

10/17/24

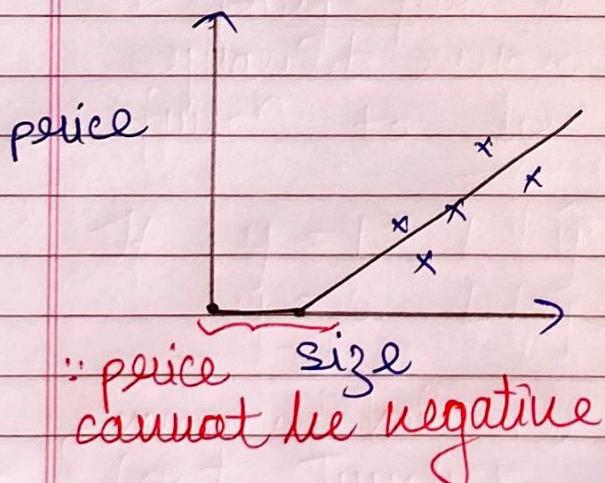
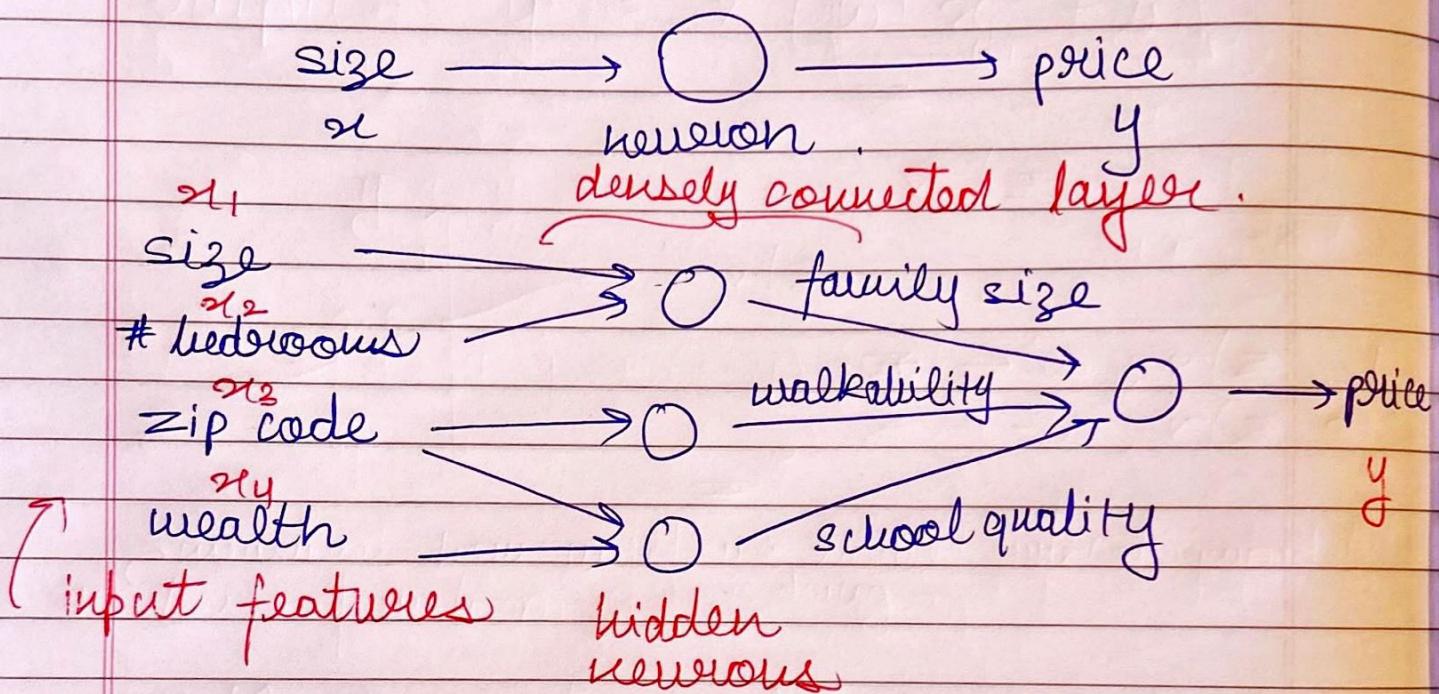
classmate

Date _____

Page _____

DEEP LEARNING

- * What is a neural network?
Stacking together many neurons gives us a neural network.



ReLU function
(Rectified Linear Unit)

- * Common Neural Networks

CNN
→ images

RNN
→ sequential data

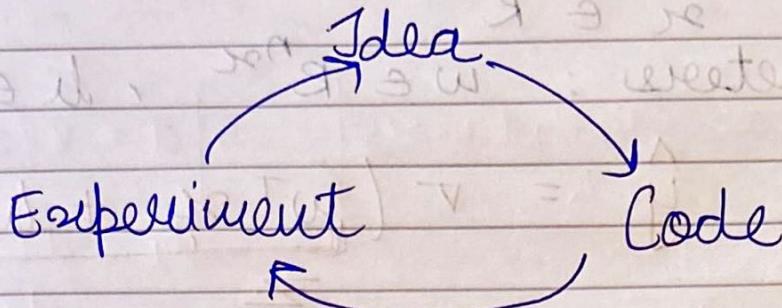
Hybrid / custom
→ advanced models.

