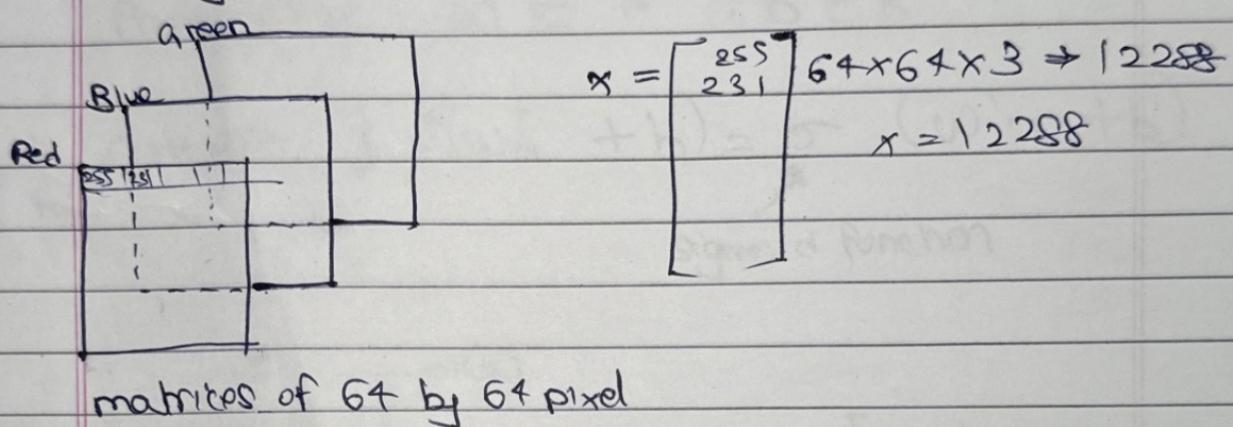
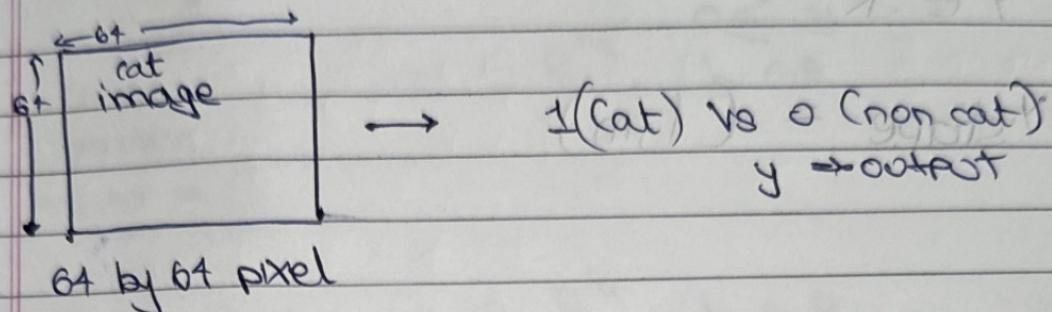


5-8-2024

WEEK - 2

Date
Page

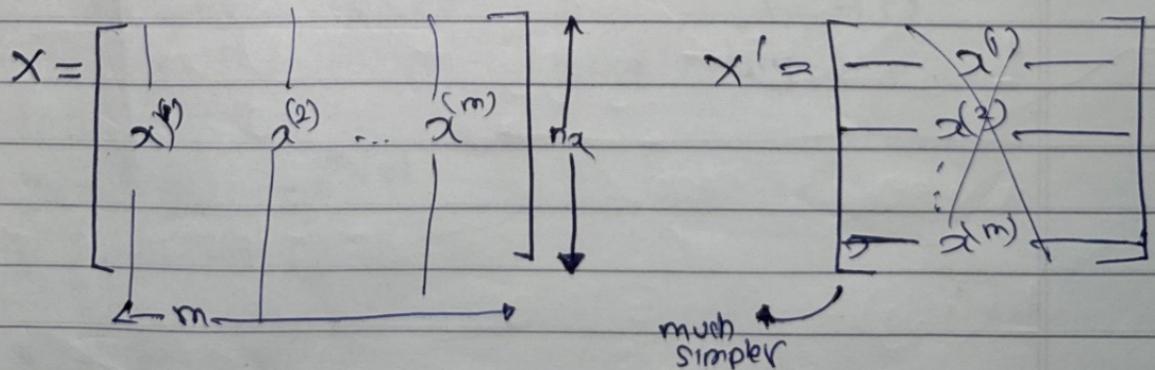
BINARY CLASSIFICATION :-



matrices of 64 by 64 pixel

Notation :-

(x, y) at R^{n_x} , $y \in \{0, 1\}$
 m -training examples : $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

 $M = M_{train}$ $M_{test} = \# \text{ test examples}$  $X \in R^{n_x \times m}$ $X.\text{shape} = (n_x, m)$

Date _____
Page _____

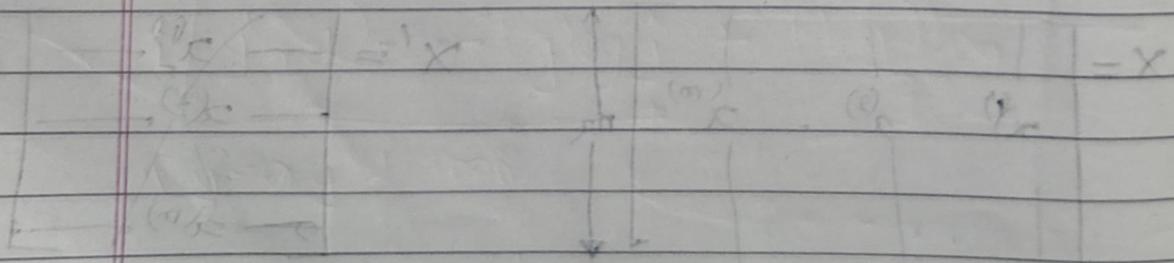
$$y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

~~$y \in \mathbb{R}^m$~~

$$y.\text{shape} = (1, m)$$

$$\begin{matrix} 228 & 228 & 228 \\ 228 & 228 & 228 \\ 228 & 228 & 228 \end{matrix} = x$$

$$\text{number of elements} = n$$



$$(n \times m) = \text{number of elements}$$

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

LOGISTIC REGRESSION :-

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

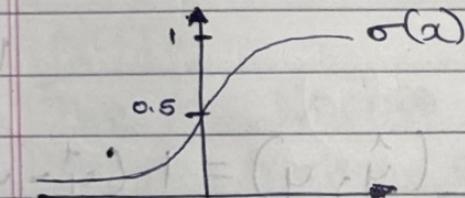
$x \in \mathbb{R}^{n_2}$, $0 \leq \hat{y} \leq 1$

input feature vector x

Parameters: $w \in \mathbb{R}^{n_2}$, $b \in \mathbb{R}$

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

sigmoid function



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

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

if z large $\sigma(z) \approx \frac{1}{1+0} = 1$

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

$$x_0 = 1, x \in \mathbb{R}^{n_2+1}$$

$$y = \sigma(w^T x)$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{n_2} \end{bmatrix}$$

LOGISTIC REGRESSION COST FUNCTION

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

Given $\{(x^{(i)}, y^{(i)})\}, (x^{(m)}, y^{(m)})$, want $\hat{y}^{(i)} \approx y^{(i)}$

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

$$(d + b^T w) \rightarrow (d + z^T w) \in \mathbb{R}$$

$x^{(i)}$
 $y^{(i)}$
 $z^{(i)}$

i-th example

$$\text{Loss (error) function: } L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

Loss function L is a function we will need to define to measure how good our output \hat{y} is when the true label is y

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

if $y=0$ then $\hat{y}=1$

$$L(\hat{y}, y) = -\log \hat{y} + 0 \\ = -\log \hat{y}$$

want $\log \hat{y}$ large, want \hat{y} large

if $y=0$:

$$L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow \text{want } \log 1-\hat{y} \text{ large}$$

want \hat{y} small

$$\text{cost function: } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

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

NOTE:-

The loss function computes the error for a single training example; the cost function is the average of the loss function of the entire training set.

