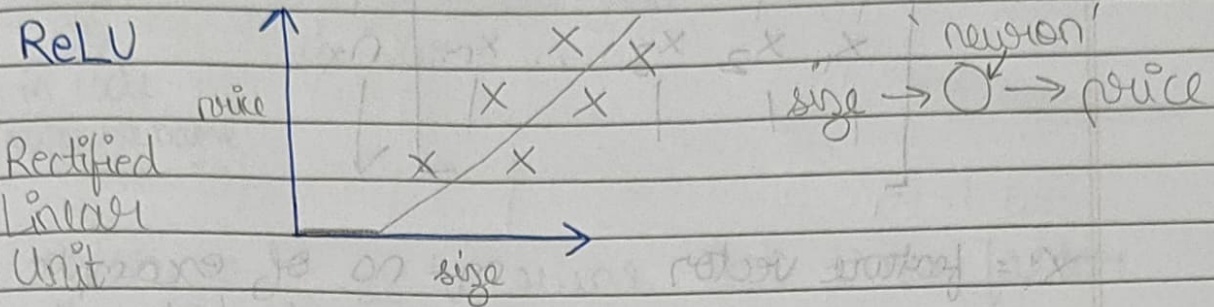
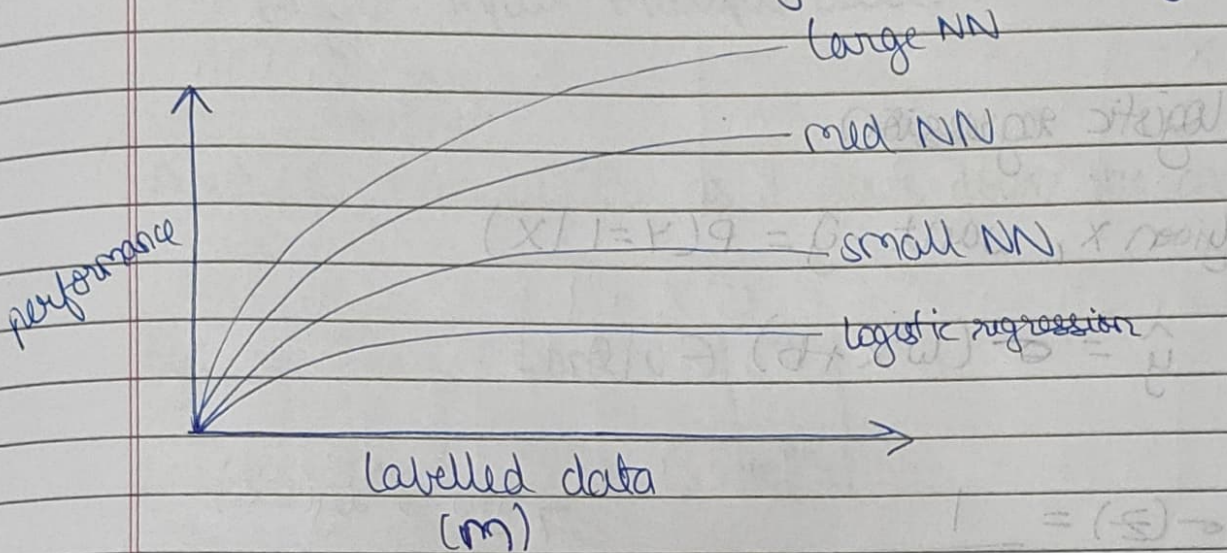


Course I



Supervised Learning

Input, Output pair given in training data.



