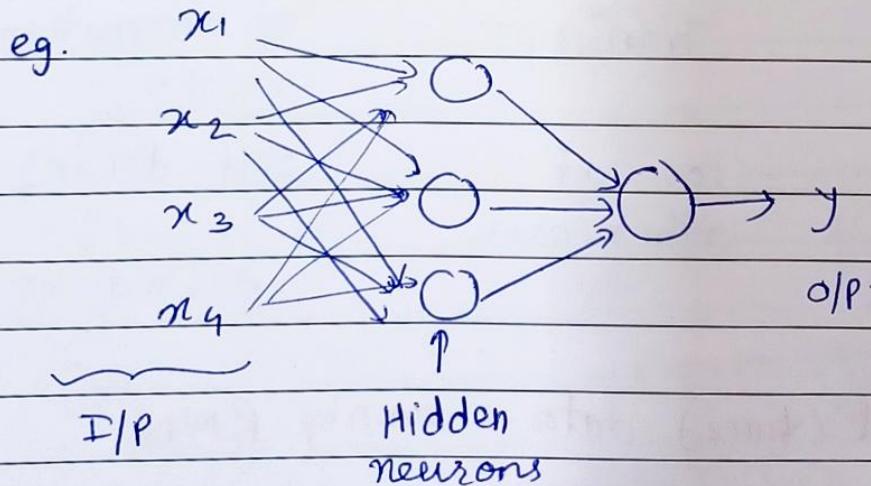
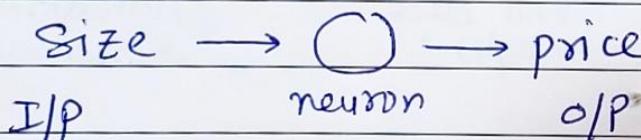
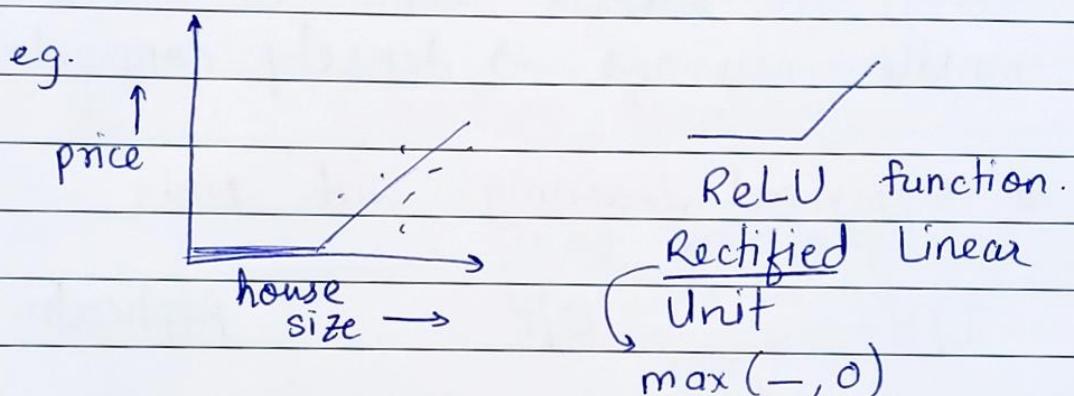


Date : 19/7/24

MON TUE WED THU FRI SAT SUN

# NEURAL NETWORKS & DEEP LEARNING (C-1)

## # Neural Network



→ (Koala example)

Date : .....

MON TUE WED THU FRI SAT SUN  
□ □ □ □ □ □ □

N.N. has ability to decide by itself what I/P features each of the hidden neuron wants/requires to do its computation.

Every I/P feature conn. to every middle neurons  $\Rightarrow$  densely connected.

## # Supervised Learning with NN.

I/P	O/P	Applicatn.
Ad, User Info.	Will click (O/I)?	Advertising $\hookrightarrow$ Standard NN
Image	object	Tagging $\hookrightarrow$ CNN
Audio	Transcript	Sp.-recognition $\hookrightarrow$ RNN
Img, Radar	Pos. of other cars	Self- driving $\hookrightarrow$ Custom/ Hybrid NN

Temporal (time) data  $\rightarrow$  mainly RNN.

↑ Data & ↑ Scale (i.e. more layers of NN)



↓  
↑ Performance

Date: .....

MON TUE WED THU FRI SAT SUN

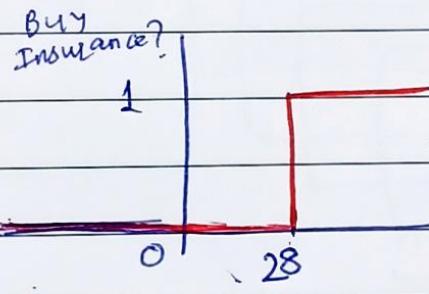
## Structured Data

⇒ eg.	Size	#bedrooms	Price
			{

## Unstructured Data

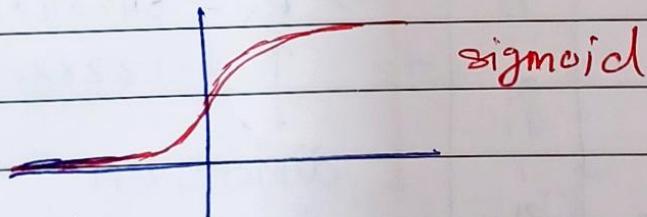
⇒ Audio, Image, Text

## # Basics - Activation functions.

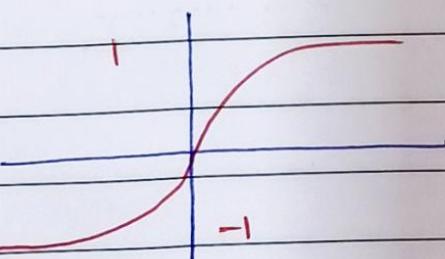


→ eg. People above 28y.  
will buy & below  
won't.  
But X.