$$(w) \hat{y}^{(i)} - w = : w$$

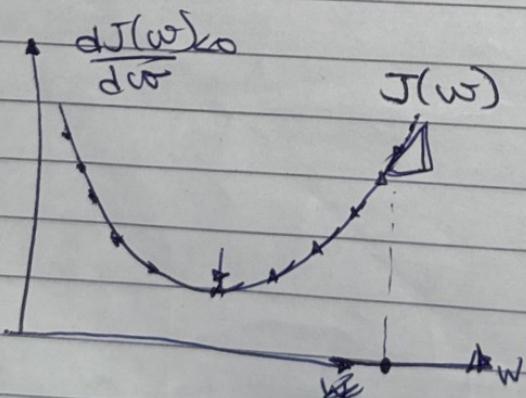
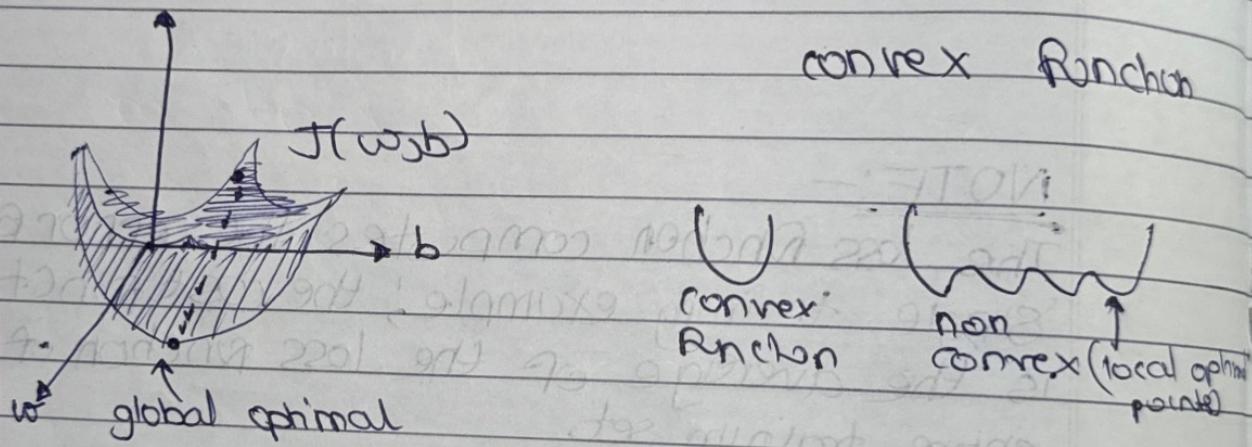
most of it to zero & cub

#GRADIENT DESCENT:

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

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

Want to find w, b that minimize $J(w, b)$



Repeat ?

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

3

$d w \Rightarrow$ derivative term

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

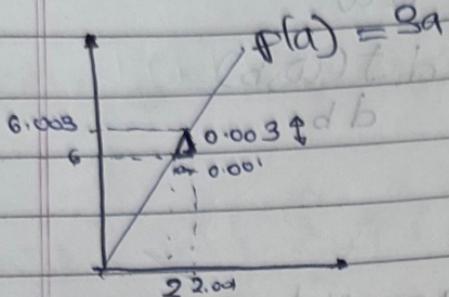
$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

$$\frac{\partial J(w, b)}{\partial w_k}$$

$\partial \rightarrow$ partial derivative

$d \rightarrow$ only 1

DERIVATIVES



$$\text{slope} = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.001} = 3$$

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

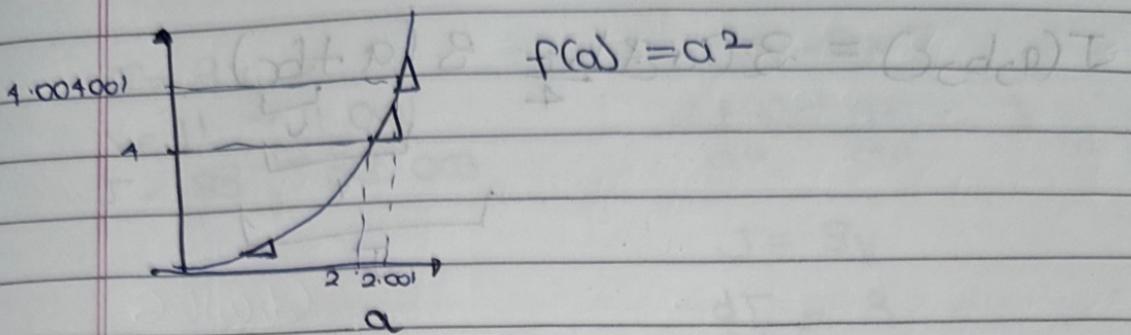
$$a = 2.001 \quad f(a) = 6.003$$

$$a = 5 \quad f(a) = 15$$

$$a = 5.001 \quad f(a) = 15.003$$

$$\text{slope} = 3$$

MORE DERIVATIVE EXAMPLES (MOMO #1)



different slope

$$\therefore \text{slope} = f'(a) = 2a$$

infinitesimal \rightarrow very very small

$$f(a) = a^3$$

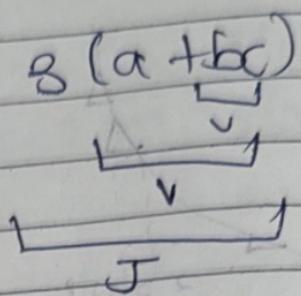
$$f'(a) = 3a^2$$

$$f(a) = \log(a)$$

$$f'(a) = 1/a$$

COMPUTATION GRAPH

$$J(a, b, c) = \frac{3(a+bc)}{4}$$



$$v = bc$$

$$v = a + v$$

$$J = 3v$$

$$s = a$$

$$3 = b$$

$$2 = c$$

$$\boxed{v = bc}$$

$$6$$

$$11$$

$$7$$

$$9$$

$$10$$

$$8$$

$$12$$

$$13$$

$$14$$

$$15$$

$$16$$

$$17$$

$$18$$

$$19$$

$$20$$

$$21$$

$$22$$

$$23$$

$$24$$

$$25$$

$$26$$

$$27$$

$$28$$

$$29$$

$$30$$

$$31$$

$$32$$

$$33$$

$$34$$

one-step of backward propagation on a computation graph yields derivative of final output variable

$$\delta_D = (0)^{\beta}$$

$$\delta_{\partial S} = (0)^{\beta}$$

$$(0)^{\beta} = (0)^{\beta}$$

$$(0)^{\beta} = (0)^{\beta}$$

DERIVATIVES WITH COMPUTATION GRAPH

$$J = 3v$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$f(a) = 3a$$

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

$$J = 3v$$

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

$$a = 5 \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dF}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

chain rule

$$\frac{dJ}{dv} = 3 \quad \frac{dv}{da} = 1$$

d · final output variable

d · variable

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{db}$$

$$= 3 \times 2 \\ = 6$$

$$b = 8 \rightarrow 8.001 \quad c = 2$$

$$v = bc \rightarrow 6.5 \rightarrow 6.002$$

$$\frac{dv}{db} = 2$$

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

$$\frac{dJ}{dc} = \frac{dJ}{dv} \cdot \frac{dv}{dc} = 9$$

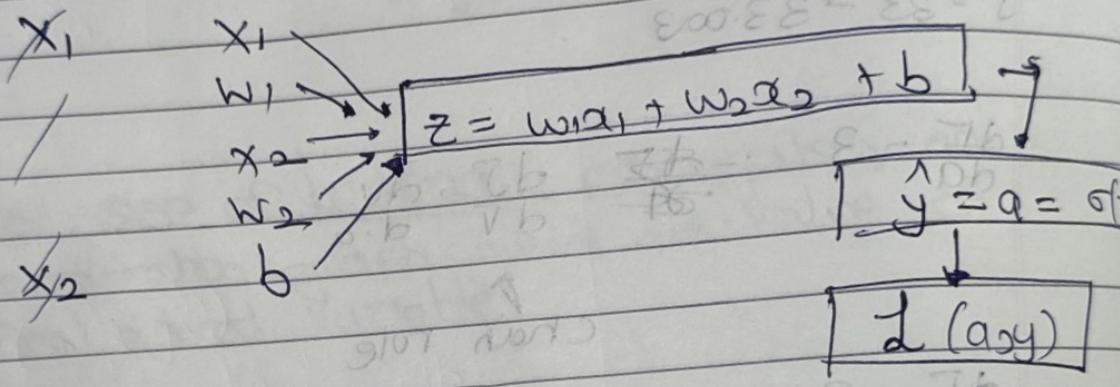
LOGISTIC REGRESSION / GRADIENT DESCENT

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

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

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

a = output of logistic regression
 y = ground truth label



$$\frac{da}{dz} = \frac{dL(a, y)}{da}$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$S = 0 \quad 100 \cdot 0 \leftarrow 0 = d$$

$$dL = \frac{dL}{dz} = \frac{dL(a, y)}{dz}$$

$$= a - y - \frac{ab}{ub} \cdot a(1-a)$$

$$= \frac{da}{dz} \cdot \frac{dL}{da}$$

$$\frac{ab}{ub}$$

$$P = \frac{ab}{ub} \cdot \frac{Ib}{ub} = \frac{Ib}{ub}$$

$$\frac{db}{dw_1} = "dw_1" = \alpha_1 dz$$

$$dw_2 = \alpha_2 dz$$

$$db = dz$$

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

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

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

$$\text{Diagram } (b) \text{ shows } \frac{db}{w_1 w_2} = (\alpha_1 w_1 + \alpha_2 w_2) dz$$

$$(w_1, w_2) = \alpha_{w_1 b}$$

$$\frac{db}{w_1 w_2} = \frac{(\alpha_1 w_1 + \alpha_2 w_2)}{w_1 w_2}$$

$$0 = db_c \quad 0 = w_1 b_c \quad 0 = w_2 b_c \quad 0 = T_c$$

$$d + \alpha_1 w_1 + \alpha_2 w_2 = 0$$

$$(\alpha_1 w_1 + \alpha_2 w_2) - \alpha_1 b_c = 0$$

$$[(\alpha_1 w_1 + \alpha_2 w_2) - \alpha_1 b_c] + \alpha_1 b_c = 0$$

$$\left. \begin{array}{l} \alpha_1 w_1 + \alpha_2 w_2 = 0 \\ \alpha_1 b_c = 0 \\ \alpha_1 = 0 \end{array} \right\} \begin{array}{l} \alpha_2 w_2 = 0 \\ \alpha_2 = 0 \end{array}$$

$$m = 1 db \quad m = 1 w_1 \quad m = 1 w_2$$

GRADIENT DESCENT ON M EXAMPLES!