ReLU makes learning faster than sigmoid

∵ the gradient in sigmoid tend to 0 hence lowering rate of learning

Hence, ReLU \gg sigmoid as activation function of NN

$$X = \begin{bmatrix} | & | & | & \dots & | \\ x_1 & x_2 & x_3 & \dots & x_m \\ | & | & | & \dots & | \end{bmatrix} \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

$\xleftarrow{\quad m \quad}$

x_i = feature vector

m = no of example
 n = no of features

$$Y = [y_1 \ y_2 \ \dots \ y_m] \updownarrow 1 \quad y_i = \{0, 1\}$$

$\xleftarrow{\quad m \quad}$

★ logistic regression

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

$$\hat{y} = \sigma(w^T x + b) \in [0, 1]$$

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

$y=0 \rightarrow \log(1-\hat{y})$
 $y=1 \rightarrow \log(\hat{y})$

$$\text{loss} = y \log \hat{y} + (1-y) \log(1-\hat{y}) \quad (\text{individual case})$$

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

avg total loss

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

$$\text{dvar} \rightarrow \text{d}(\text{Final Output Var})$$

$$\Downarrow$$

$$\text{d}(\text{Var})$$
 in code
variable name

vectorize as it runs 300 times faster than non vectorize code.

SIMD: parallel computation.

↳ single input multiple data.

AVOID FOR LOOPS as it slows down the code

Eg: $v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$ want $u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$

Non vector

Vector

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

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

for i in range(n):

→ this explicit for is removed and code is made simpler and faster.

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

explore the numpy built in function ~~lib~~ lib.

$$dw = \text{np.zeros}((n_x, 1))$$

$$dw += x^i dz^i$$

$$dw[1] = m$$

these changes help
remove one of the
explicit for loop.

$$\star z^1 = w^T x^1 + b$$

$$z^2 = w^T x^2 + b$$

$$z^3 = w^T x^3 + b$$

$$a^1 = \sigma(z^1)$$

$$a^2 = \sigma(z^2)$$

$$a^3 = \sigma(z^3)$$

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & \dots & x_m \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$X_{n_x \times m}$ matrix

$$Z = \begin{bmatrix} z^1 & z^2 & \dots & z^m \end{bmatrix} = w^T X + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

$\xleftarrow{\quad m \quad} \quad \quad \quad \xleftarrow{\quad m \quad}$

$$= \begin{bmatrix} w^T x^1 & w^T x^2 & \dots & w^T x^m \end{bmatrix}$$

$$+ \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

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

$\xleftarrow{\quad 1 \times m \quad} \quad \quad \quad \xrightarrow{\quad}$

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

real number

python converts
broadcasting to it to $1 \times m$ matrix

$$A = [a^1 \ a^2 \ \dots \ a^m] = \sigma(Z) \rightarrow \text{in assignment}$$

★ Vectorizing Logistic Regression $(S) - A = B$

$$dz^i = a^i - y^i$$

$$Y - A = S b$$

$$dZ = [dz^1 \ dz^2 \ \dots \ dz^m] \quad (Sb) \text{ } 1 \times m \text{ matrix}$$

$$A = [a^1 \ a^2 \ \dots \ a^m]$$

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

$$dZ = A - Y = [a^1 - y^1 \ a^2 - y^2 \ \dots \ a^m - y^m]$$

1 step mat
calculated.

$$\text{np.sum}(dZ) = \text{sum of all elements.}$$

Using vector for all operations, removes the need for explicit for loops

Also, it make the code simpler, increase readability and faster and efficient.

★ Vectorized implementation.

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

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dW = \frac{1}{m} X dZ^T$$

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

put this in

a for loop

to calculate

multiple iteration

Broadcasting.