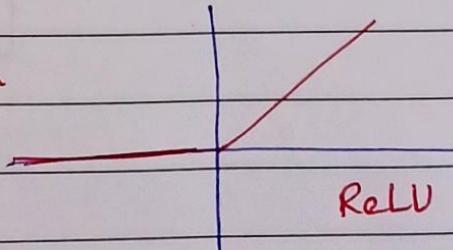
STEP Functrn



sigmoid



tanh



ReLU

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

Date : .....

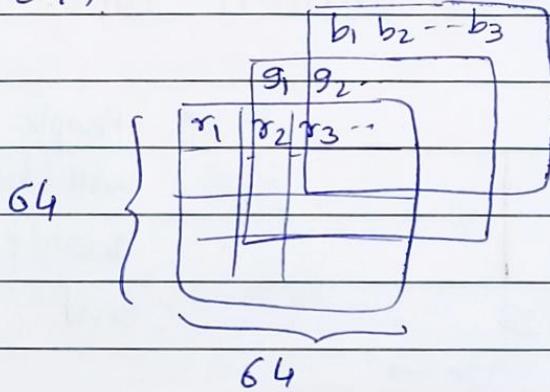
MON TUE WED THU FRI SAT SUN

Hidden Layer mostly - ReLU or Leaky ReLU.  
 O/P  $\rightarrow$  sigmoid mostly.

## # Binary classification

O/P  $\rightarrow$  0 or 1

e.g. cat image  $\rightarrow$  cat OR non-cat  
 $(64 \times 64)$



Feature vector

$$X = \begin{bmatrix} r_1 \\ g_1 \\ b_1 \\ r_2 \\ g_2 \\ b_2 \\ \vdots \\ r_{64} \\ g_{64} \\ b_{64} \end{bmatrix} \Rightarrow n_x = 64 \times 64 \times 3$$

$\uparrow = 12288$   
dimension.

Date : .....

MON TUE WED THU FRI SAT SUN

$(x, y)$  = Training eg.

$$x \in \mathbb{R}^{n_x}, y = \{0, 1\}$$

$m \rightarrow$  # training examples :

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

$$X = \left[ \begin{array}{c|c|c|c|c} & & & & \\ \hline & | & | & | & | \\ \hline x^1 & x^2 & \dots & x^m & \\ \hline & | & | & | & | \\ \hline \end{array} \right] \Bigg\} n_x$$

$\underbrace{\hspace{10em}}$   
 $m$

compact notation to store  $m$   
training examples.

For  $Y$  also, stack in cols -

$$Y = [y^1, y^2, \dots, y^m]$$

## # Logistic Regression

→ Supervised ML Algo. used for binary class. tasks

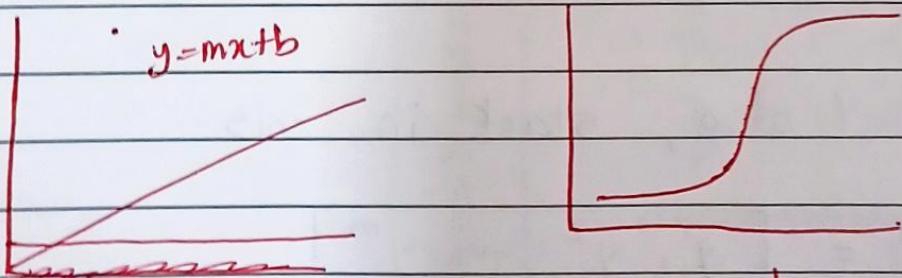
Given I/P  $x$ ,  $\hat{y} = P(y=1|x)$   
probability

Params:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

O/I/P:  $\hat{y} = w^T x + b$  ↗  
is not betw [0, 1]

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

$\underbrace{z \mapsto}_{\text{limits 0 to 1}}$



$$y = \frac{1}{1 + e^{-(mx+b)}}$$

no. of weights ( $w$ ) = no. of I/P features -

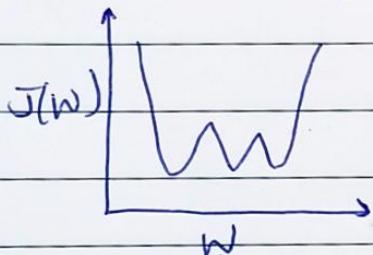
## # Loss Function:

→ Measures how far an estimated value is from its true value.

→ For general Linear Regr:

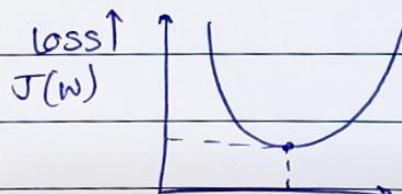
$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

→ But for Logistic Regr. ↗ don't work.



Gradient Descent

Curve with multiple local minima



optimal Grad. Desc.

Loss fun applied to single training eg.  
cost fn: Cost of parameters.

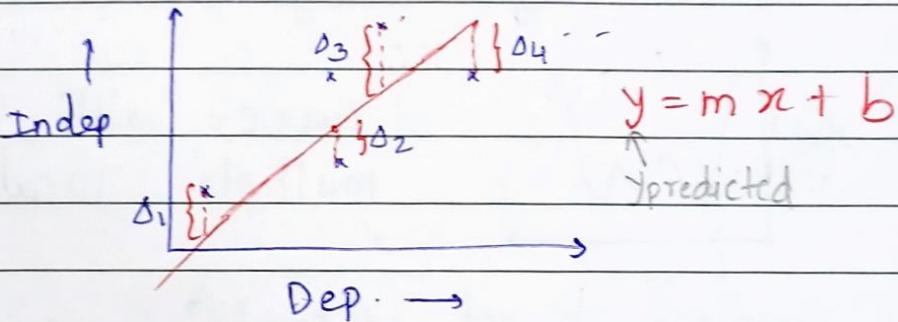
$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

$$\rightarrow J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

## # Linear Regression

e.g. House area → Price.