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

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

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

$$d_{w_i}^{(i)} = (x^{(i)}, y^{(i)})$$

$$(x^{(i)}, y^{(i)}) \\ d_{w_1}^{(i)}, d_{w_2}^{(i)}, db^{(i)}$$

$$J=0, dw_1=0, dw_2=0, db=0$$

for $i=1 \text{ to } m$

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

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

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

$$dw_1 + = x_1^{(i)} dz^{(i)} \quad \uparrow n=2$$

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

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

$$J / = m$$

$$dw_1 / = m ; dw_2 / = m ; db / = m$$

$$\Delta w_1 = \frac{\partial J}{\partial w_1}$$

$$w_1' = w_1 - \alpha \Delta w_1$$

$$w_2' = w_2 - \alpha \Delta w_2$$

$$b' = b - \alpha \Delta b$$

VECTORISATION:-

Q] what is vectorization?

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

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

$$w^T x - d = 0$$

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

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

Non-vectorized

$$z = 0$$

for i in region(n-x):

$$z += w[i]^* x[i]$$

$$z += b$$

Vectorized

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

GPU $\not\sim$ SIMD

CPU $\not\sim$ single instrn

multiple data

MORE VECTORISATION EXAMPLES:-

Whenever possible, avoid explicit for-loops

$$v = Av$$

$$v_i = \sum A_{ii} v_j$$

NON-VECTORIZED

$$v = np.zeros((n, 1))$$

for i ...

for j ...

$$v[i] += A[i][j] * v[j]$$

$$v = np.dot(A, v)$$

vectors and matrix-valued functions

Say you need to apply the exponential operation on every element of a matrix/vector

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

import numpy as np
v = np.exp(v)

np.log(v)

np.abs(v)

np.maximum(v, 0)

$v^{1/2}$ $1/v$

$$\rightarrow v = np.zeros((n, 1))$$

→ for i in range(n):

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

VECTORISATION LOGISTIC REGRESSION

$$\rightarrow \hat{z}^{(1)} = w^T x^{(1)} + b$$

$$\rightarrow a^{(1)} = \sigma(\hat{z}^{(1)})$$

$$\hat{z}^{(2)} = w^T x^{(2)} + b$$

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

$$\hat{z}^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(\hat{z}^{(3)})$$

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

(n × m)
n × m
R

$$\begin{bmatrix} \hat{z}^{(1)} \\ \hat{z}^{(2)} \\ \vdots \\ \hat{z}^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix}$$

$1 \times m$

$$\text{softmax boundary} = \left[\begin{bmatrix} w^T x_1 + b \\ w^T x_2 + b \\ \vdots \\ w^T x_n + b \end{bmatrix} \right]$$

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

↖ (1, 1) R

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

Brookings

$$(v) \text{ dot } u$$

$$(v) \text{ dot } u$$

$$(v) \text{ dot } u$$

$$((1, 1)) 20195 \text{ qn} = 0$$

$$((1, 1)) 20195 \text{ qn} = 0$$

$$((1, 1)) 20195 \text{ qn} = 0$$

$$v \cdot u$$

VECTORIZING LOGISTIC REGRESSION'S GRADIENT

OUTLINE -

$$\begin{aligned} \frac{\partial}{\partial z^{(1)}} = a^{(1)} - y^{(1)} \\ \frac{\partial}{\partial z^{(2)}} = a^{(2)} - y^{(2)} \end{aligned}$$

$$dz = \begin{bmatrix} \frac{\partial}{\partial z^{(1)}} & \frac{\partial}{\partial z^{(2)}} & \frac{\partial}{\partial z^{(3)}} & \dots & \frac{\partial}{\partial z^{(m)}} \end{bmatrix}$$

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

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$dz = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix}$$

$$d\omega = 0$$

$$d\omega^1 = x^{(1)} dz^{(1)}$$

$$d\omega^2 = x^{(2)} dz^{(2)}$$

$$(D_{5b} + D_r) = + ab$$

:

$$d\omega^l = m$$

$$db = 0$$

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

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

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

$$db^l = m$$

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

$$= \frac{1}{n} np \cdot sm(dz)$$

$$d\omega = \frac{1}{n} X dz^T$$

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

$$\nabla J = \begin{bmatrix} \sum_{i=1}^m x^{(i)} d\epsilon^{(i)} + \dots + x^{(m)} d\epsilon^{(m)} \end{bmatrix}$$

n x 1

Implementing Logistic Regression

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

for $i = 1$ to m ,

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

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

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

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

$$\begin{cases} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \end{cases}$$

$$db += d\epsilon^{(i)}$$

$$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m$$

$$db = db/m$$

$$(5b) m \cdot \frac{1}{m} =$$

$$5b \times \frac{1}{m} = wb$$

$$\begin{bmatrix} m \\ 5b \end{bmatrix} \begin{bmatrix} dw_1 \\ dw_2 \\ db \end{bmatrix} = \begin{bmatrix} 0 \\ wb \\ 1 \end{bmatrix}$$

for iteration in range(1000):

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

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

$$A = \sigma(z)$$

$$dz = A - y$$

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

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

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

($\alpha = 0.001$) since $A = 10^7$
 $(+1) \text{epoch } (m) \backslash A^{1000} = \text{optimal}$

101	001	001	1
001	-	001	0
001	-	001	8
001	-	001	7
101	-	-	-

BROADCASTING IN PYTHON :-

Calories from Carbs, proteins, fats in 100g of different foods:

	Apples	Beef	Egg	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	165.0	99.0	0.9
Total	59.0	185.0	155.0	75.0

$$\frac{56}{59} \times 100 \Rightarrow 94.4\%$$

$$w - w = 0$$

$$dbx - d = d$$

Calculate % of calories for Carb, Protein, fat
Can you do this without an explicitly
for loop?

$$cal = A.sum(axis=0)$$

$$\text{percentage} = 100 * A / (\text{cal.reshape}(1, 4))$$

Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \cdot 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

GENERAL PRINCIPLE

$$(m \times n) + (l \times n) \quad m \neq l \quad (m \times n)$$

$$/ \quad (m \times l) \quad m \neq (m \times n)$$

$$(m \times l) + R$$

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

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

MATLAB/OCTAVE = bsx function

A NOTE ON PYTHON NUMPY VECTORS

$a = np.random.randn(5)$ } don't use
 $a.shape = (5,)$ }
rank 1 array "rank 1 array"

$a = np.random.randn(5, 1) \rightarrow a.shape = (5, 1)$

$a = np.random.randn(1, 5) \rightarrow a.shape = (1, 5)$

assert (a.shape == (5, 1))

