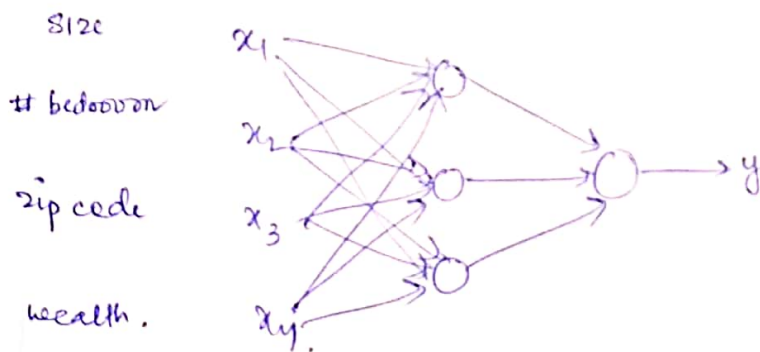
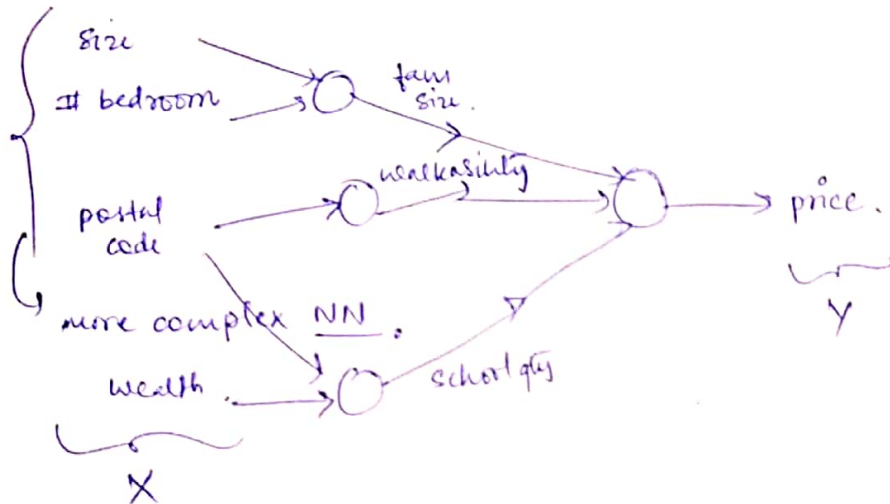
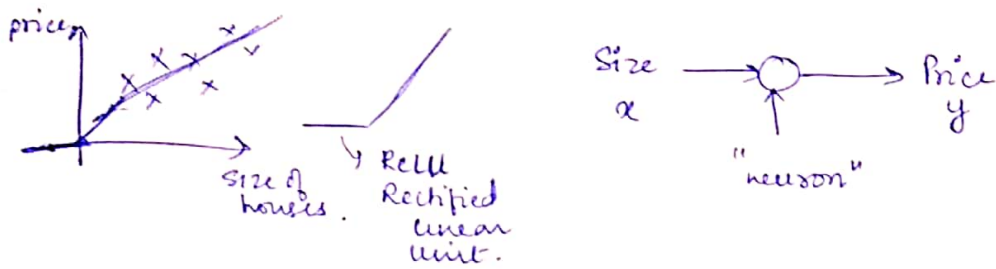


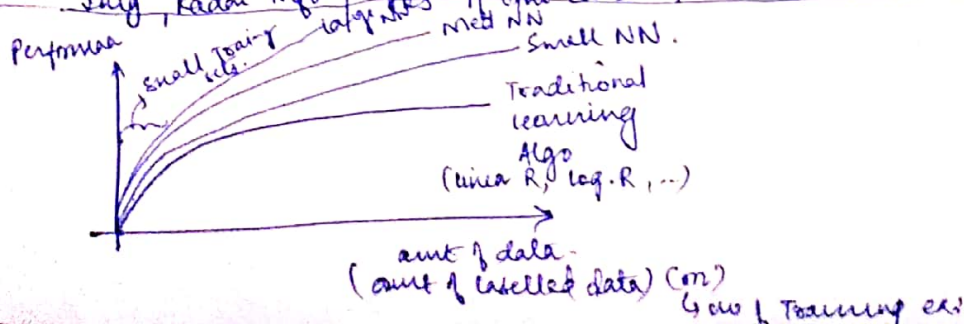
# I. Neural N/w & Deep Learning:

\* What is a NN?