Useful when the relationship betn dependent & independent var can be assumed to be linear.



Lin. Reg. with multiple variables.

$$\text{price} = m_1 \times \text{area} + m_2 \times \text{bedrooms} + m_3 \times \text{age} + b$$

Mean Squared Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \text{predicted})^2$$

$n$  = no. of data pts avail

Also called Cost Function

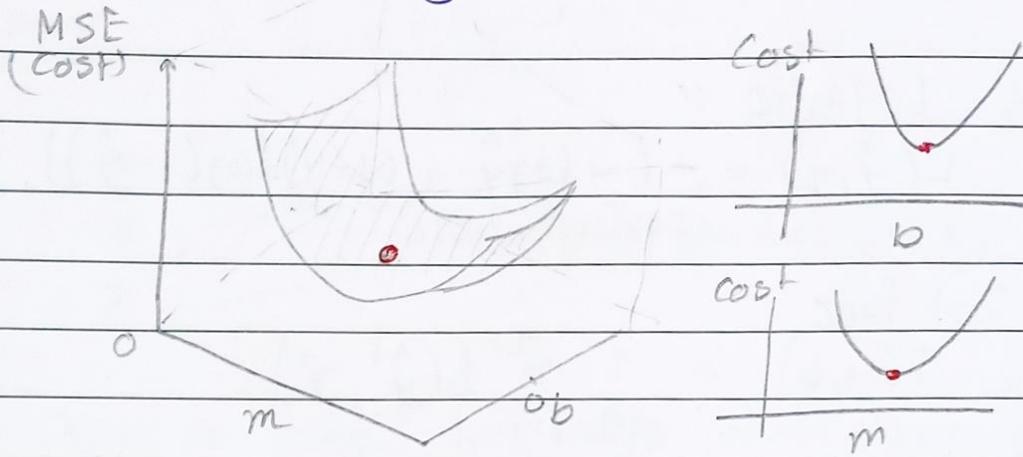
Date : .....

MON TUE WED THU FRI SAT SUN

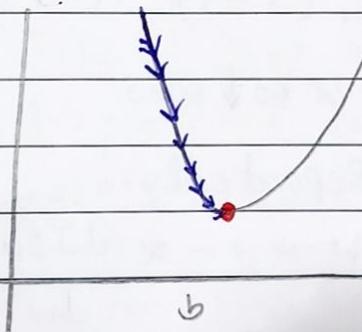


## # Gradient Descent

↳ algo to find best fit line for given training dataset



cost



We will take  
baby steps to  
move towards  
minima & slowly  
decrease the  
length of babystep.

Achieved using derivative

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

$$\frac{\partial}{\partial m} = \frac{2}{n} \sum_{i=1}^n -x_i (y_i - (mx_i + b))$$

$$\frac{\partial}{\partial b} = \frac{2}{n} \sum_{i=1}^n - (y_i - (mx_i + b))$$

Date : .....

MON TUE WED THU FRI SAT SUN  
      

$$m = m - \text{Learning rate} \cdot \frac{\partial}{\partial m}$$

$$b = b - \text{Learning rate} \cdot \frac{\partial}{\partial b}$$

In Logistic :

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

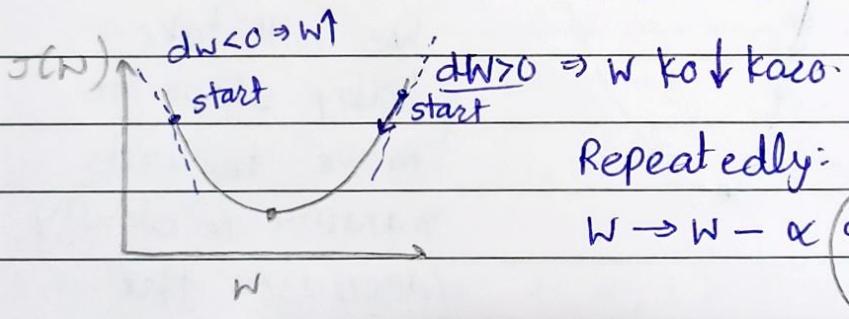
↳ for each individual data

Cost func

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

for entire training

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



Repeatedly:

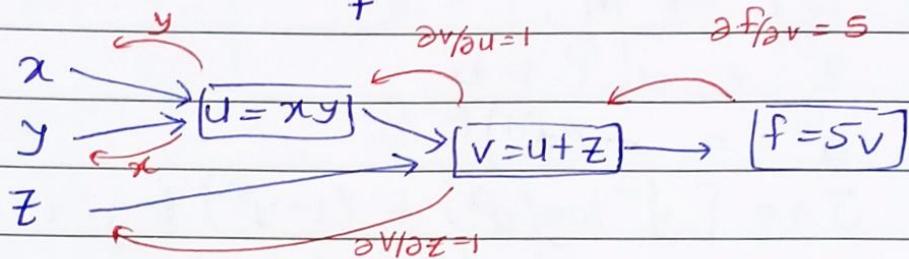
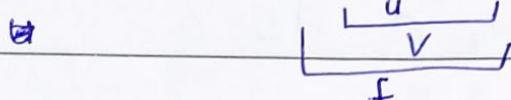
$$w \rightarrow w - \alpha \frac{dJ(w)}{dw}$$

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

## # Computation Graphs :

$$f(x, y, z) = 5(xy + z)$$



This is useful in cases where we need to optimize final o/p variable.