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

$$(m, n) + \frac{(1, n)}{x} = (m, n) + (m, n) = (m, n)$$

copy row 1, m times

$a = \text{np.random.randn}(5)$

$a.\text{shape} = (5,) \rightarrow$ rank 1 vector

} X do not use this

$a = \text{np.random.randn}(5, 1) \rightarrow$ USE THIS

$\text{assert}(a.\text{shape} == (5, 1))$

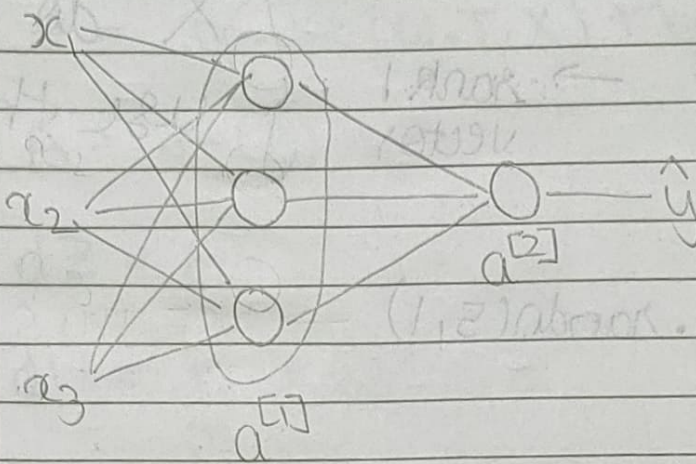
$$\begin{aligned} \# \log P(\text{label in trainingset}) &= \log \prod_{i=1}^M P(y^i | x^i) \\ &= \sum_{i=1}^M \log P(y^i | x^i) \\ &= - \sum_{i=1}^M L(\hat{y}^i, y^i) \end{aligned}$$

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

The training set is assumed to be i.i.d

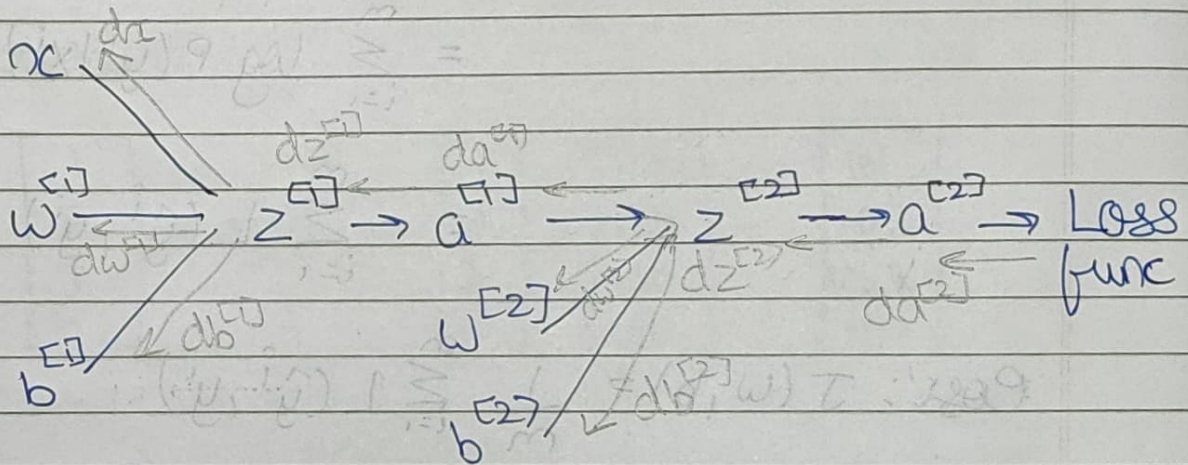
↓
independent and identically distributed.

Overview



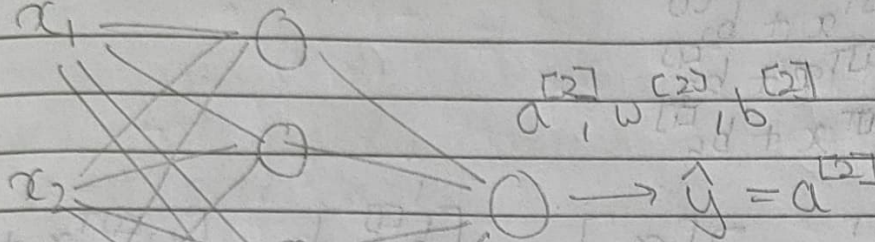
$[1] \rightarrow$ layer 1
 $[2] \rightarrow$ layer 2

(1) \rightarrow example 1
 (2) \rightarrow example 2



★ NN Representation.