80	80	101	801	801	801	801	801	801
202	202	202	202	202	202	202	202	202

QUICK TOUR OF JUPYTER/PYTHON NOTEBOOK

8/8/2024

WEEK-3 :- SHALLOW NEURAL NETWORKDate _____
Page _____

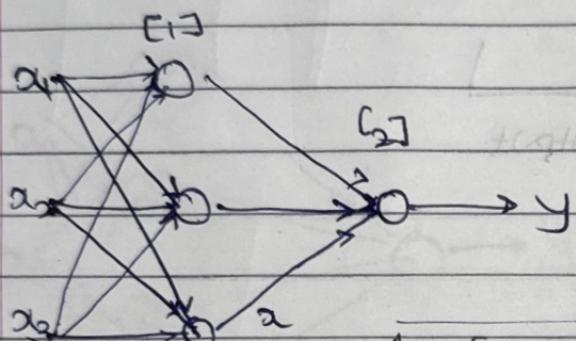
NEURAL NETWORK OVERVIEW:-

$$a_1 \\ a_2 \rightarrow \text{O} \rightarrow y = a$$

a
b

$$w \quad z = w^T a + b \rightarrow a = \sigma(z) = L(a|y)$$

b



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

$[i] \leftarrow$ layer

$\mathcal{E}^{(i)} \leftarrow$ ith example

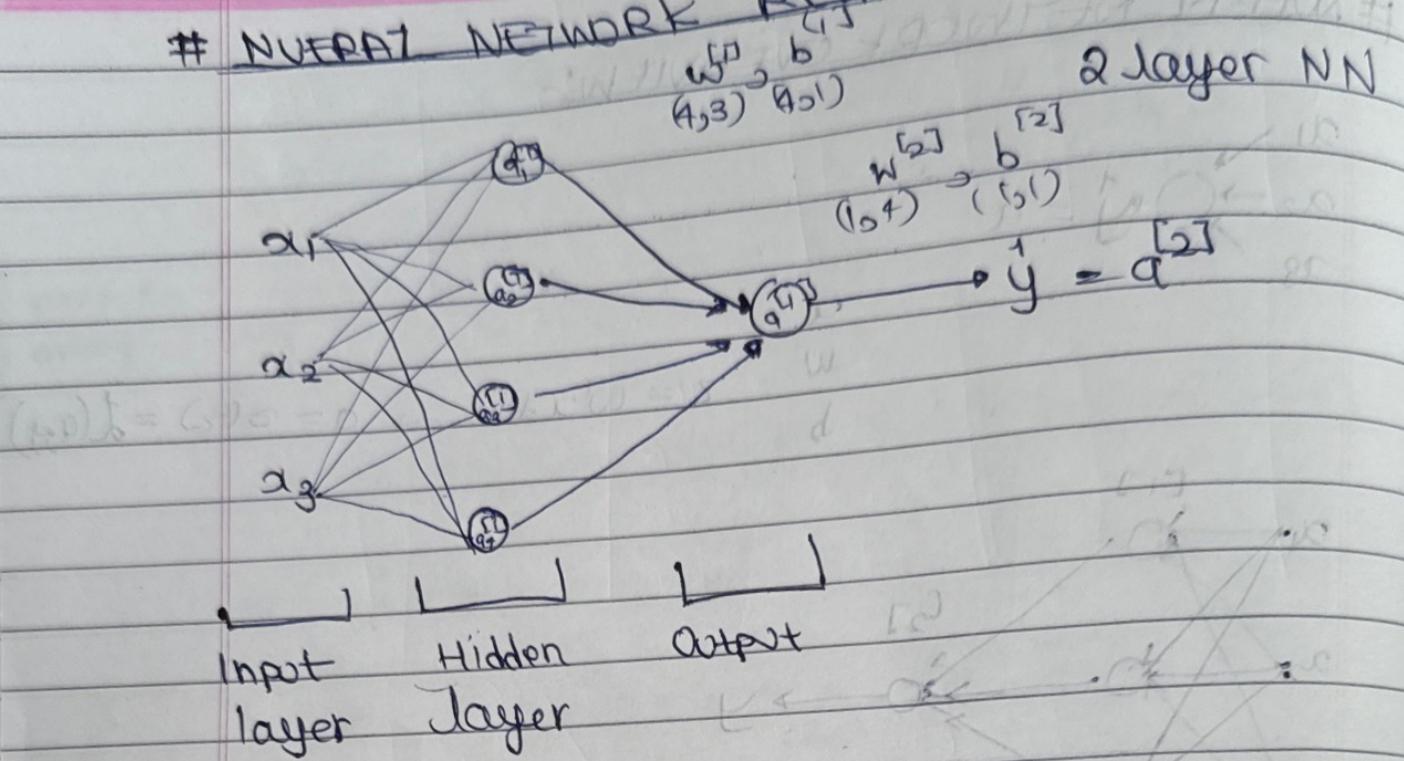
$$a^{[2]} = \sigma(z^{[2]})$$

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

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

$$L(a^{[2]}, y)$$

NEURAL NETWORK REPRESENTATION :-



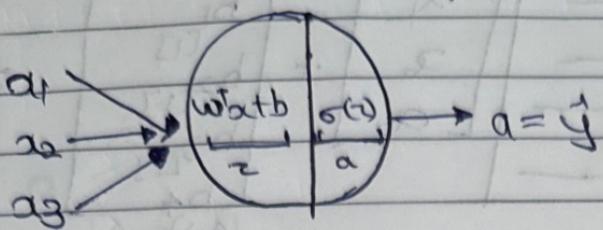
$$a^[[1]] = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$$w^[[1]] = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

$$(w^[[1]]) a^[[1]] = o^[[1]]$$

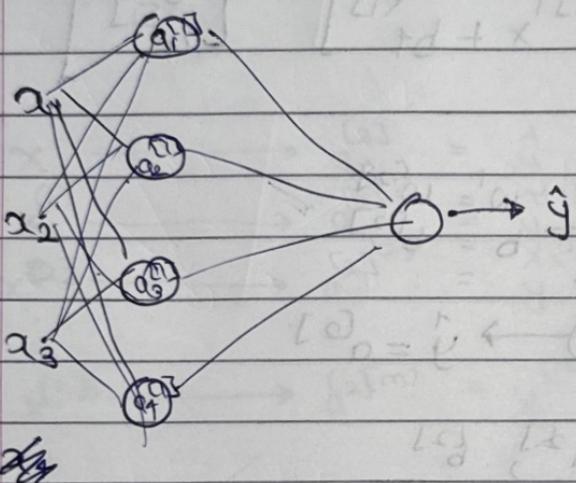
$$\begin{bmatrix} o^[[1]] \\ b^[[1]] \end{bmatrix}$$

COMPUTING A NEURAL NETWORKS OUTPUT:-



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

$$a = \sigma(z)$$



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, q_1^{[1]} = \sigma(z_1^{[1]})$$

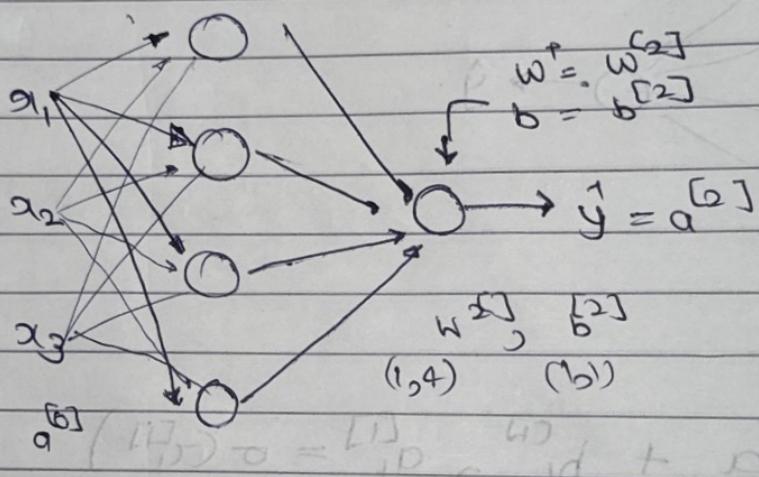
$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, q_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, q_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, q_4^{[1]} = \sigma(z_4^{[1]})$$

$$z^{[1]} = \begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

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



$$(1,2) z = w^T a_0 + b$$

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

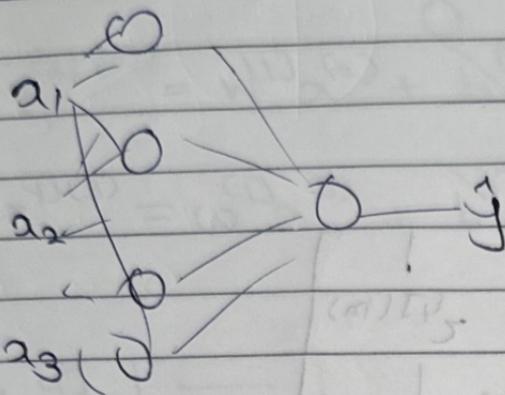
$$(1,2) a = w^T a_0 + b$$

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

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

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

VECTORISING ACROSS MULTIPLE EXAMPLES :-



$$\begin{aligned}
 x &\rightarrow a_1^{(1)} = y_1 \\
 x^{(1)} &\rightarrow a_1^{(2)} = y_2 \\
 x^{(2)} &\rightarrow a_1^{(3)} = y_3 \\
 &\vdots \\
 x^{(n)} &\rightarrow a_1^{(m)} = y^{(m)}
 \end{aligned}$$

for $i=1$ to m

$$\begin{aligned}
 z^{(1)(i)} &= w^{(1)}_1 x^{(1)(i)} + b^{(1)} \\
 z^{(2)(i)} &= w^{(2)}_1 x^{(2)(i)} + b^{(2)} \\
 a^{(2)(i)} &= \sigma(z^{(2)(i)})
 \end{aligned}$$

$$\begin{aligned}
 z^{(2)(1)} &= w^{(2)}_1 x^{(2)(1)} + b^{(2)} \\
 a^{(2)(1)} &= \sigma(z^{(2)(1)})
 \end{aligned}$$

$$\begin{aligned}
 x &= \begin{bmatrix} 1 & 1 & \dots & 1 \\ a_1^{(1)} & a_1^{(2)} & \dots & a_1^{(m)} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \\
 & \quad (n, m)
 \end{aligned}$$

$$\begin{bmatrix} z \\ z \end{bmatrix} = w^T x + b$$

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

$$\begin{bmatrix} z \\ z \end{bmatrix} = w^T x + b$$

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

$$z = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ \vdots & \vdots & \vdots \\ 1 & 1 & 1 \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$A^{[2]} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ g_1(1) & g_1(2) & \dots & g_1(m) \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$$P = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_n)$$

$$\begin{aligned} L_1 &= \alpha_1 P^{-1} \\ L_2 &= \alpha_2 P^{-1} \\ L_3 &= \alpha_3 P^{-1} \\ \vdots & \vdots \\ L_n &= \alpha_n P^{-1} \end{aligned}$$

$$\begin{aligned} P^{-1} &= M^{-1} \\ M &= \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \end{aligned}$$

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} = X$$

(contd.)