i.e.  $J(w, b)$  cost function ✓

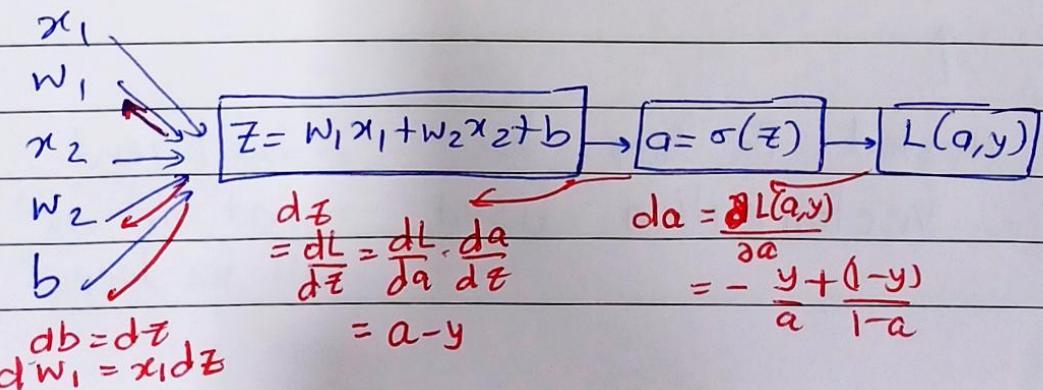
Forward prop  $\rightarrow$  calc.  $J$  values.

Back prop  $\rightarrow$  calc. derivatives

$$\therefore \frac{\partial f}{\partial x} = \frac{\partial f}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial x} = 5 \times 1 \times y = 5y$$

$\frac{dJ}{da}$  or "da" in short

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



But this was only for 1 training ex.

For m examples:

$$J=0; dw_1=0; dw_2=0; db=0$$

for  $i=1 \rightarrow m$ :

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

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

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

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

this is not separate for indiv. eg.

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \right\} \begin{array}{l} n=2 \\ \text{for loop run 2 times.} \end{array}$$

to eliminate for-loop:  
 $dw = np.zeros((n, 1))$   
 $dw += x^{(i)} dz^{(i)}$

$$J / = m$$

$$dw_1 / = m \quad dw_2 / = m \quad db / = m$$

$$\left. \begin{array}{l} w_1 := w_1 - \alpha dw_1 \\ w_2 := w_2 - \alpha dw_2 \\ b := b - \alpha db \end{array} \right\} \begin{array}{l} \text{update} \\ \text{parameters (weights)} \end{array}$$

↑

Nested for loop here  $\times$  Inefficient  
 $\therefore$  Vectorization used. → got rid of all for loops.

np.zeros((0, 1)) creates 2D arr with

0 rows 1 col  $\Rightarrow$  no elem.

→ To initialize empty array

Date : .....

MON TUE WED THU FRI SAT SUN

## # Vectorization

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

$$w = \begin{bmatrix} i \end{bmatrix}$$

$\mathbb{R}^{n_x}$

$$x = \begin{bmatrix} i \end{bmatrix}$$

$\mathbb{R}^{n_x}$

To calc. dot product:

Non-vectorized

$$z = 0$$

for i in range  $n_x$

$$z += w[i] \cdot x[i]$$

$$z += b$$

Vectorized

$$z = np.dot(w, x) + b$$

GPU & CPU  $\rightarrow$  Parallelization instructions  
(SIMD)

## # Vectorizing Log. Reg.

Date : .....

MON TUE WED THU FRI SAT SUN



$$z_1 = w^T x^1 + b \quad - \quad a_1 = \sigma(z^1)$$

$$X = \begin{bmatrix} | & | & | \\ x^1 & x^2 & \cdots & x^m \\ | & | & | \end{bmatrix} \Rightarrow (n_x, m)$$

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \cdots & z^{(m)} \end{bmatrix}$$

$$= \underbrace{w^T X}_{(1, n_x)} + \underbrace{b}_{(1, m)}$$

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

↑ automatically converts to (1, m)

$$A = [a^{(1)}, a^{(2)}, a^{(3)}, \dots, a^{(m)}] = \sigma(Z)$$

For backpropagation vectorizing :

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

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

Date : .....

MON TUE WED THU FRI SAT SUN  
      

$$\begin{aligned} d\bar{Z} &= [dz^{(1)} \dots dz^{(m)}] \\ &= A - Y \end{aligned}$$

$$db = \frac{1}{m} np \cdot \text{sum}(dZ)$$

$$\begin{aligned} dw &= \frac{1}{m} X(dZ)^T \\ &= \frac{1}{m} \left[ \begin{array}{c} x_i^{(1)} \dots x_i^{(m)} \end{array} \right] \left[ \begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right] \end{aligned}$$

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

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

## # Broadcasting

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}$$

$$= \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

Date : .....

MON TUE WED THU FRI SAT SUN



$$\begin{array}{c} (m,n) \\ \text{matrix} \end{array} \quad \begin{array}{c} + \\ \times \end{array} \quad \begin{array}{c} (1,n) \\ (m,1) \end{array} \quad \longrightarrow \quad \begin{array}{c} (m,n) \\ (m,n) \end{array}$$

## # Numpy vectors

$a = np.random.randn(5)$

O/P:  $[ \dots, \dots, \dots, \dots, \dots ]$

$a.shape \rightarrow (5,)$

Rank 1 Array

DON'T USE LIKE THIS

$a.T$  same.

→ It doesn't behave like either  
row / column vector.