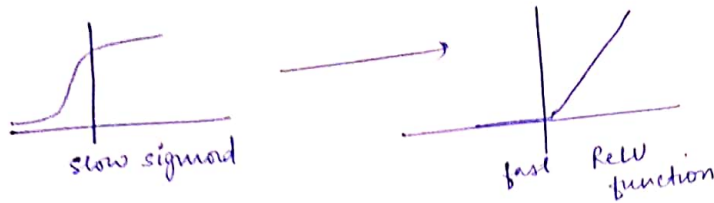


\* Supervised learning with Neural N/w:

Ex	Input ( $x$ )	o/p ( $y$ )	Applic <sup>n</sup>
	Ad, user info	click on ad? (0/1)	Online Advert $\rightarrow$ standard NN
	Image	Obj (1, ..., 1000)	Photo tagging $\rightarrow$ CNN
	Audio	Text transcript	speech recog $\rightarrow$ RNN
	$\vdots$	$\vdots$	$\vdots$
	Sing, Radar info	offr cars	Auto drive $\rightarrow$ Hybrid/complex:



DL:



⊕ logistic Regression as Neural N/w :

\* Binary classification :



unroll pixel values into a feature vector.

$$x = \begin{bmatrix} 255 \\ 255 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$$

if img is  $64 \times 64$ , 3 vectors.  
 $\therefore n_x = 64 \times 64 \times 3 = 12288$   
 $x \rightarrow y$ .

Notation :

$(X, y)$  : single T.S.  $\in \mathbb{R}^{n_x}$ ,  $y \in \{0, 1\}$

m training example.

$m = m_{\text{train}}$

$m_{\text{test}} = \# \text{ test examples}$ .

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

$\xleftarrow{m}$

$\uparrow n_x$   
 $\downarrow$

$$X \in \mathbb{R}^{n_x \times m}$$

$$Y = [y^{(1)}, y^{(2)} \dots y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

\* Logistic Regression

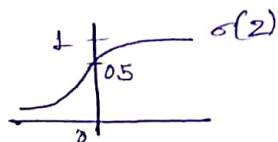
$X$ , want  $\hat{y} = P(y=1|x)$

$$X \in \mathbb{R}^{n_x}$$

params  $w \in \mathbb{R}^{n_x}$

$$b \in \mathbb{R}$$

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



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

if  $z$  large,  $\sigma(z) \approx 1$   
 $z$  small  $< 0$ ,  $\sigma(z) \approx 0$

$$x_0 = 1 \quad x \in \mathbb{R}^{n \times 1}$$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \begin{cases} b. \\ w. \end{cases}$$


→ logistic Regression cost func<sup>n</sup>:

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

given  $\{(x^{(0)}, y^{(0)}), \dots, (x^{(m)}, y^{(m)})\}$   $\hat{y}^{(i)} \approx y^{(i)}$

$$z^{(i)} = w^T x^{(i)} + b$$

$$\begin{Bmatrix} x^{(i)} \\ y^{(i)} \\ z^{(i)} \end{Bmatrix} \rightarrow i^{th} \text{ example}$$

loss (error func<sup>n</sup>):  $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$  

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

$y=1$  ;  $L(\hat{y}, y) = -\log \hat{y} \leftarrow$  want  $\log \hat{y}$  large, want  $\hat{y}$  large.

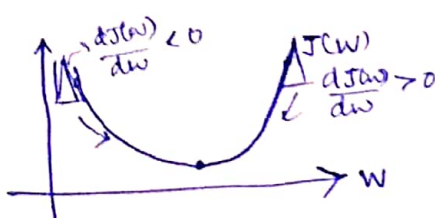
$y=0$  ;  $L(\hat{y}, y) = -\log (1-\hat{y}) \leftarrow$  want  $\log (1-\hat{y})$  large, want  $\hat{y}$  small.