$$a^{[0]} = X \quad a^{[1]}, w^{[1]}, b^{[1]}$$

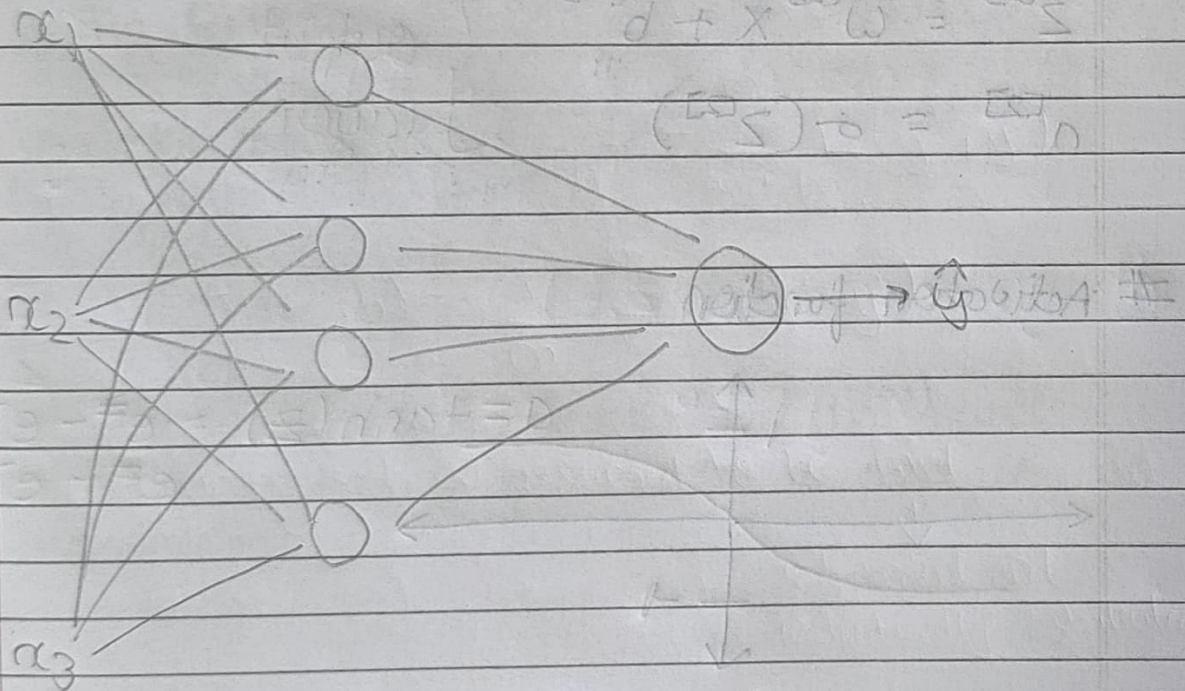
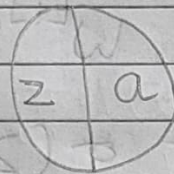