Instead,  $a = np.random.randn(5, 1)$

To check: assert  $(a.shape == (5, 1))$

If there is Rank 1:  $a.reshape(5, 1)$

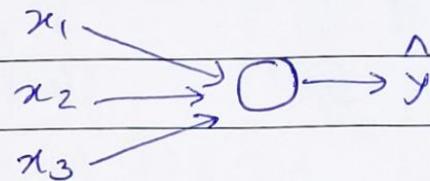
$$\hat{y} = P(y=1 | x)$$

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

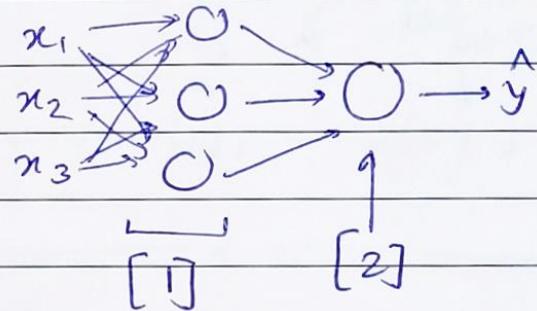
$$\text{If } y=0 : P(y|x) = 1 - \hat{y}$$

$$\Rightarrow P(y|x) = (\hat{y})^y \cdot (1-\hat{y})^{(1-y)}$$

## # Neural Network Representation:

Logistic Regr:

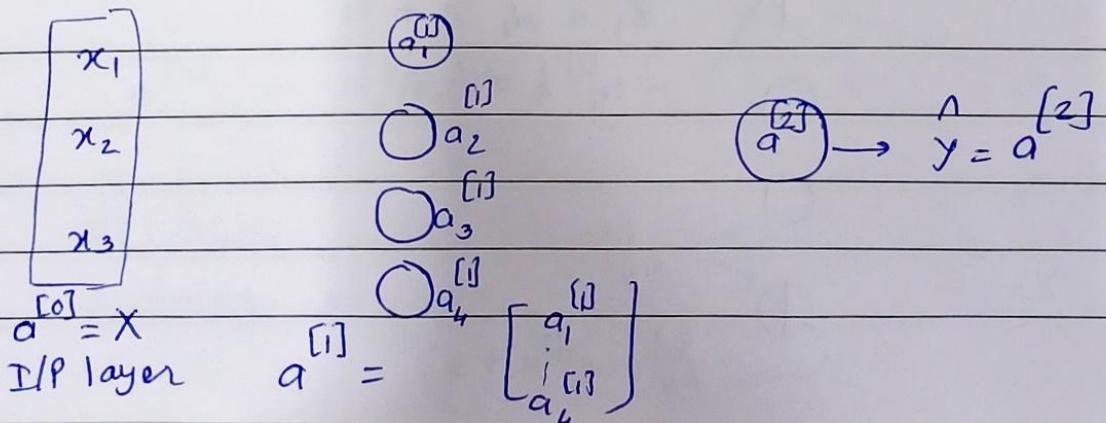
$$x \begin{matrix} \xrightarrow{w} \\ \xrightarrow{b} \end{matrix} Z = w^T x + b \rightarrow a = \sigma(Z) \rightarrow L(a, y)$$

Neural Net:

$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \\ L(a^{[2]}, y) &\leftarrow a^{[2]} = \sigma(z^{[2]}) \leftarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \end{aligned}$$

## # 2 Layer NN representation:

↳ 1 hidden + 1 o/p.



## Hidden layer parameters:

$$w^{[1]} \quad \& \quad b^{[1]}$$

↓                  ↓

shape:  $(4, 3)$        $(4, 1)$

$b_{CZ}$        $b_{CZ}$  3 I/P features  
4 nodes      in hidden layer.

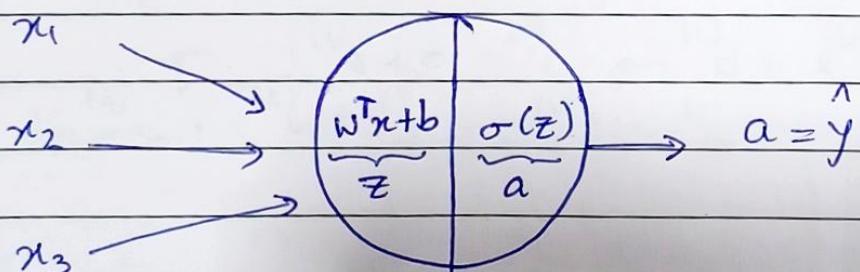
## O/P Layer Parameters:

$$w^{[2]} \quad \& \quad b^{[2]}$$

↓                  ↓

$(1, 4)$        $(1, 1)$

$$x \rightarrow$$



Logistic

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}$$

$$a_4^{[1]}$$

$[l] \leftarrow$  layer no-

$a^l_i \leftarrow$  node in layer

Date : .....

MON TUE WED THU FRI SAT SUN

Vectorising again:

$$\underbrace{\begin{bmatrix} \cdots & w_1^{[l]T} \\ \cdots & w_2^{[l]T} \\ \cdots & w_3^{[l]T} \\ \cdots & w_4^{[l]T} \end{bmatrix}}_{W^{[l]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[l]} \\ b_2^{[l]} \\ \vdots \\ b_4^{[l]} \end{bmatrix}$$

$$\therefore z^{[l]} = \begin{bmatrix} z_1^{[l]} \\ \vdots \\ z_h^{[l]} \end{bmatrix} \quad a^{[l]} = \begin{bmatrix} a_1^{[l]} \\ \vdots \\ a_h^{[l]} \end{bmatrix} = \sigma(z^{[l]})$$

Given IIP

$$z^{[1]} = w^{[1]} \cdot a^{[0]} + b^{[1]}$$

(4,1) (4,3) (3,1) (4,1)

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

(4,1) (4,1)

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

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

All this was for a single training example.