Cost function (J) for entire T.S.

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

\* Gradient descent :

want to find  $w$  and  $b$  that minimise  $J(w, b)$



Repeat  $\{$

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

$\}$

learning rate

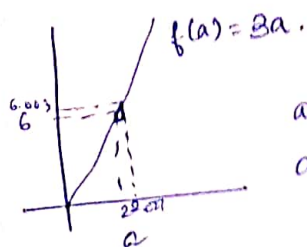
$$\Rightarrow w := w - \alpha dw$$

"dw"

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

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

## \* Derivatives :



$$a = 2 \quad f(a) = 6$$

$$a = 2.001 \quad f(a) = 6.003$$

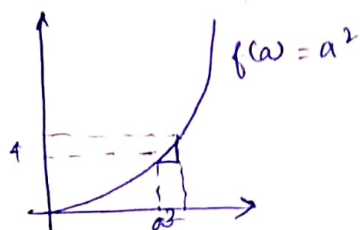
Slope (derivative) of  $f(a)$  at  $a = 2$  is '3'.

↓  
h/w.

$$\frac{0.003}{0.001} = 3$$

$$\frac{df(a)}{da} = 3 = \frac{d}{da} f(a).$$

for derivatives, nudging by infinitesimally small amount



$$a = 2 \quad f(a) = 4$$

$$a = 2.001 \quad f(a) \approx 4.004$$

$$\left. \frac{df(a)}{da} \right|_{a=2} = 4.$$

## \* Computation graph

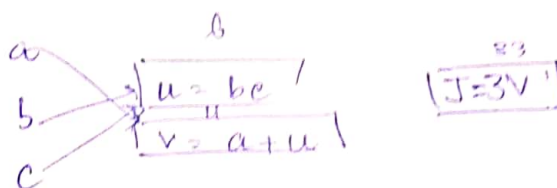
$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

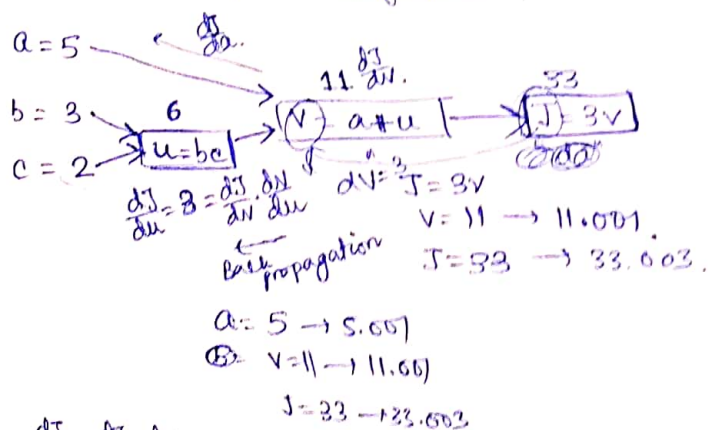
$$J = 3v$$

$$= 3(5 + 3 \times 2) = 33$$



one step of ( $\leftarrow$ ) prop. on computation graph yields derivative of final op variable.

## \* Derivatives with computation graph :-



$$f(a) = 3a$$

$$\frac{df(a)}{da} = \frac{dJ}{da} = 3$$

$$J = 3v$$

$$\frac{dJ}{dv} = 3$$

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} \rightarrow \text{chain rule.}$$

$$\frac{d(\text{final output variable})}{d(\text{var})} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