Output

2 layer NN

input

hidden

values
are not
seen

$$z_i^{[l]} = w_i^{[l]} x + b_i^{[l]}$$

$$a_i^{[l]} = \sigma(z_i^{[l]})$$

l=layer i=number

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

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

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

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

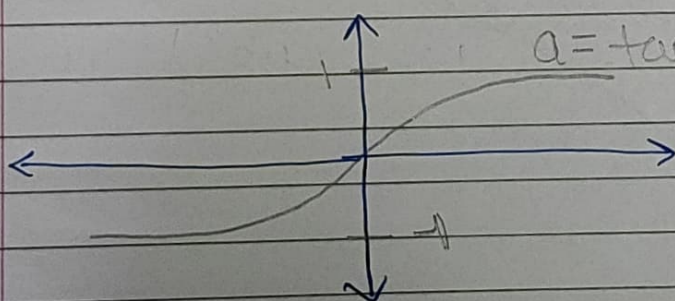
} hidden layer

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

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

} output layer

Activation function



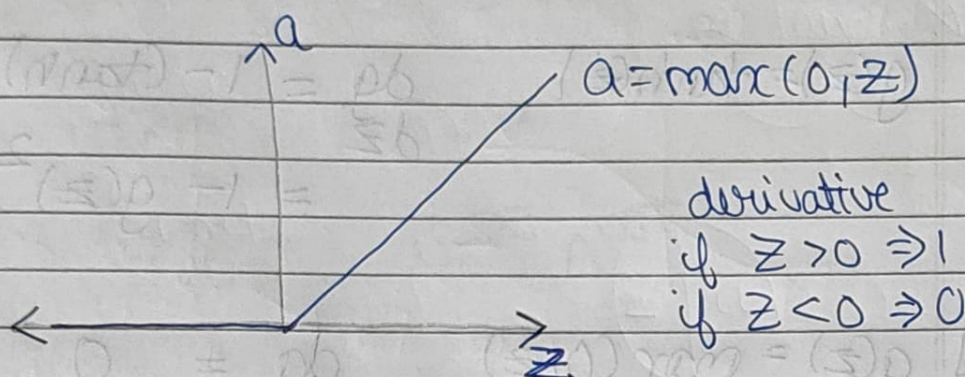
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$\tanh(z)$ works better than $\sigma(z)$

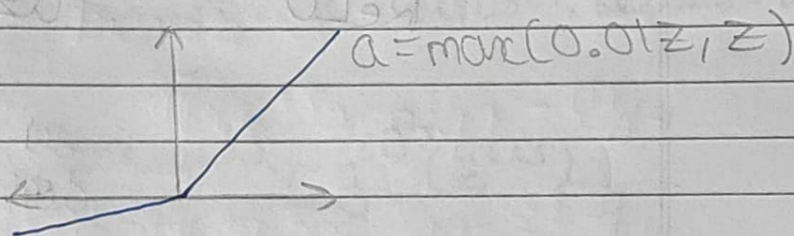
Sigmoid for binary classification for output layer

tanh gives derivative as 1 for large values and hence better than sigmoid but still not good enough

∴ ReLU = $\max(z, 0)$ is used



Leaky ReLU



better than ReLU
but doesn't
make huge difference

hidden layer mai linear use nhi krte
because ulta output linear hi dega so its
redundant

★ Derivative of activation function.

$$\begin{aligned} \text{a)} \quad g(z) &= \frac{1}{1+e^{-z}} & \frac{dg}{dz} &= \frac{+e^{-z}}{(1+e^{-z})^2} \\ &= \sigma(z) & &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ & & &= g(z)(1-g(z)) \end{aligned}$$

$$\begin{aligned} \text{b)} \quad g(z) &= \tanh(z) & \frac{dg}{dz} &= 1 - (\tanh)^2 \\ & & &= 1 - g(z)^2 \end{aligned}$$

$$\begin{aligned} \text{c)} \quad g(z) &= \max(0, z) & \frac{dg}{dz} &= \begin{matrix} 0 & z < 0 \\ 1 & z > 0 \end{matrix} \\ &\text{ReLU} & & \text{undefined } z = 0 \end{aligned}$$

← undefined
set to 1 in code
doesn't affect anything

$$\begin{aligned} \text{d)} \quad g(z) &= \max(0.01z, z) & \frac{dg}{dz} &= \begin{matrix} 1 & z \geq 0 \\ 0.01 & z < 0 \end{matrix} \\ &\text{leaky ReLU} & & \end{aligned}$$

★ gradient descent for neural networks. ★

Forward : $z^{[1]} = w^{[1]}x + b^{[1]}$
 $A^{[1]} = g^{[1]}(z^{[1]})$

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

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

Back : $dz^{[2]} = A^{[2]} - Y$

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

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

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

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

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

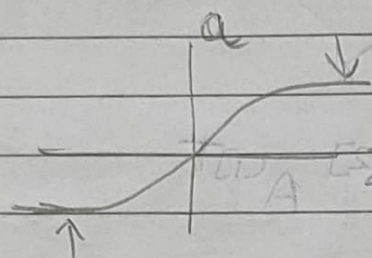
★ Random initialization is better than 0

bez all nodes will perform ~~same~~ same calc and doesn't make sense.

So random initialization.