↑  
hidden  
units

↑ training eg. →

Date : .....

MON TUE WED THU FRI SAT SUN

for  $i = 1$  to  $m$ ,

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

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

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

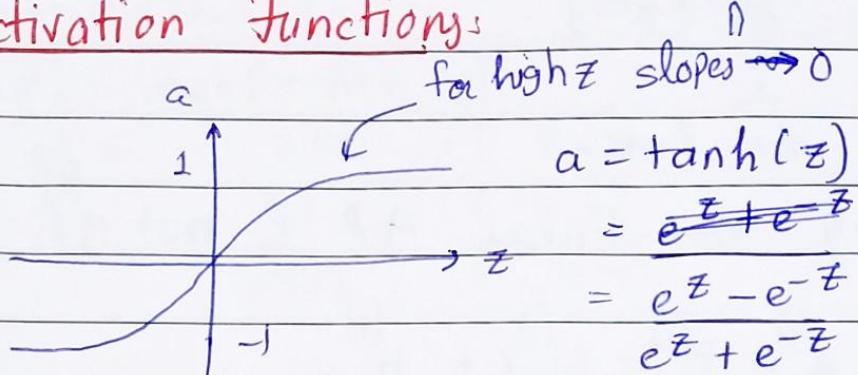
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$x = \begin{bmatrix} & & \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ & & \end{bmatrix} \in (n_n, m)$$

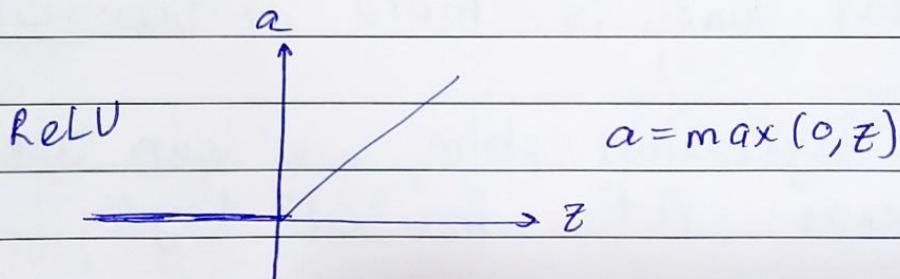
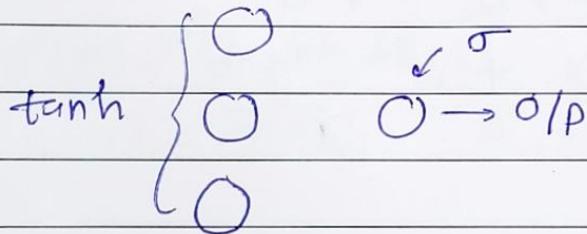
$$z^{[1]} = \begin{bmatrix} z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & \dots & | \end{bmatrix}$$

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

## # Activation Functions:



almost  
tanh, always superior than sigmoid

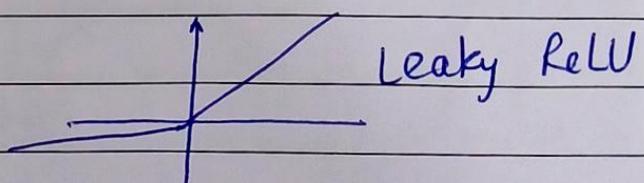


If binary classif.  $\rightarrow$  in its o/p layers  
generally  $\sigma$  used.

In all other layers  $\rightarrow$  ReLU can be  
used.

Try all a.f.

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



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

↑

Act. func.

Why non-linear A.F & not  $g^{[1]} = z$

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

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

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

$$= \underbrace{(w^{[2]}w^{[1]})}_{w'} x + \underbrace{w^{[2]}b^{[1]} + b^{[2]}}_{b'}$$

Linear func. is more or less USELESS

For regression pbm, we can use  
Linear A.F. for OLP layer

\* ReLU  $\rightarrow g(z) = \max(0, z)$

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

## ★ Leaky ReLU:

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

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

## † tanh func.

$$g(z) = \tanh z$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

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

## ★ Sigmoid func:

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

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

$$= a(1 - a)$$

## # Gradient descent for NN:

$$dW^{[l]} = \frac{dJ}{dW^{[l]}} \quad , \quad db^{[l]} = \frac{dJ}{db^{[l]}}$$

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

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

keepdims  
- prevents Rank 1 array

Date : .....

MON TUE WED THU FRI SAT SUN

## Back Propagation:

$$dZ^{[2]} = A^{[2]} - Y$$

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

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

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

$$dW^{[2]} = dZ^{[2]} a^{[1]T}$$

$$db^{[2]} = dZ^{[2]}$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} \cdot g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dZ^{[1]} \cdot x^T$$

$$db^{[1]} = dZ^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

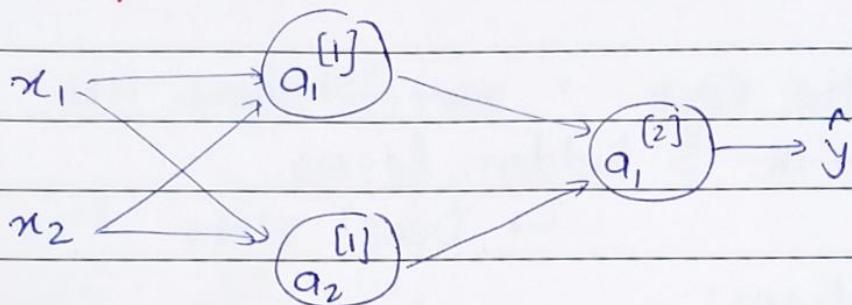
$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

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

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} x^T$$

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

# # Why initialise weights randomly?



$$\text{If } w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

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

Both these units will compute for the same function.