Structured Data

Proper data set, every feature is defined

Unstructured Data

Images, audio, random texts.



as building an effective NN is an **iterative process**.

* BINARY CLASSIFICATION:

input image : 64×64

feature vector: matrice $X = \begin{bmatrix} 255 \\ 251 \\ \vdots \\ 255 \\ 251 \\ \vdots \end{bmatrix} = 64 \times 64$
 $n_x = 12288$ (here)
 \xrightarrow{X} image of cat 0/1
 $\xrightarrow{\text{no. of}}$

m = training examples

Basic Notations

$$\star x \in \mathbb{R}^{n_x \times m} \quad (\text{$n_x \times m$ dimensional matrice})$$

$$x = \left[\begin{array}{c|c|c|c}
1 & 1 & \dots & 1 \\
\hline x^{(1)} & x^{(2)} & \dots & x^{(m)}
\end{array} \right] \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \\ m \end{matrix}$$

$$\star x \in \mathbb{R}^{k \times m} = \left[y^{(1)} \ y^{(2)} \ \dots \ y^{(m)} \right]$$

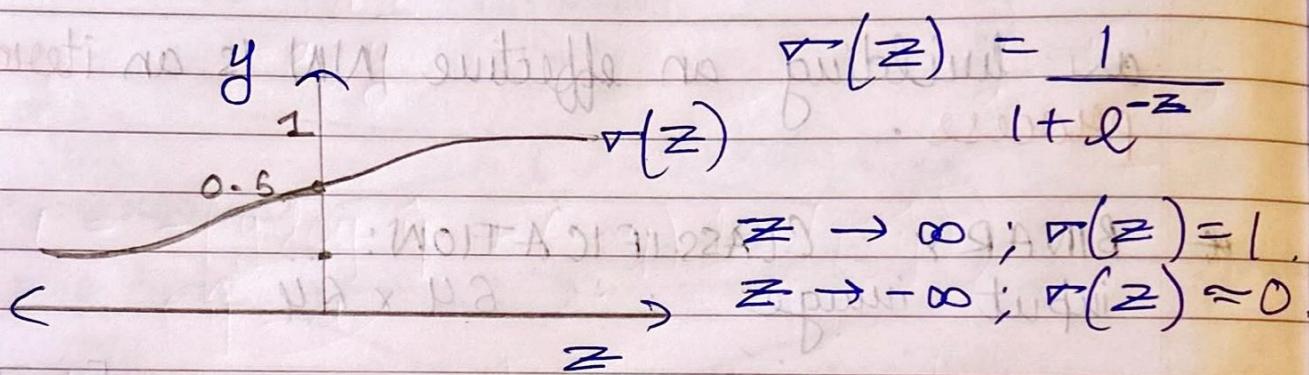
* LOGISTIC REGRESSION:

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

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

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

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



→ measures how well our model func' (sigmoid func') is doing

* Loss function:

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

→ opting for a sigmoid function over squared error here as the want the loss func' to be convex.

$$\text{If } y=1 \quad L(\hat{y}, y) = -\log \hat{y}$$

$$y=0 \quad L(\hat{y}, y) = -\log(1-\hat{y})$$

→ measures how well parameters w, b are doing

* Cost Function:

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

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

* GRADIENT DESCENT:

→ want to find w, b that minimizes $J(w, b)$

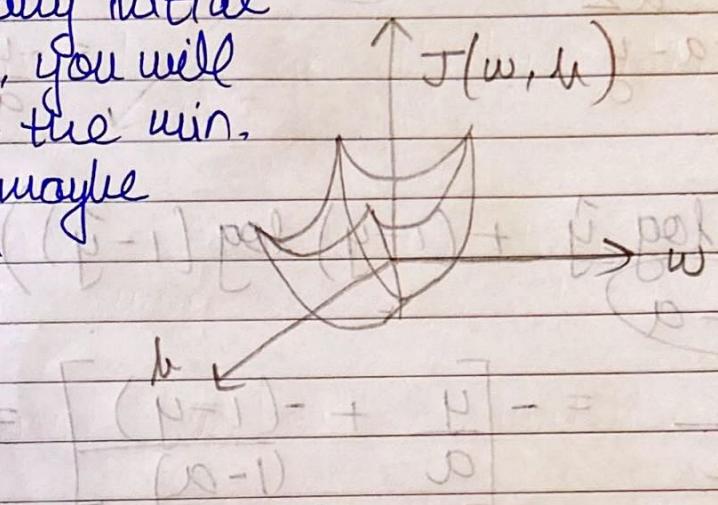
$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

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

With any initial value, you will reach the min. (path maybe diff.).

Why gradient descent?