EXPLANATION FOR VECTORISATION IMPLEMENTATION

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

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

$$z^{(3)} = w^{(1)}x^{(3)} + b^{(1)}$$

$$w^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)}x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)}x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)}x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)} \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = z^{(1)}$$

$(500 \times 3) \times m = n$

U109

$$d + \frac{1}{2}W = 175$$

$$(10)^{10} \cdot 10 = 10 \Rightarrow$$

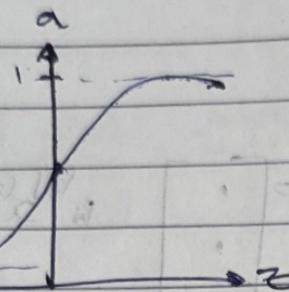
$(5000 \times 3) \times m$

F ACTIVATION FUNCTIONS

Sigmoid is called activation function

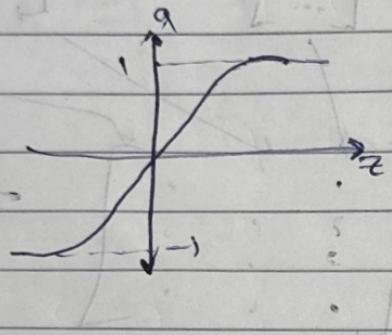
$g(z) \rightarrow$ non-sigmoid

only
for
binary
classification



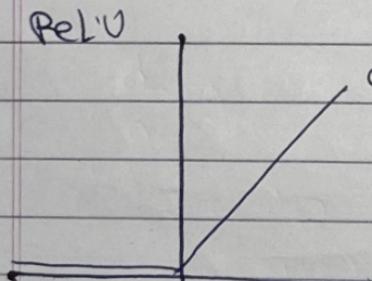
$$a = \frac{1}{1+e^{-z}}$$

strictly
superior



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

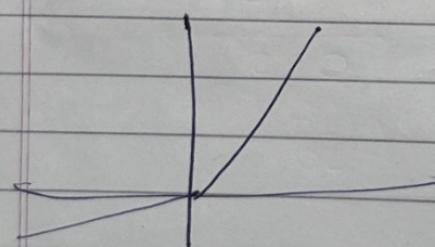
most
commonly
used



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

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

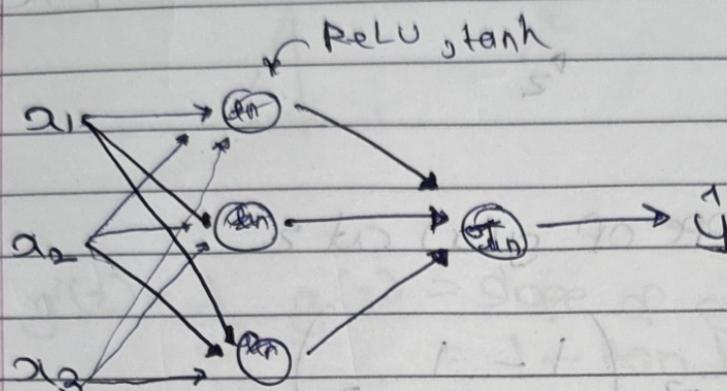
$$\text{or } a^{[j]} = g^{[j]}(z^{[j]})$$



$$\max(0, 0.01z)$$

NN
WHY DO YOU NEED ~~*~~ NON-LINEAR ACTIVATION
FUNCTION?

Activation function



Given x :

$$\begin{aligned} z^{(0)} &= w^{(0)}x + b^{(0)} \\ a^{(1)} &= g(z^{(1)}) = g(z^{(0)}) \\ z^{(1)} &= w^{(1)}a^{(0)} + b^{(1)} \\ a^{(2)} &= g(z^{(2)}) = g(z^{(1)}) \end{aligned}$$

non-linear activation function

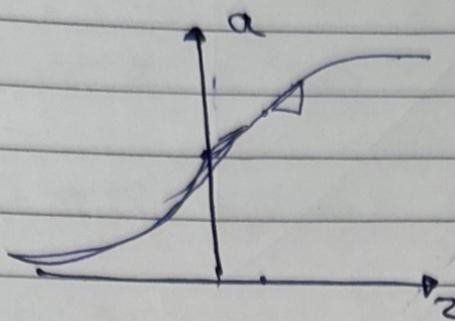
$$\begin{aligned} a^{(1)} &= f^{(1)} = w^{(1)}x + b^{(1)} \\ a^{(2)} &= f^{(2)} = w^{(2)}a^{(1)} + b^{(2)} \\ a^{(2)} &= w^{(2)}(w^{(1)}x + b^{(1)}) + b^{(2)} \end{aligned}$$

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

w' b'

$$= w'x + b'$$

DERIVATIVES OF ACTIVATION FUNCTIONS



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

$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$

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

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

$$z=0 \quad g(z) \approx 1$$

$$\frac{d}{dz} g(z) \approx 1(1-1) = 0$$

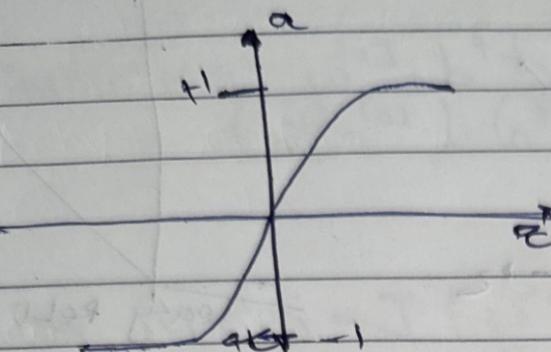
$$z=-10 \quad g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0 \cdot (1-0) = 0$$

$$z=0 \quad g(z) = \frac{1}{1+\frac{1}{2}} = \frac{1}{\frac{3}{2}} = \frac{2}{3}$$

$$d + x'w =$$

tanh activation function



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

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

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

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

$$z = 10 \quad \tanh(z) \approx 1$$

$$g'(z) \approx 0$$

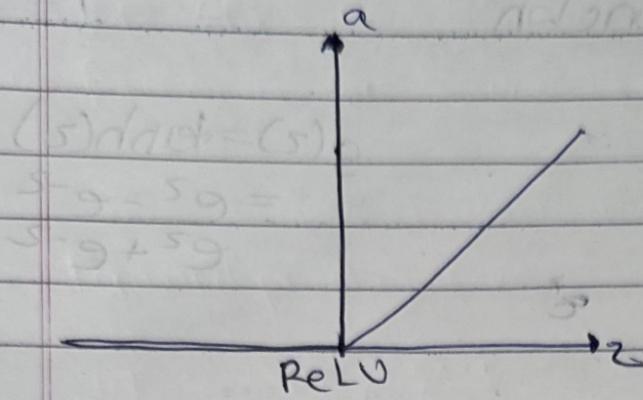
$$z = -10 \quad \tanh(z) \approx -1$$

$$g'(z) \approx 0$$

$$z = 0 \quad \tanh(z) = 0$$

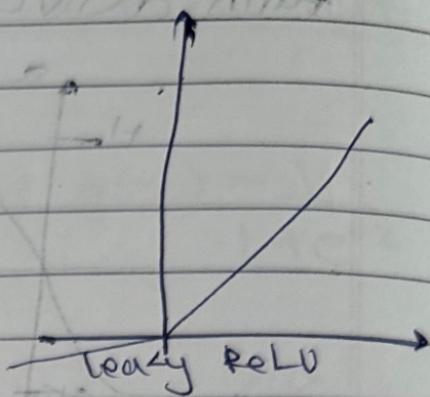
$$g'(z) = 1$$

ReLU and leaky ReLU



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

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



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

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

$$l \approx (5) \text{not } \alpha_1 = 5$$

$$\partial l \approx (5)^1 \alpha$$

$$\alpha \approx (5) \text{not } \alpha_1 = 5$$

$$l = (5)^1 \alpha$$

GRADIENT DESCENT FOR NEURAL NETWORKS :-

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 $(n^{[1]}, n^{[0]}) \quad (n^{[2]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$

$$\begin{aligned} \text{Cost Function} &= J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) \\ &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) \end{aligned}$$