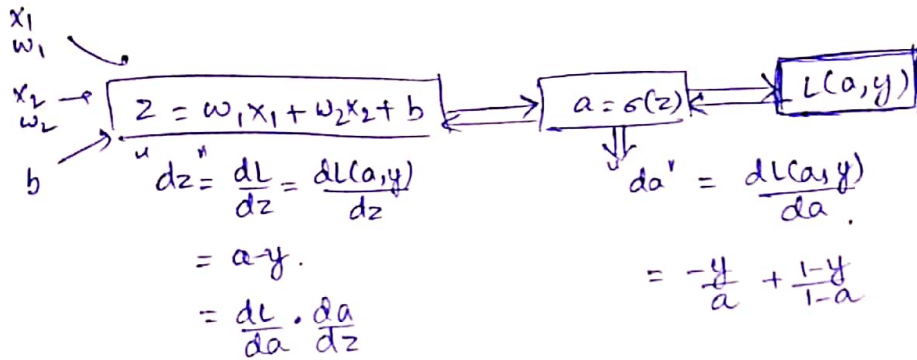


# \* logistic regression Grad. Descent.

$$z = w^T x + b.$$

$$\hat{y} = a = \sigma(z)$$

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



$$\frac{dL}{dw_1} = "dw_1" = x_1 dz$$

$$"dw_2" = x_2 dz$$

$$"db" = "dz"$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db.$$

\* for 'm' training examples:

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

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

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_1}}_{dw_1^{(i)} = (x_1^{(i)} - y^{(i)})}$$

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

for  $i=1$  to  $m$

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J^{(i)} = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

~~for~~

$$J /= m$$

$$dw_1 /= m$$

$$dw_2 /= m$$

$$db /= m$$

soln from for loops less eff.

→ vectorisation is soln.

# \* Vectorisation :

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

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w \in \mathbb{R}^{n_x}$$

$$x \in \mathbb{R}^{n_x}$$

vectorised impl<sup>n</sup>:

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

{ GPU } SIMD - single inst. multiple data.  
 { CPU } helps in parallelism.

- avoid for loops wherever possible.

$$u = Av.$$

$$u_i = \sum_j A_{ij} v_j$$

$$u = \text{np.zeros}(n, 1)$$

for i = ---

for j = ---

} long time.

vector version

$$u = \text{np.dot}(A, v)$$

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

non vectorised.

$$u = \text{np.zeros}(n, 1)$$

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

vectorised.

import numpy as np.

$$u = \text{np.exp}(v)$$

$$\text{np.log}(v)$$

$$\text{np.abs}(v)$$

$$\text{np.max}(v, 0)$$

⋮

vectorised impl of ~~logistic~~ logistic regression.

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \quad \begin{matrix} (n_x, m) \\ \mathbb{R}^{n_x \times m} \end{matrix}$$

$$Z = [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = w^T X + [b \ b \ b \ \dots \ b]_{1 \times m}$$

$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & w^T x^{(3)} + b & \dots \end{bmatrix}$$

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

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(Z)$$

$$y(1, 1) \in \mathbb{R}$$

"Broadcasting"

vectorizing logistic Regr<sup>n</sup> gradient of p:

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

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

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

$$dZ = A - Y$$

$$dw = 0$$

$$dw += x^{(1)} dz^{(1)}$$

$$dw += x^{(2)} dz^{(2)}$$

$$\vdots$$

$$dw/m$$

$$db = 0$$

$$db += dz^{(1)}$$

$$db += dz^{(2)}$$

$$\vdots$$

$$db += dz^{(m)}$$

$$db/m$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

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

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

$$= \frac{1}{m} \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]_{n \times 1}$$

revised: logistic regression.

$$Z = W^T X + b$$

$$= \text{np.dot}(W, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

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

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

\* Broadcasting in Python.

Ex:

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fats	1.8	135.0	99.0	8.9

$$= A_{(3,4)}$$

% calories: without explicit for loop.

$$\text{cal} = A \cdot \text{sum}(\text{axis}=0) \rightarrow \text{vertical sum}$$

$$\text{percentage} = \frac{100 * A / \text{cal} \cdot \text{reshape}(1, 4)}{\text{python broadcasting}}$$

divide by 1x4 matrix.

$$(3,4)/(1,4) \quad \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$







$$\underbrace{\begin{bmatrix} w_1^{[1]} \\ w_2^{[1]} \\ w_3^{[1]} \\ w_4^{[1]} \end{bmatrix}}_{W^{[1]} \quad (4,1)} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

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

given i/p's.