→ it is a convex func. Thus, has one minimum only (global)



* COMPUTATION GRAPH:

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

$$\frac{da}{a} = 3$$

$$b \rightarrow [u = bc] \xrightarrow{du = 3} V = a + u$$

$$c \rightarrow$$

$$V = a + u \xrightarrow{dv = 3} [J = 3V]$$

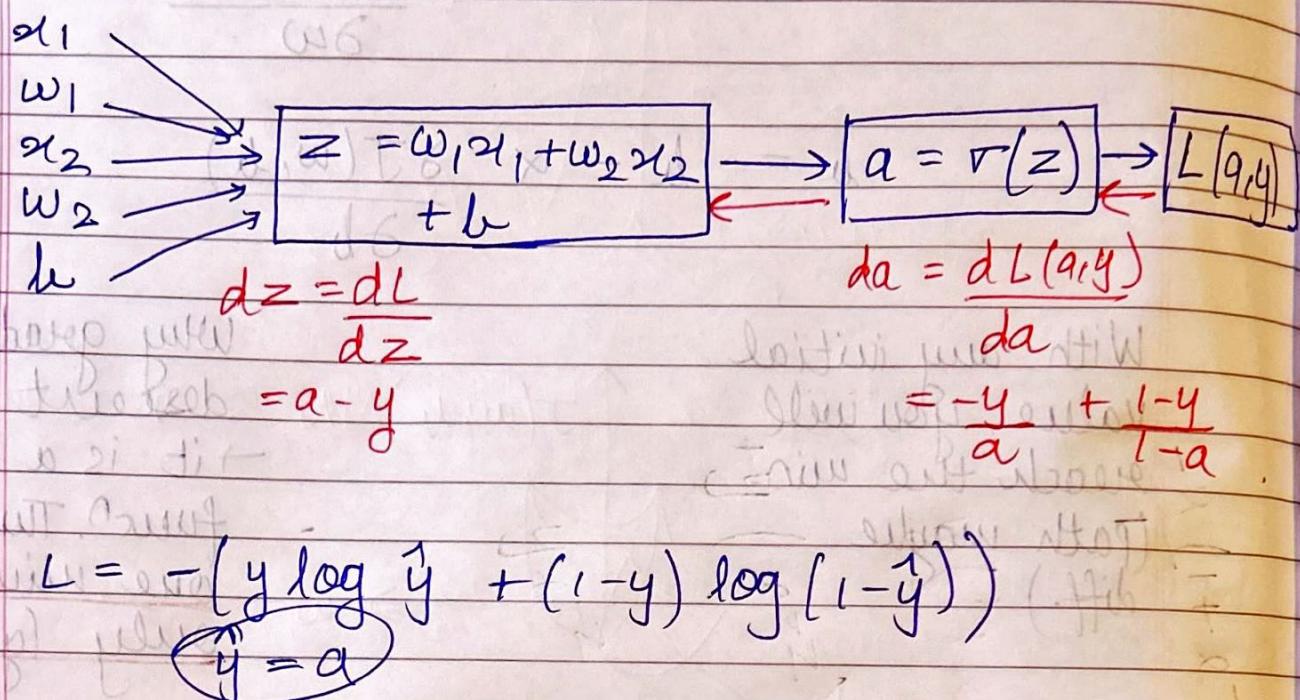
derivative

We observe, through a left to right pass we computed the value of J .

To compute derivatives, we move from right to left, i.e., in technical terms use backpropagation to calculate derivative.

* d_{var} in code $\equiv J$ (here).
in code, $\frac{dJ}{da} = da = 3$ (here)

* Gradient Descent for Logistic Regression (using computation graph)



$$\frac{\partial a}{\partial z} = \frac{\partial L}{\partial a} = -\left[\frac{y}{a} + \frac{-(1-y)}{(1-a)} \right] = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$\frac{\partial z}{\partial z} = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}$$

$$= \left(-\frac{y}{a} + \frac{(1-y)}{(1-a)} \right) \frac{\partial r(z)}{\partial z}$$

$$= \left(-\frac{y}{a} + \frac{(1-y)}{(1-a)} \right) \times \underbrace{\frac{\partial}{\partial z} \frac{1}{1+e^{-z}}}_{\frac{\partial}{\partial z} \frac{1}{1+e^{-z}}}$$

$$\frac{\partial a}{\partial z} = \frac{1}{(1+e^{-z})^2} (e^{-z})$$

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

$$\frac{\partial a}{\partial z} = a^2 \frac{(1-a)}{z} = a(1-a).$$

$$\therefore dz = \frac{d L}{d z} = \frac{(-y + a - ay) a (1-a)}{a(1-a)}$$

$$dz = \underline{a - y}$$

- * Logistic Regression on m examples
 - using for loop

for $i := 1 \rightarrow m$

$$z^{(i)} =$$

$$a^{(i)} =$$

$$J_t =$$

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

$$dw_1 + =$$

$$dw_2 + =$$

$$db + =$$

considering only
 w_1 & w_2)

$$J_t = m$$

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

use

(greater datasets are involved)

As we advance in DL, for loops can be time-consuming, thus we use **vectorization**.

- * VECTORIZATION

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

- * Neural Network Programming Guideline
 - whenever possible, ~~avoid explicit~~ for loops.
 - Use built in functions of Numpy & Python, makes the code more efficient.

* Vectorizing Logistic Regression :

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

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

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

↑ transpose of w

↑ column vector

Gradient Descent

Pseudo Code

we can

use a for loop here

if we want to implement gradient

descent multiple times.

$$A = r(z)$$

$$dz = A - y$$

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

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

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

$$b := b - \alpha dw$$

Broadcasting in Python.