Gradient descent

Repeat {

computer prediction $(\hat{y}^{(j)}, (-1, \dots, m))$

$$d w^{[1]} = \frac{\partial J}{\partial w^{[1]}} \rightarrow d b^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

$$\begin{aligned} w^{[1]} &:= w^{[1]} - \alpha d w^{[1]} \\ b^{[1]} &:= b^{[1]} - \alpha d b^{[1]} \end{aligned}$$

$$(a^{(j)})^{(l)} \quad \frac{1}{m} \sum_{i=1}^m s_h^{(i)} \quad (a^{(l)})^{(h)}$$

$$w^{(l)} = \frac{1}{m} \sum_{i=1}^m (a^{(i)})^{(l)} \quad (a^{(l)})^{(h)}$$

$$1 \times \frac{m}{m} s_h^{(i)} = \frac{1}{m} s_h^{(i)}$$

$$w^{(l)} = \frac{1}{m} \sum_{i=1}^m a^{(i)} s_h^{(i)}$$

~~for 200~~

formulas for computing derivatives

Forward propagation

$$\begin{aligned} z^{[1]} &= w^{[1]} X + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \\ z^{[2]} &= w^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) = \sigma(z^{[2]}) \end{aligned}$$

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

Back propagation:

$$\begin{aligned} dz^{[2]} &= A^{[2]} - y \\ dw^{[2]} &= \frac{1}{n} dz^{[2]} A^{[1]T} \end{aligned}$$

$$db^{[2]} = \frac{1}{n} np \cdot \text{sum}(dz^{[2]}, \text{axis}=1, \text{keepdim}=\text{True})$$

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

$(n^{[1]}, m)$ $(n^{[2]}, m)$
dot product

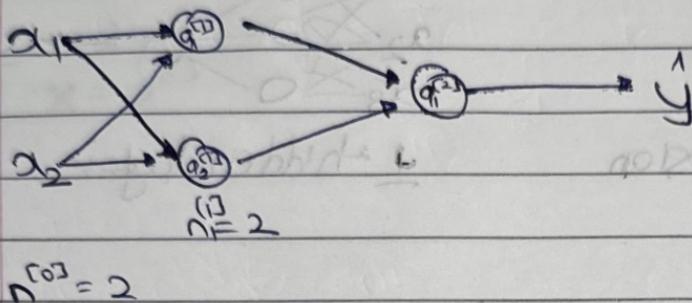
$$dw^{[1]} = \frac{1}{n} dz^{[1]} X^T$$

$$db^{[1]} = \frac{1}{n} np \cdot \text{sum}(dz^{[1]}, \text{axis}=1, \text{keepdim}=\text{True})$$

RANDOM INITIALIZATION:-

what happens if you initialize weights to zero?

$$w^{[0]} = [0 \ 0]$$



$$\begin{matrix} x_1 & 0 \\ x_2 & 0 \\ x_3 & 0 \end{matrix} \rightarrow \begin{matrix} y_1 & 0 \\ y_2 & 0 \\ y_3 & 0 \end{matrix}$$

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

$$a_1^{[1]} = a_2^{[1]} \quad d_{21}^{[1]} = d_{31}^{[1]}$$

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

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

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

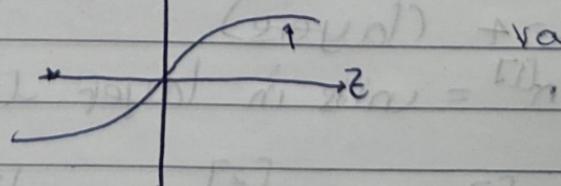
$$w^{[1]} = np.random.rand(2,2) \times 0.01$$

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

$$w^{[2]} = -np.random.rand(1,2)$$

$$b^{[2]} = 0$$

We like initializing with small values



0.01 constant
works for shallow

$$z^{[1]} = w^{[1]}(z^{[0]}) + b^{[1]}$$

8/8/2024

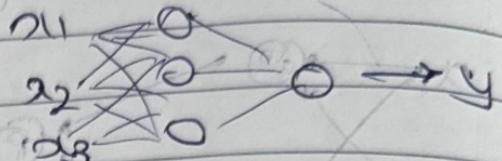
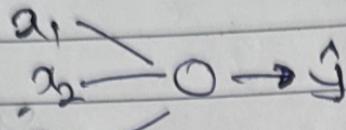
WEEK-4 :- DEEP NEURAL NETWORK

DEEP L-LAYER NEURAL NETWORK

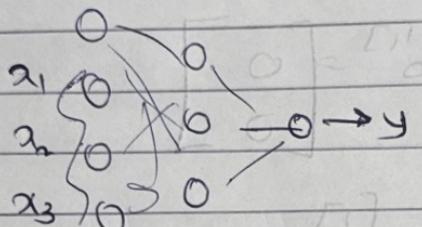
What is deep neural network?

single

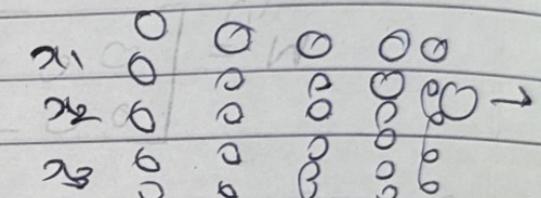
logistic regression



1 hidden layer

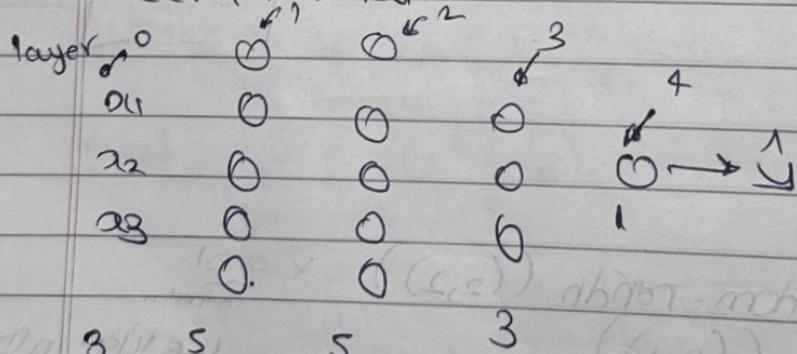


2 hidden layers



5 hidden layers
(deep)

Deep NN notation



$L = 4$

4 layer NN

$L = 4$ (layers)

$n^{[1]} = \text{unit in layer 1}$

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

$a^{(l)}$ = activations in layer x on input x

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

at $w^{(l)}_j = \text{weight for } z_j^{(l)}$
 $b^{(l)}$

$$x = d^{(l)}$$

$$d^{(l)}$$

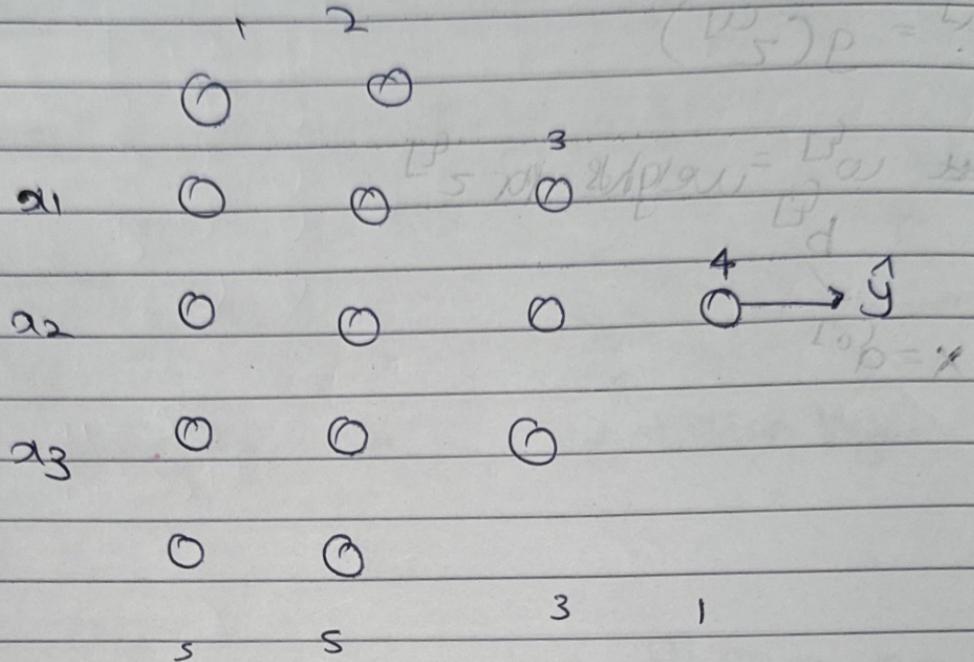
$$d^{(l)} + x^{(l)} w^{(l)} = \frac{1}{5} x \\ (w_5) d^{(l)} = \frac{1}{5} x$$

$$d^{(l)} + \frac{1}{5} d^{(l)} w^{(l)} = \frac{1}{5} \\ (w_5) d^{(l)} = \frac{1}{5}$$

$$d^{(l)} + \frac{1}{5} d^{(l)} w^{(l)} = \frac{1}{5} \\ \hat{v} = (\frac{1}{5}) \frac{d^{(l)}}{w^{(l)}} = \frac{1}{5}$$

$$d^{(l)} + \frac{1}{5} d^{(l)} w^{(l)} = \frac{1}{5} \\ (w_5) \frac{d^{(l)}}{w^{(l)}} = \frac{1}{5}$$