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

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

next layer.

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

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

$$X \longrightarrow a^{[2]} = \hat{y}$$

$$x^{(1)} \longrightarrow \hat{y}^{(1)} = a^{[2]}(1)$$

$$x^{(2)} \longrightarrow \hat{y}^{(2)} = a^{[2]}(2)$$

$$\vdots$$

$$x^{(m)} \longrightarrow \hat{y}^{(m)} = a^{[2]}(m)$$

$a^{[2]}(i)$   
↓  
layer 2

example i

vectorize:

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad (n \times m)$$

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

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

for  $i=1$  to  $m$

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

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

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

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

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

$$A^{[1]} = \sigma(z^{[1]})$$

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

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

training ex

↓  
Hidden units.

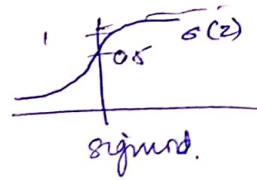
## \* Activation func<sup>n</sup>:

$$\sigma(z^{[1]}) \leftrightarrow g(z^{[1]})$$

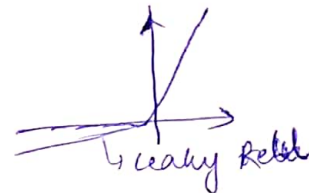
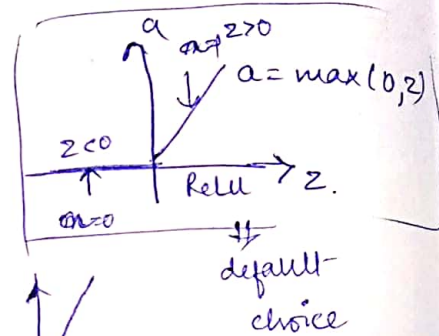
$$\sigma(z^{[2]}) \leftrightarrow g(z^{[2]})$$

activation func<sup>n</sup> can be diff for diff layers.

$y \in \{0, 1\}$   
use sigmoid



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



## \* Why nonlinear Activation func<sup>n</sup>?

- if not, then

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

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

$\therefore$  linear implementation only.

$\therefore$  might not have any hidden layers.

## \* Derivatives of Activation func<sup>n</sup>.

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

$$\frac{dg(z)}{dz} = \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right)$$

$$= g(z)(1 - g(z))$$

$$z \approx 10$$

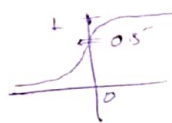
$$g(z) \approx 1$$

$$g'(z) = \frac{d}{dz} g(z) \approx 0$$

$$z \approx -10$$

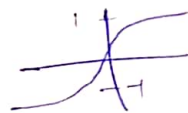
$$g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0$$



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

$$\frac{dg(z)}{dz} = 1 - (\tanh(z))^2$$



$\rightarrow$  ReLU & leaky ReLU ..

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

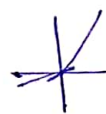
$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$



technically undefined at zero.

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

$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$$



\* gradient descent for neural N/w:

Params:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$(n^{[1]}, n^{[0]}) \rightarrow (n^{[1]}, 1)$   
 $(n^{[2]}, n^{[1]}) \rightarrow (n^{[2]}, 1)$

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

Cost func:  $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^n L(\hat{y}^i, y^i)$

$\uparrow$   
 $a^{[2]}$

Grad. descent

Repeat {

compute pred  $(\hat{y}^i, \dots, i=1, \dots, m)$

$$dw^{[1]} = \frac{dJ}{dw^{[1]}}$$

$$db^{[1]} = \frac{dJ}{db^{[1]}}$$

$$w^{[1]} := w^{[1]} - \alpha dw^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha db^{[1]}$$

Formulas for computing derivatives.

forward pass:

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

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

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

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

Back propag<sup>n</sup>:

$$dz^{[2]} = a^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} a^{[1]T}$$

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

$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

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

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

\* Random initialisations

initialising  $w$  to all zeros is problematic

$$a_1^{[1]} = a_2^{[1]}$$

$$dz_1^{[1]} = dz_2^{[1]}$$

$\therefore$  2 hidden units can be completely identical.