* BROADCASTING IN PYTHON:

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

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

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

can use $+/-/*/ /$

General Principle: $\begin{array}{c} + \\ - \\ * \\ / \end{array} (1, n) \rightsquigarrow (m, n)$
 $(m, 1) \rightsquigarrow (m, n)$

* main form of programming in neural network
when broadcasting which is useful in

TIP: with python/numpy vectors don't

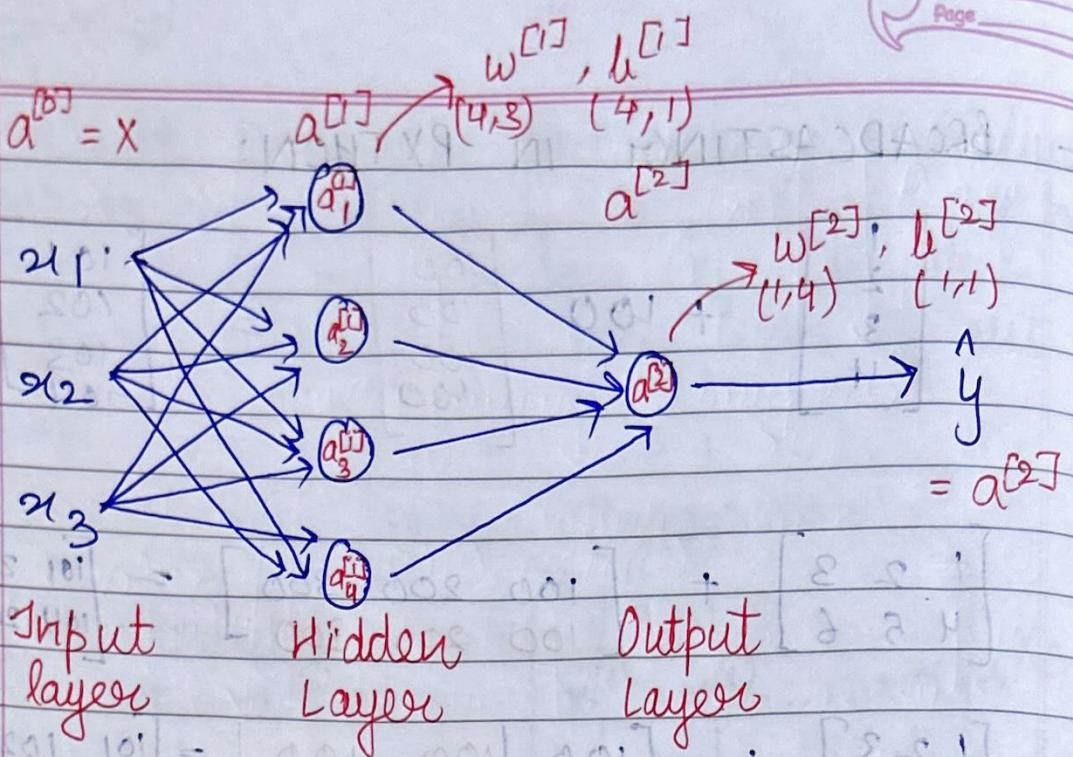
use 1D array, rather use:

$a = np.random.rand(5, 1)$

here $a.shape = (5, 1)$ & NOT (5)

* Neural Network Representation:

- The things in the hidden layer are not visible to the training set. (i.e., the value of y is not visible to the NN while computing y)



2 layer NN (considered conventionally, as input layer does not count as an official layer).

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

first parameter of first training sample

for each feature (x_1, x_2)
(x_3) the value of the parameters

$$(w_1, w_2, w_3)$$

w_4 is

different.

This value

of z is 4th parameter of

also different
1st + 2nd +
+ ... + 4th of with

Note:

$$z = w_1 x_1 + b_1$$

$$(z_n \times w_n \text{ not } z_l)$$

↳ 4th parameter of 2nd training eg

$$\begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

nth parameter of with training sample
 $\times m^{th} x_1$

$$[w_{11} \quad w_{12} \quad w_{13} \quad w_{14}]$$

classmate

Date _____

Page _____

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = r(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = r(z_2^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = r(z_4^{[1]})$$

$w^{[1]} \quad (4, 3)$
 (n_x, m)

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

$(m, 1)$
 $(n_x, 1)$

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

$b^{[1]} \quad (4, 1)$
 $(n_x, 1)$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = r(z^{[1]})$$

* Vectorizing across multiple examples:

~~$$a^{[1](i)} = w^{[1]} \cdot x^{(i)} + b^{[1]}$$

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

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

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

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

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

Forward
Propagation

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

n features

m training samples

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

nth value of corresponding m

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

n₂ hidden units

iterating all training samples (m) horizontally

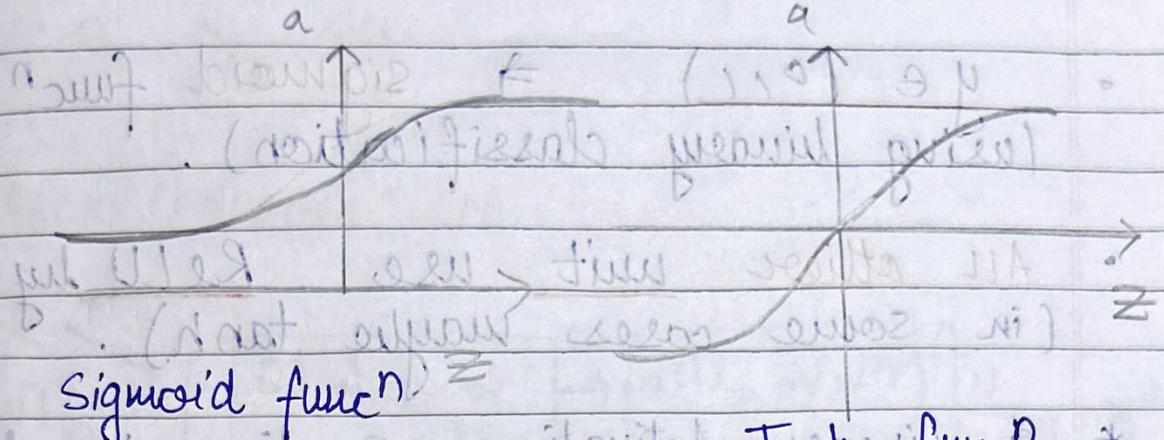
Eg: $a_1^{[i](1)}$ → first hidden unit of the first training example