FORWARD PROPAGATION IN A DEEP NETWORK



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

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

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

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

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[3]} = g^{[3]}(z^{[3]}) = \hat{y}$$

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

$$a^{[4]} = g^{[4]}(z^{[4]})$$

Vectorized

$$\# \quad z^{[1]} = w^{[1]} A^{[1]} + b^{[1]} \\ A^{[1]} = g^{[1]}(z^{[0]})$$

$$z^{[2]} = w^{[2]} A^{[2]} + b^{[2]} \\ A^{[2]} = g^{[2]}(z^{[1]})$$

$$z = \begin{bmatrix} 1 & | & 1 & | & 1 \\ | & | & | & | & | \\ 1 & | & 1 & | & 1 \end{bmatrix}$$

#GETTING YOUR MATRIX DIMENSIONS RIGHT!

$$\begin{matrix} & 0 & 0 \\ a_1 & 0 & 0 \\ a_2 & 0 & 0 \\ a_3 & 0 & 0 \\ 0 & 0 & \end{matrix}$$

$$L + U = A \quad b = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$(L^T)U = b$$

$$L^T A = b$$

$$(L^T)^{-1} b = A$$

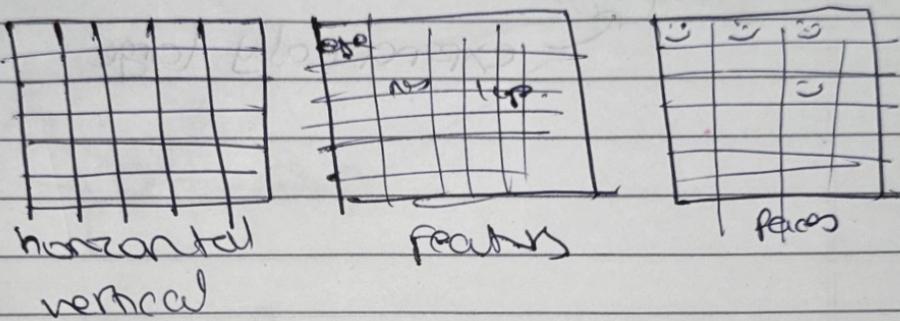
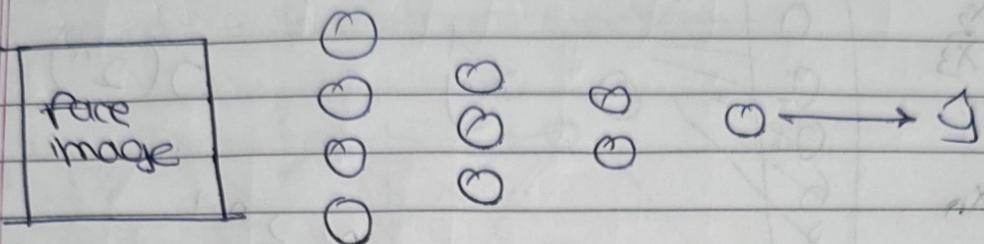
$$w_1 = \begin{pmatrix} 5 \\ 3 \end{pmatrix}$$

$$b_1 = \begin{pmatrix} 5 \\ 1 \end{pmatrix}$$

:

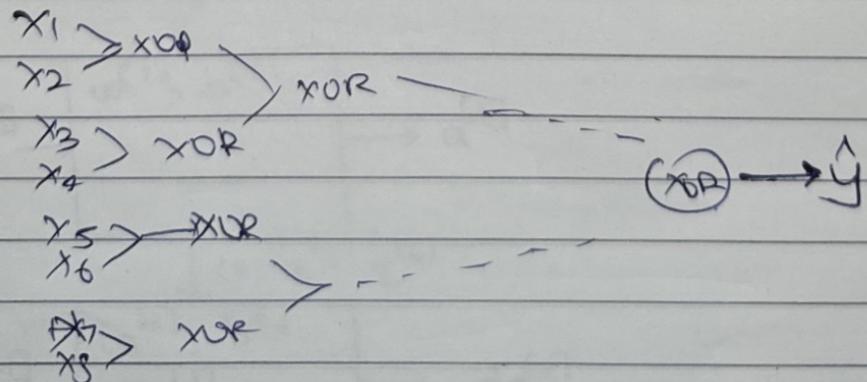
$$\begin{vmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{vmatrix} = 5$$

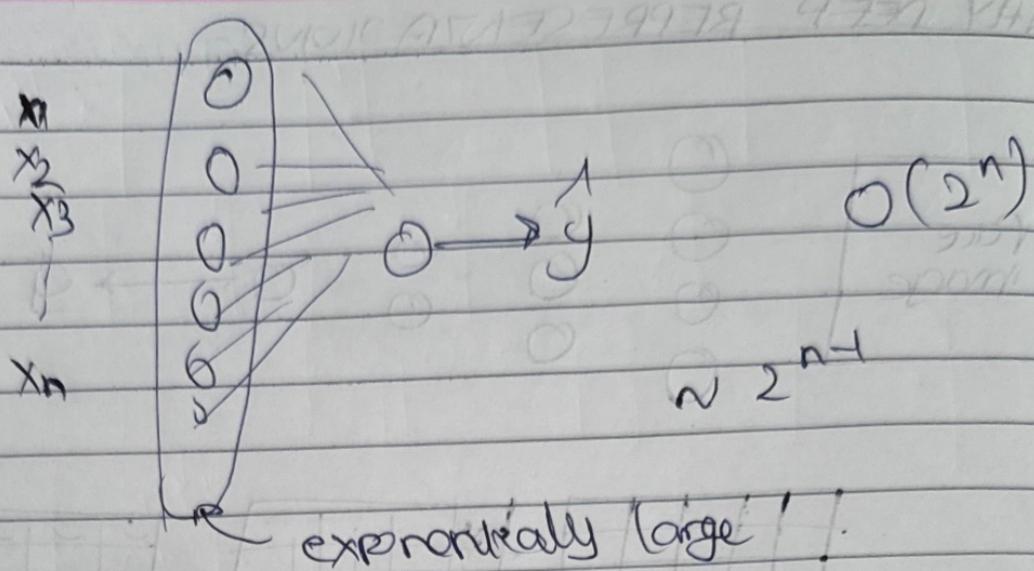
WHY DEEP REPRESENTATIONS?

CIRCUIT THEORY of DEEP LEARNING!

Informally; there are functions you can compute with a small 1-layer deep neural network that shallow networks require exponentially more hidden units to compute

$$y: x_1 \text{ XOR } x_2 \text{ XOR } x_3 \dots \text{ XOR } x_n$$





CALCULATING THE NUMBER OF DIFFERENT ARRANGEMENTS

Recall that there are 16 possible two-bit codes with a column (subset) such that each code has exactly one bit set to 1. There are 16 possible ways to combine the two bits.