$\therefore$

$$w^{[1]} = \text{np.random.randn}((2, 2)) * \underbrace{0.01}_{\text{so that}}$$

$$b^{[1]} = \text{np.zeros}((2, 1))$$

so that in sigmoid derivative don't  $\rightarrow 0$

$$\left. \begin{array}{l} z^{[1]} = w^{[1]}x + b^{[1]} \\ \hookrightarrow \text{is small} \\ a^{[1]} = g^{[1]}(z^{[1]}) \end{array} \right\} \text{so that}$$

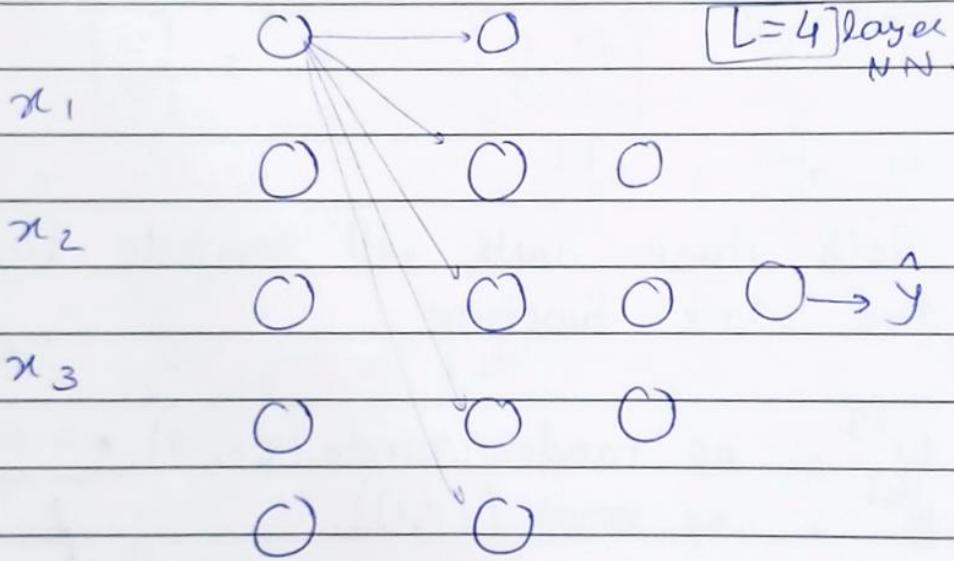
# # Deep Neural Networks

Logistic Regr → very "shallow" NN.

NN with 5 hidden layers

↪ "Deep" NN.

Notations:



$$n^{[l]} = \text{no. of units in layer } l.$$

$$n^{[1]} = 5 \quad n^{[2]} = 5 \quad n^{[3]} = 3$$

$$n^{[4]} = n^{[L]} = 1 \quad n^{[0]} = n_x = 3$$

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

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

$w^{[l]}$  = weights for  $z^{[l]}$

$$\begin{aligned} x &= a^{[0]} \\ \hat{y} &= a^{[L]} \end{aligned}$$

Date : .....

MON TUE WED THU FRI SAT SUN



## ## Forward Prop in DNN (Deep)

① For single tr. ex.  $\rightarrow$  all  $z, a, b$  are small.

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

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

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

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

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

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

② Vectorized:

$$\begin{aligned} z^{[1]} &= w^{[1]}x + b^{[0]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned} \quad \left. \right\} \text{For } L=1$$

$$\begin{aligned} z^{[2]} &= w^{[2]}A^{[1]} + b^{[1]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) \end{aligned} \quad \left. \right\} \text{For } L=2$$

General:  $z^{[l]} = w^{[l]}A^{[l-1]} + b^{[l]}$

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

A for loop  $l=1$  to  $L$  is necessary //

$Z \rightarrow$  called as vector activations.

Date : .....

MON TUE WED THU FRI SAT SUN

## # Knowing dimensions of Parameters $w^{[l]}$ & $b^{[l]}$ :

o o

eg.

x<sub>1</sub> o

o o o o o →

x<sub>2</sub> o o o o

o o o →

o o

L: 0 | 1 | 2 | 3 | 4 | 5  
 $n^{[0]} = 2$     $n^{[1]} = 3$  - - -

$$\check{Z}^{[l]} = \cancel{W^{[l]}} X + b^{[l]}$$

↓      ↓  
 $(3,1)$      $(2,1)$       General:  
or  $(n^{[l]}, 1)$    or  $(n^{[0]}, 1)$   
↓  
 $(3,2)$        $b^{[l]}: (n^{[l]}, 1)$   
or  $(n^{[l]}, n^{[0]})$

↓  
or  $(3,1)$   
 $(n^{[l]}, 1)$

$Z$  & a same dimension.

Date : .....

MON TUE WED THU FRI SAT SUN

For Backprop:  $dW^{[l]}$ : same as  $W^{[l]}$   
 $= (n^{[l]}, n^{[l-1]})$

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

For vectorized Implementation:

$$\text{capital } Z^{[l]} = W^{[l]} X + b^{[l]}$$

↓                  ↓                  ↓  
 $(n^{[l]}, m)$      $(n^{[l]}, n^{[0]})$      $(n^{[0]}, m)$   
↓                  ↓  
 $(n^{[l]}, m)$                   ↓ broadcasting

General: small  $Z^{[l]}, a^{[l]} = (n^{[l]}, 1)$

caps.  $Z^{[l]}, A^{[l]} = (n^{[l]}, m)$

$A^{[0]} = X = (n^{[0]}, m)$

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

Date : .....