$a_2^{[i](1)}$ → second hidden unit of the first training example.

$a_1^{[i](2)}$ → first hidden unit of the second training example.

$\star \left\{ a_n^{i} \right\}$ → nth hidden unit of the ith training example.

* ACTIVATION FUNCTIONS :

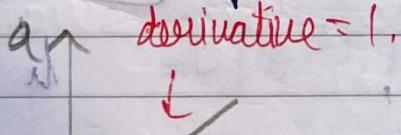


$$\text{Sigmoid funcn} = a = \frac{1}{1 + e^{-z}}$$

$$\text{Tanh funcn} = a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$\tanh > \text{sigmoid}$ [except] for when
 $y \in (0, 1)$.

- The activation functions of different layers can be different.
- Disadvantage : if $y \rightarrow \infty$ or $z \rightarrow -\infty$, the gradient derivative becomes very slow, thus affecting gradient descent.
- Thus another function that's used is : ReLU (Rectified Linear Unit) funcn :



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

thus derivative is 0 or 1.

Preferably, leaky ReLU > ReLU
 $(\because z \text{ is } -ve)$

leaky ReLU $a = \max(0, 0.01z)$

* Rules to select Activation Functions:

- $y \in (0, 1)$ \Rightarrow sigmoid funcⁿ
(using binary classification).
- All other unit use ReLU by default
(in some cases maybe tanh).

* A Linear Activation Function is almost futile in a hidden NN, as we're proceeding in the NN but not performing any complex funcⁿ.

* Given x :

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

$$a^{[1]} = z^{[1]}$$

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

$$a^{[2]} = z^{[2]}$$

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

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

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

$$(w, b) \text{ known} = P$$

$$\therefore o = W' x + b'$$

The only time we can use a linear activation funcⁿ is in the output layer.
Eg: in the housing price prediction eg.

$n \times 1$

$1 \times n$

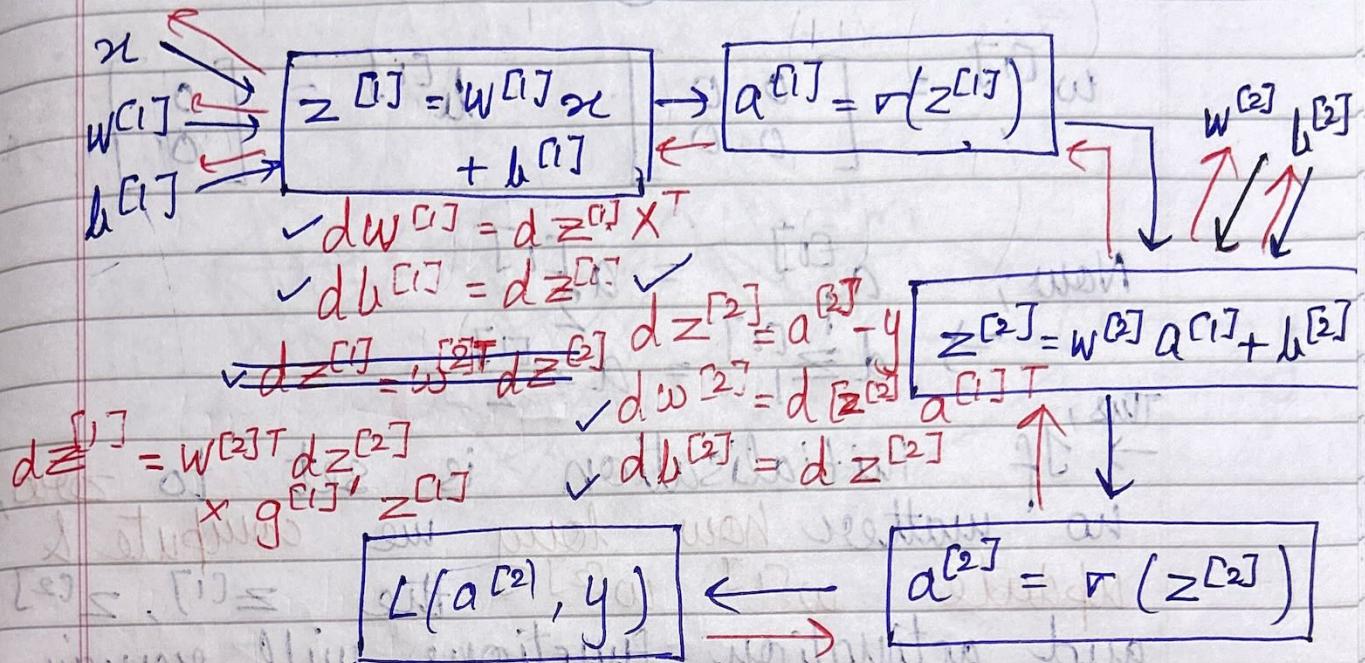
classmate

Date _____

Page _____

- Non-Linear Activation funcn is a critical part of NN.