$$2^2 = 4 \times 4 = 16$$

Let n

$$2^n < 2^{n+1}$$

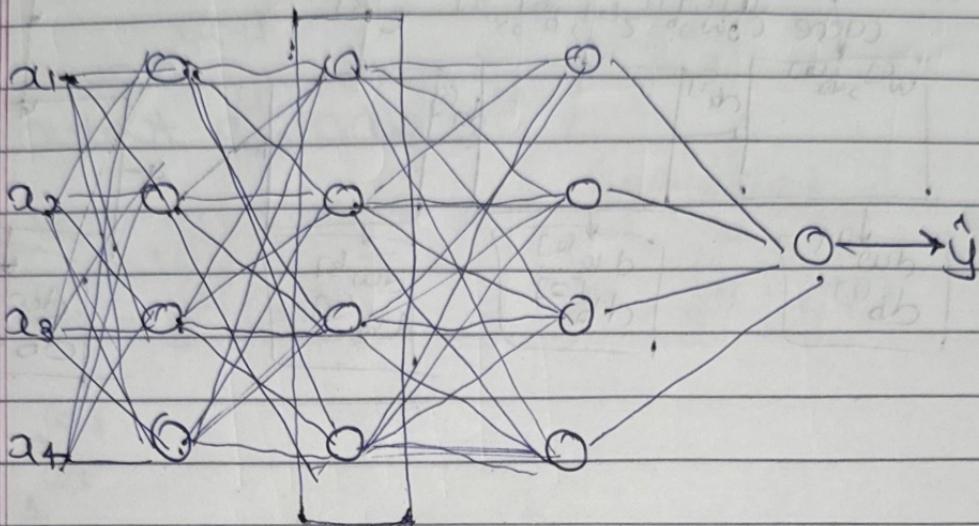
$$4096 < 8192$$

$$8192 < 16384$$

$$16384 < 32768$$

BUILDING BLOCKS OF DEEP NEURAL NETWORKS

forward & backward functions

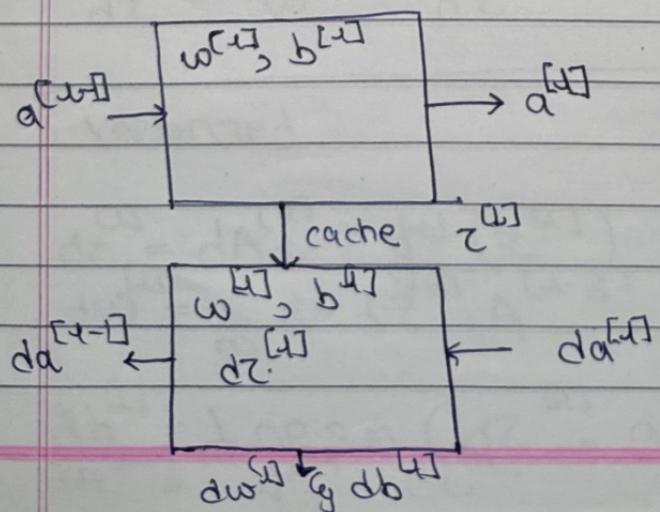


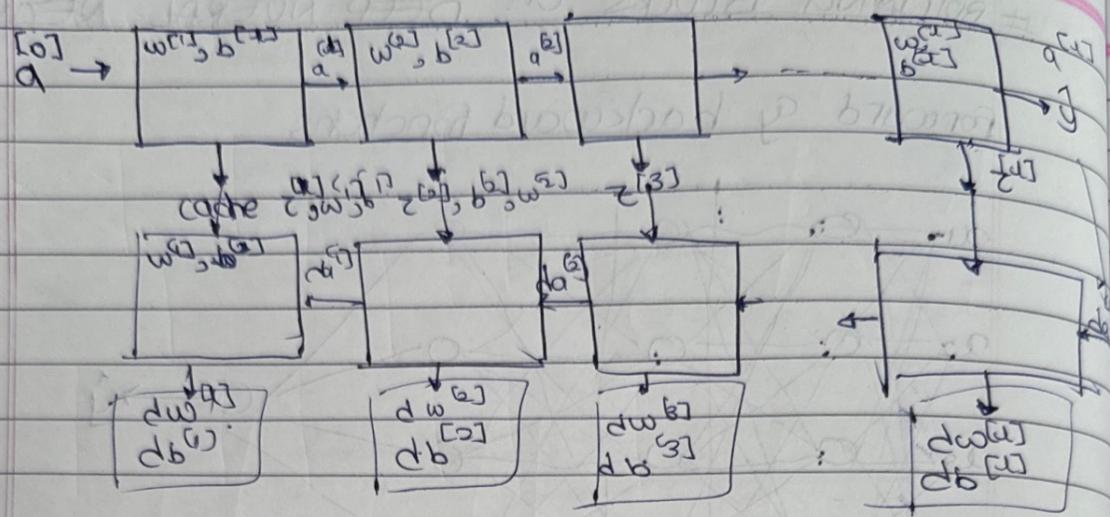
Layer l: $w^{[l]}, b^{[l]}$

→ forward: Input $a^{[l-1]}$ → output $a^{[l]}$
 $z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]}$ cache $z^{[l]}$
 $a^{[l]} = g(z^{[l]})$

→ Backward: Input $da^{[l]}$, output $da^{[l-1]}$
cache($z^{[l]}$), $dw^{[l]}$, $db^{[l]}$

layer l

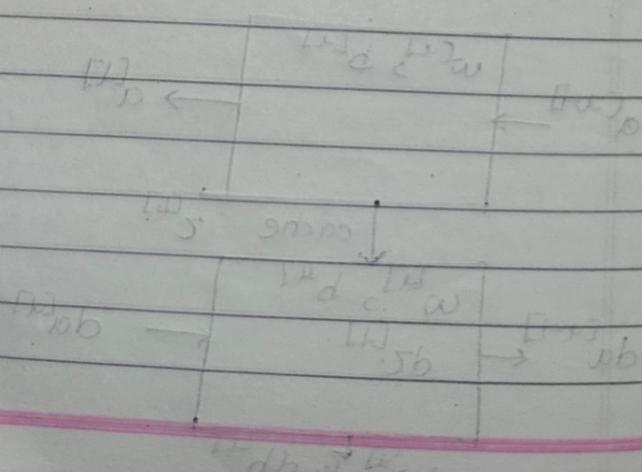




Bus
 PC : I -> P
 Register : R -> P
 Data : D -> P
 Address : A -> P
 Control : C -> P
 (D + A) = S
 (R + S) = D
 (C + S) = C

Bus
 in trame (P)
 Bus
 Bus

I need



FORWARD AND BACKWARD PROPAGATION :-

forward propagation for layer l

→ input $a^{[l-1]}$ → output $a^{[l]}$, cache ($z^{[l]}$)

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

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

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

Vectorized

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

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

 $\begin{matrix} [l] \\ a \\ 0 \\ [l] \\ A \end{matrix}$

$$x = A^0 \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow$$

backward propagation for layer l

→ Input $da^{[l]}$ → Output $da^{[l-1]}, dz^{[l]}, dw^{[l]}, db^{[l]}$

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

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

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

$$da^{[l-1]} = w^{[l]} \cdot dz^{[l]}$$

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

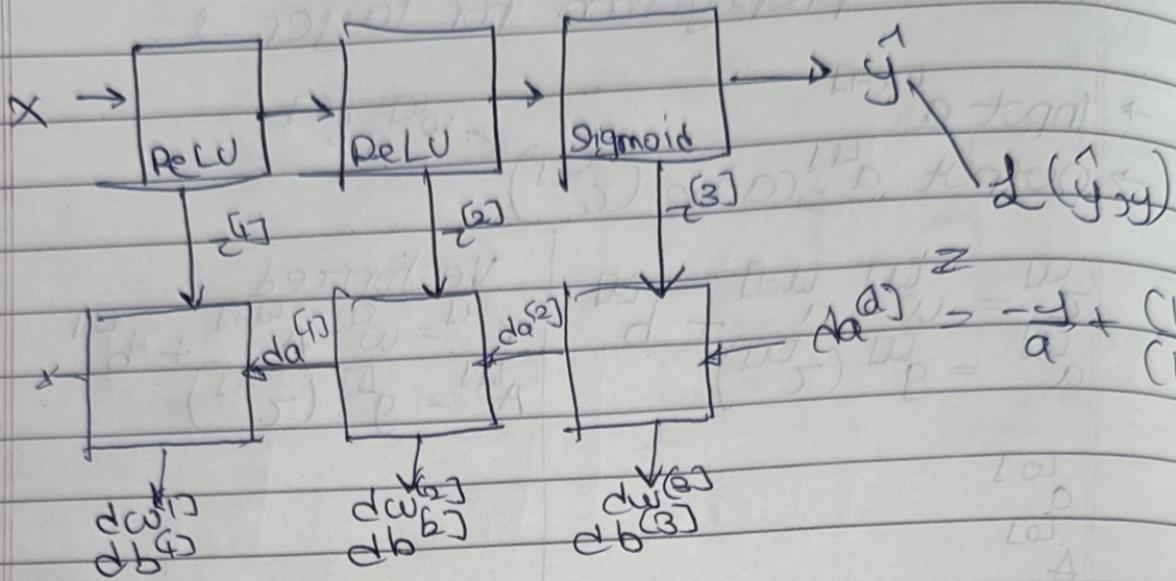
vectorized

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

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

$$db^{[l-1]} = \frac{1}{m} \sum w^{[l+1]T} \cdot dz^{[l+1]}, \text{ axis=1; keepdims=True}$$

Summary



$$dA^{(l)} = \left(-\frac{y^{(l)}}{a^{(l)}} + \frac{(1-y^{(l)})}{1-a^{(l)}} \right)$$

~~1 round of backpropagation~~

$$\frac{(D-1)}{D} * w_0 = b_0$$

$$\frac{D-1}{D} * b_0 = \frac{D}{D-1} * w_0$$

$$5b = db$$

$$\frac{D-1}{D} * b_0 * \frac{D}{D-1} * w_0 = db$$

$$\frac{(D-1)}{D} * w_0 * \frac{D}{D-1} * w_0 = db$$

herefore

$$\frac{(D-1)}{D} * w_0 * A_b = db$$

$$A * \frac{D-1}{D} * w_0 = db$$

m

$$(2)(7) - \min(99, 1) = 2 * D - Sh \quad m \leq 99 \quad l = db$$

PARAMETERS VS HYPERPARAMETERS

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

Hyperparameters:- Learning rate α