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

after iterations, all rows equal in  $w^{[1]}$

symmetric

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$\text{sol}^n$ :  $w^{[1]} = \text{np.random.randn}((2, 2)) * 0.01$   
 $b^{[1]} = \text{np.zeros}(2, 1)$   
 $w^{[2]} = \dots$   
 $b^{[2]} = 0$   
 use small initial values.

## ⊕ Deep neural Networks :

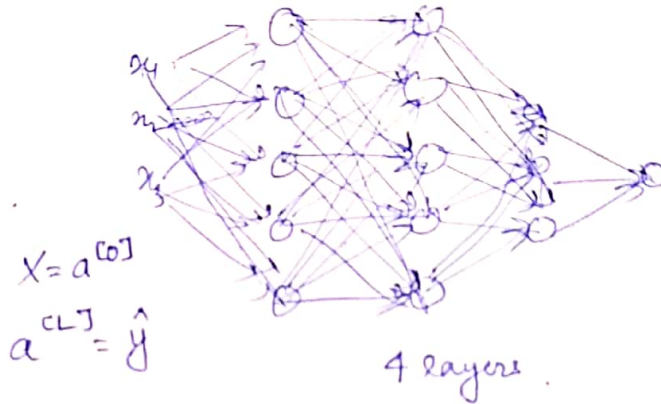
$L \rightarrow$  no of layers

$n^{[l]} \rightarrow$  no of units in 'l'.

$$n^{[1]} = 5$$

$$n^{[2]} = 5$$

$$n^{[3]} = 3$$



\* Forward propagation in a deep neural N/w:-

$$X \text{ or } z^{[1]} = w^{[1]}x + b^{[1]}$$

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

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

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

...

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

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

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

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

vectorized

$$Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}$$

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

⋮

for deep NN, a for loop has to be there  
loop goes from 1 to L

\* getting matrix dimensions right.

$$L = 5$$

$$n^{[0]} = 3$$

$$n^{[1]} = 5$$

$$n^{[2]} = 5$$

$$n^{[3]} = 2$$

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

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

$$(3, 1) \quad (2, 1)$$

$$(n^{[0]}, 1) \quad (n^{[1]}, 1)$$

$$[:] = [::] [:]$$

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

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

vectorized

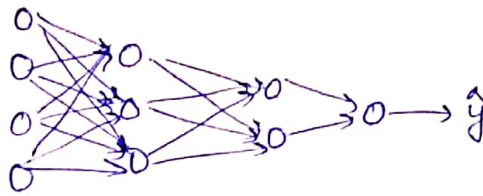
$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$(5, 1) \quad (5, 5) \quad (5, 1) \quad (5, 1)$

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

$$(n^{[0]}, m) \quad (n^{[1]}, n^{[0]}) \quad (n^{[1]}, m) \quad (n^{[1]}, 1) \rightarrow \text{broadcasting here}$$

\* Distribution and deep representation.



- face recognition

- Audio  $\rightarrow$  low level  $\rightarrow$  Phonemes  $\rightarrow$  words  $\rightarrow$  Sentence.  
audio

- Circuit theory & deep learning;  
function to compute with small & large exponential hidden layers.

deep NN that shallower NN need more

\* Building blocks of deep NN.