* BACK PROPAGATION



Dimensions:
 n_x : input features
 $n^{[0]}$: hidden features
 $n^{[1]}$: $n^{[2]}$: activation funcⁿ parameters.

$$w^{[2]} : (n^{[2]}, n^{[1]})$$

$$z^{[2]}, dz^{[2]} : (n^{[2]}, 1)$$

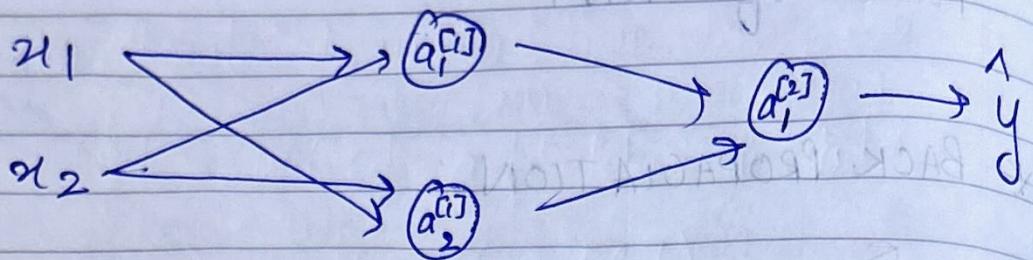
$$z^{[1]}, dz^{[1]} : (n^{[1]}, 1)$$

$$\begin{aligned} dz^{[1]} &= w^{[2]T} * g^{[1]'}(z^{[1]}) \\ &= (n^{[2]}, n^{[2]}) * (n^{[2]}, 1) * (n^{[1]}, 1) \end{aligned}$$

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

$$\therefore dz^{[1]} = (n^{[1]}, 1)$$

* INITIALISATION:



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

Now, $a_1^{[1]} = a_2^{[1]}$
 $\Rightarrow d^{[1]}_1 = d^{[1]}_2$

Thus,

→ If initialisation is set to zero, no matter how long we compute & update $w^{[1]}, w^{[2]}$, the $z^{[1]}, z^{[2]}$ and activation functions will remain symmetric (give same output).

As long w is not initialised to 0, we can initialise $b^{[1]}$ to 0, as it does not have symmetry issue.

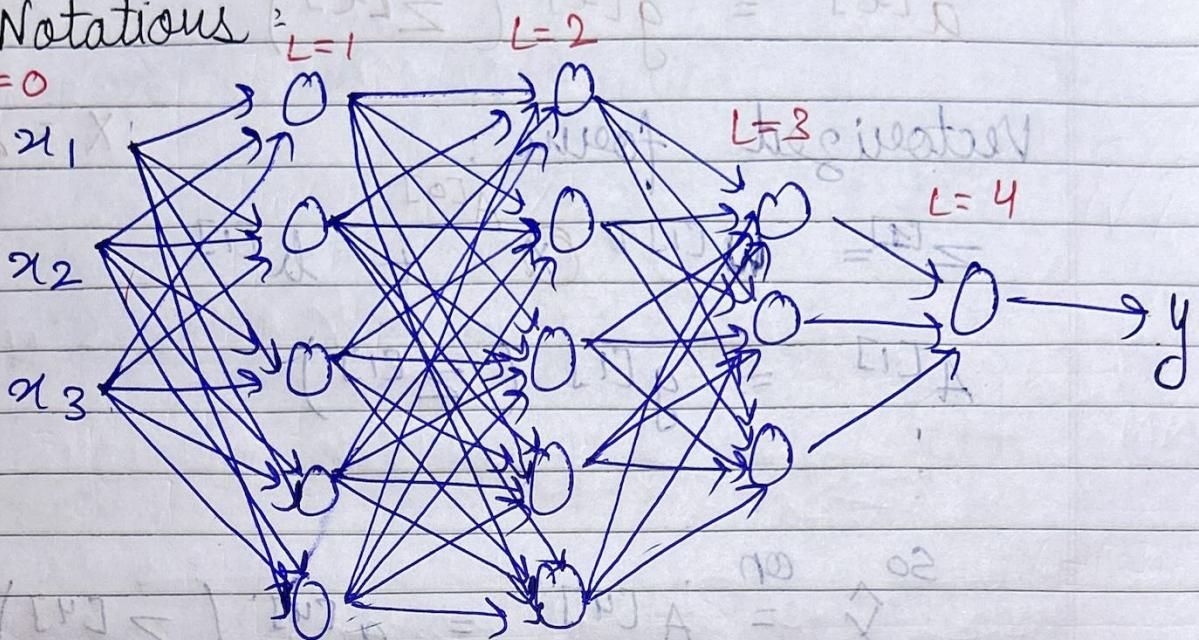
We generally randomly initialise the parameters to very small numbers

* DEEP NEURAL NETWORK.

o hidden layer; 2 layer NN, 3 layer NN, 5 ^{layer}_{NN}
shallow NN **deep NN**

Eg: Logistic Regression.

* Notations



$L = 4$ (layers)

$n^{[l]}$ = units in layer l .