MON TUE WED THU FRI SAT SUN

## # Forward & Backward functions :

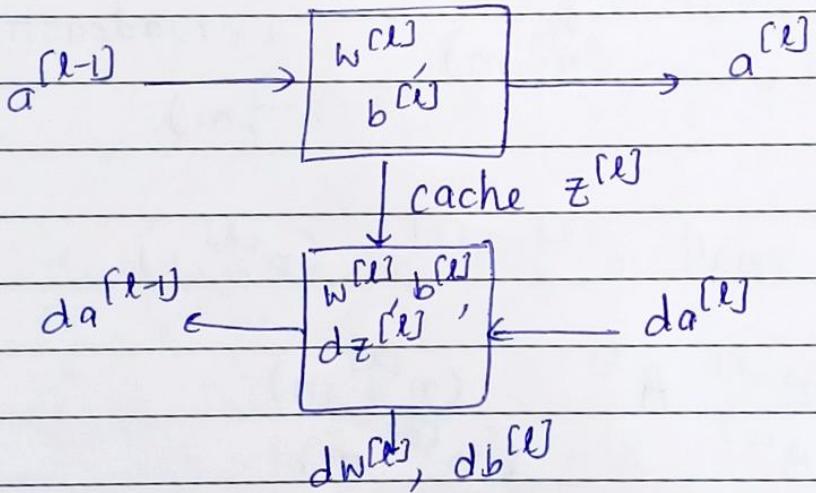
Forward: Layer  $l$ :

$$\text{I/P: } a^{[l-1]} \quad \text{O/P: } a^{[l]}$$

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

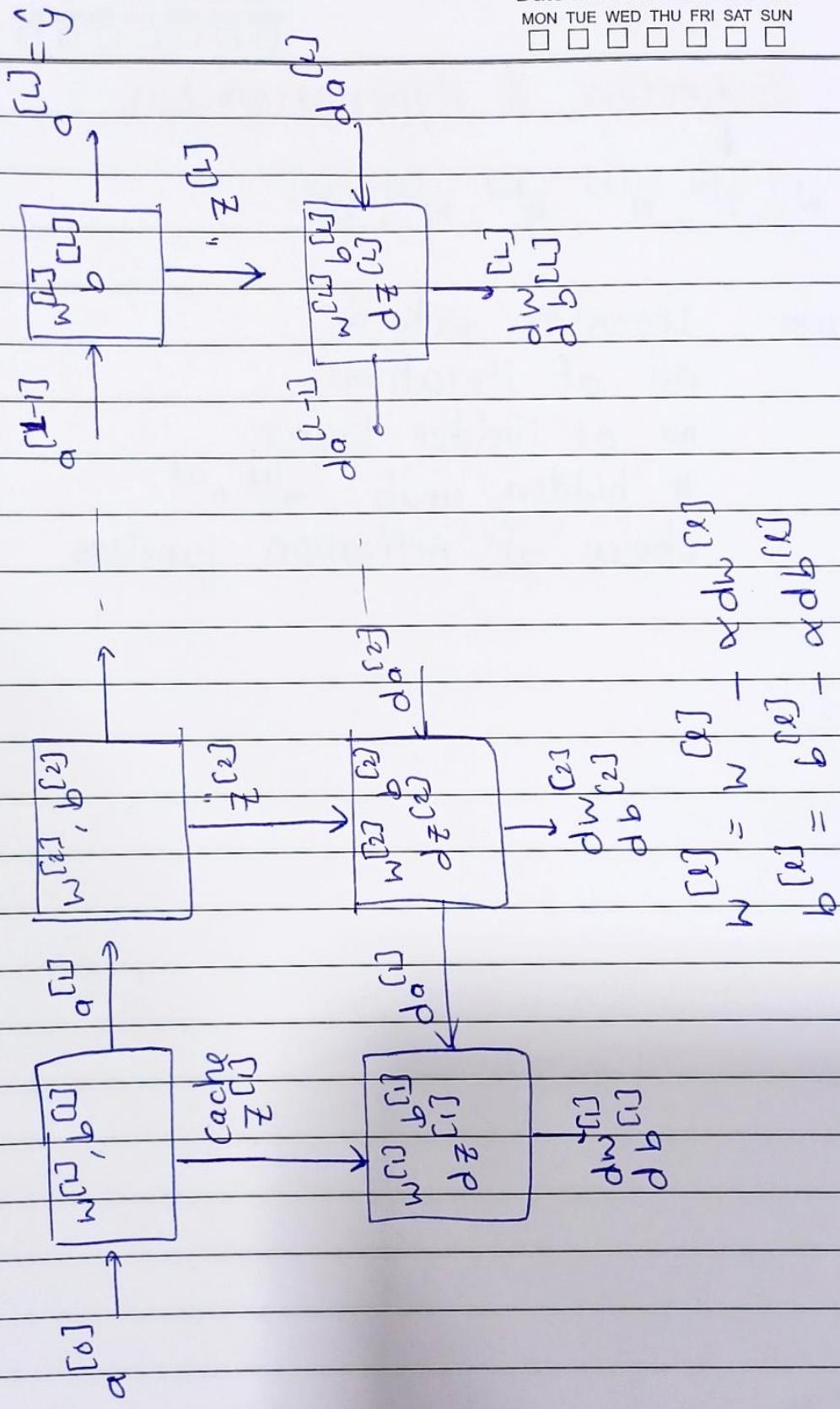
$$a^{[l]} = g^{[l]}(z^{[l]})$$

Backward: I/P:  $da^{[l]}$  O/P:  $da^{[l-1]}$



Date : .....

MON TUE WED THU FRI SAT SUN



Date : .....

MON TUE WED THU FRI SAT SUN

## # Parameters & Hyperparameters:

$w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]}, \dots$

Hypers: Learning rate  $\alpha$ ,  
no. of iterations,  
no. of hidden layers,  
# hidden units  $n^{[1]}, n^{[2]}, \dots$   
choice of Activation function.