Forward  
Layer  $l$ :  $w^{(l)}$ ,  $b^{(l)}$

i/p:  $a^{(l-1)}$

o/p:  $a^{(l)}$

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

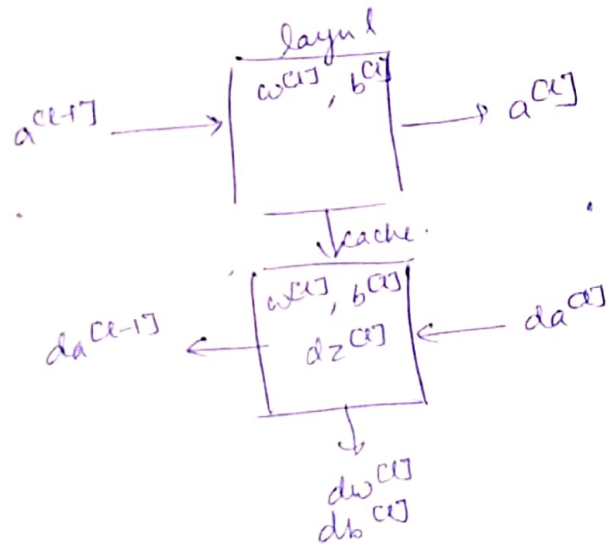
$$a^{(l)} = g^{(l)}(z^{(l)})$$

cache  $z^{(l)}$

Backward

o/p:  $da^{(l-1)}$

Cache  $[z^{(l)}]$   
 $dw^{(l)}$   
 $db^{(l)}$



\* Prop & Back Propagation.

Forward  
i/p:  $a^{(l-1)}$   
o/p:  $a^{(l)}$ , cache( $z^{(l)}$ )

vectorized:

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

$$X = A^{(0)} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \dots$$

Back

i/p:  $da^{(l)}$

o/p:  $da^{(l-1)}$ ,  $dw^{(l)}$ ,  $db^{(l)}$

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

$$dw^{(l)} = dz^{(l)} \cdot a^{(l-1)}$$

$$db^{(l)} = dz^{(l)}$$

$$da^{(l+1)} = w^{(l+1)T} dz^{(l)}$$

$$dz^{(l)} = w^{(l+1)T} dz^{(l+1)} * g^{(l)}(z^{(l)})$$

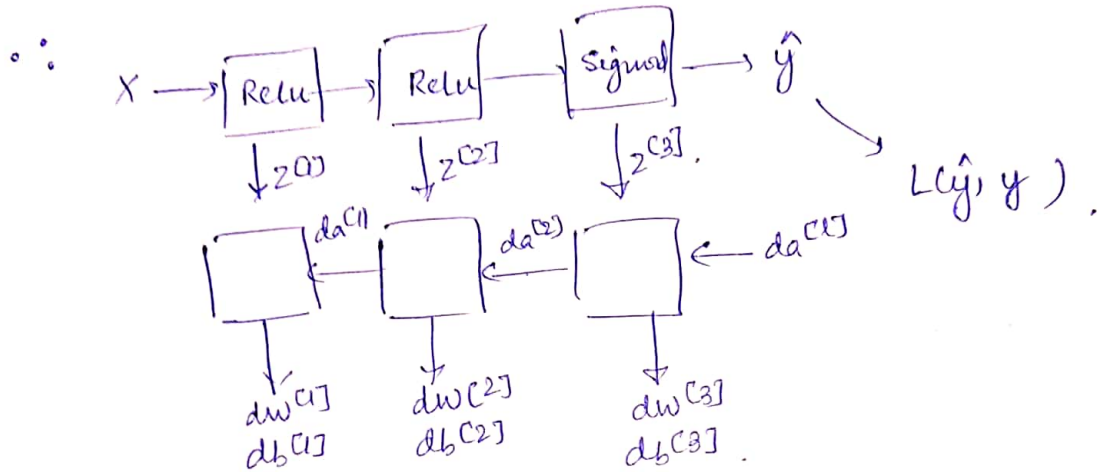
vectorized

$$dz^{[L]} = dA^{[L]} * g^{[L]'}(z^{[L]})$$

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

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

$$dA^{[L-1]} = W^{[L]T} \cdot dz^{[L]}$$



\* Parameters vs Hyperparameters .

#  
 $w, b$  .

#  
learning rate  $\alpha$

# iterations

# hidden layer  $L$

# hidden unit,  $n^{[1]}$ ,  $n^{[2]}$ , ...

choice of activation func<sup>n</sup>

\* Deep learning & human brain :