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

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

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

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

* Forward Propagation in DNN:

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

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

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

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

and so on,

Generalised formula :

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

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

Vectorized form :

$$\Sigma^{[1]} = \omega^{[1]T} A^{[0]} +$$

$$A^{[1]} = g^{[1]} (\geq^{[1]})$$

so on

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

For computing the activation of layers from 1 to n^{th} layer we need to use an explicit for loop.

* Matrix Dimensions

	1	2	3	4	5
211	0	0	0	0	0
212	0	0	0	0	0
	0	0	0	0	0

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

$(3,1) \quad (3,2) \quad (2,1)$ ignore (for now)

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

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

$$\underline{z^{[2]}} \Rightarrow (n^{[2]}, 1)$$

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

$(5,1) \quad (5,3) \quad (3,1)$

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

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

$$db^{[l]}, b^{[l]} \text{ equal } (n^{[l]}, 1)$$

On implementing vectorization:

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

* Intuition on Deep Representation:

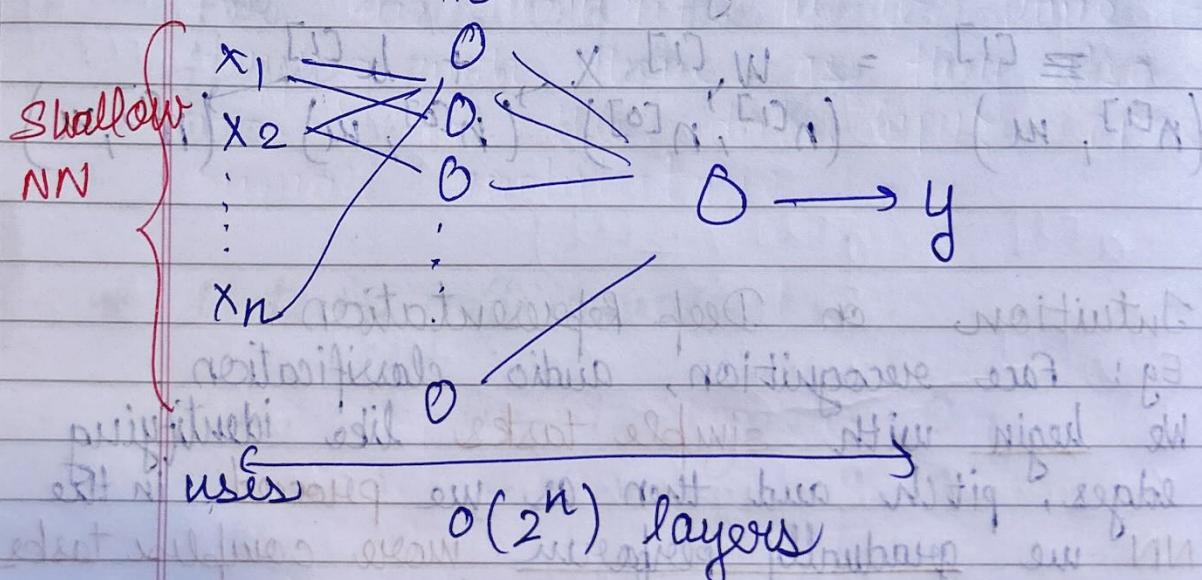
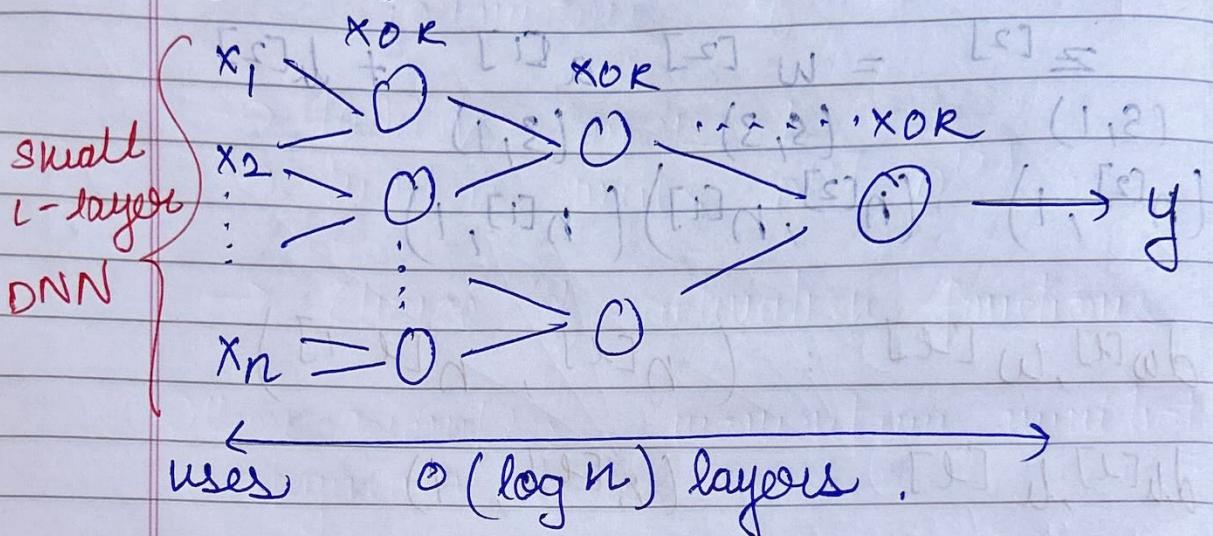
Eg: Face recognition, audio classification.

We begin with simple tasks like identifying edges, pitch and then as we proceed in the NN we gradually perform more complex tasks to determine the output

* Circuit Theory:

There are functions you can compute with 'small' L-layer DNN that shallower networks require exponentially more hidden units to compute.

$$\text{Eg: } y = (x_1 \text{ XOR } x_2) \text{ XOR } x_3 \dots \dots \text{ XOR } x_n$$



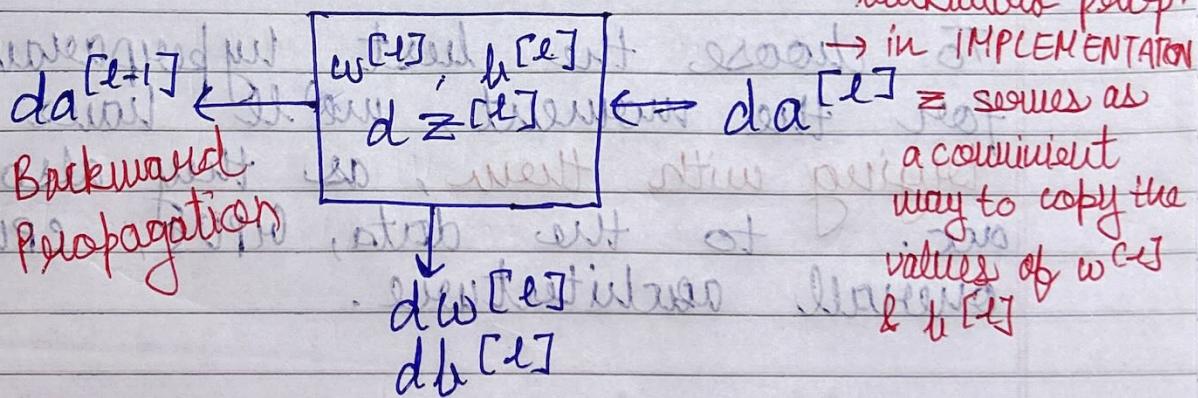
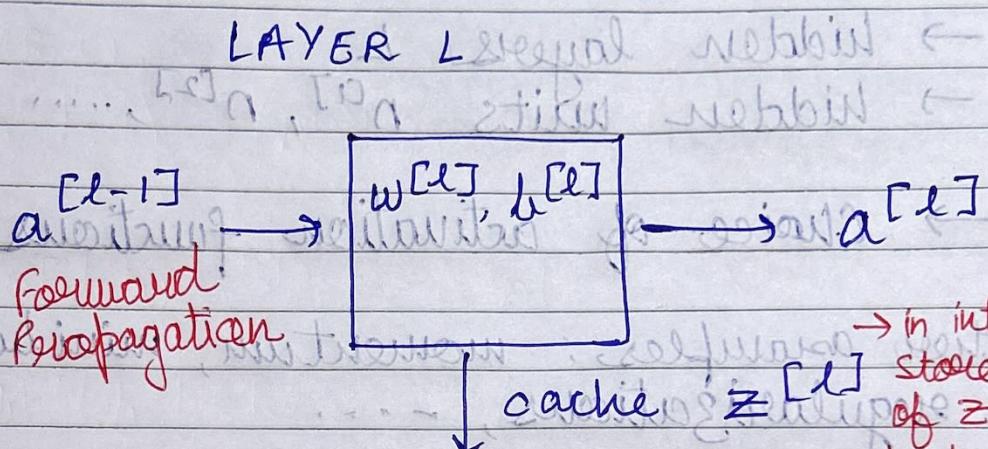
Thus, using deep neural networks really simplifies the math & efficiency of your model.

* Forward & Backward Func:

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

Forward: Input $a^{[l-1]}$, output $a^{[l]}$
cache $z^{[l]}$

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



* Backward Prop.:

Input: $da^{[l]}$

Output: $da^{[l-1]}$, $d w^{[l]}$, $db^{[l]}$
 $d z^{[l]} = da^{[l]} * g^{[l]}(z^{[l]})$ from cache.

$$d w^{[l]} = d z^{[l]} * a^{[l-1]}$$

$$d b^{[l]} = d z^{[l]}$$

$$d z^{[l]} = w^{[l+1]T} d z^{[l+1]} * g^{[l]}(z^{[l]})$$

$$da^{[l-1]} = w^{[l]T} d z^{[l]}$$

* What are hyperparameters?

- actual parameters that control the model
- Parameters : $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$
- hyper parameters which are key in the model & impact the parameters

→ learning rate α : learning

iterations

parameters of $w^{[1]}$

→ hidden layers

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

→ Choice of activation functions.

Other examples: momentum, mini batch size, regularizations, ...

To choose the best hyperparameters for the moment, we'll have to keep trying with them, as they also change acc. to the data, GPUs, CPUs & overall architecture.

After course :

- implement logistic Reg
- activation funtions
- gradient descent
- backprop

* Implementing Logistic Regression:

Input features:

x_1, x_2, x_3, x_4 .

loan - ID, gender, married, dependents, education, self-employed, applicant-income, co-applicant-income, loanAmount, loanAmount term, Credit History, Property Area, ~~loan~~

y: Loan Status.

→ id

→ income

→ amount → term