Eg:

$$W^{[1]} = \text{np.random.randn}(2, 2) * 0.01$$



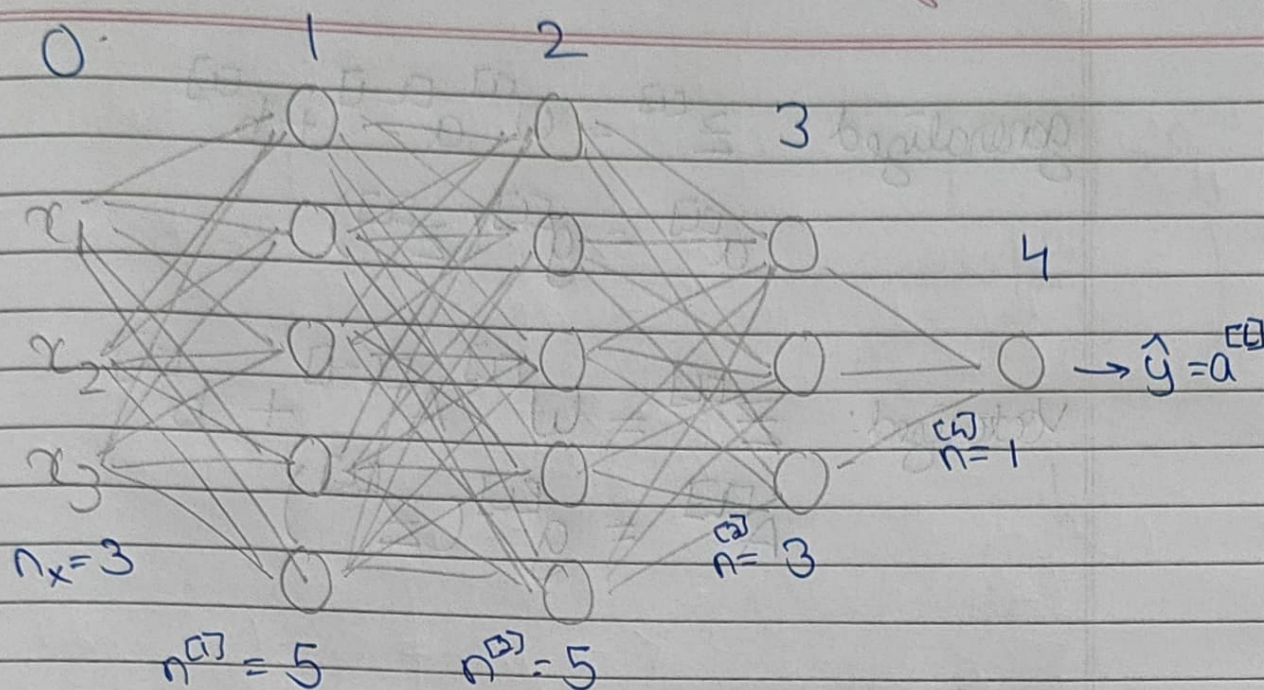
issue when
using tanh
or sigmoid

to minimize the
parameter s bez
if not done then
learning slows as
the gradient reaches 1

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

$$W^{[2]} = \text{np.random.randn}(1, 2) * 0.01$$

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



$L=4$ (#layers).

n^L = nodes in layer L

a^L = activation function = $g(z^L)$

$$x = a^0$$

★ Forward Propagation for example x

$$x: z^1 = w^1 x + b^1 = w^1 a^0 + b^1$$

$$a^1 = g^1(z^1)$$

$$z^2 = w^2 a^1 + b^2$$

$$a^2 = g^2(z^2)$$

$$\vdots$$

$$z^4 = w^4 a^3 + b^4, a^4 = g^4(z^4) = \hat{y}$$

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

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

vectorized: $Z^{[L]} = W^{[L]} A^{[L-1]} + b^{[L]}$

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

★ Matrix Dimensions

$$W^{[L]} = (n^{[L]}, n^{[L-1]})$$

$$dZ^{[L]} = Z^{[L]} = (n^{[L]}, m)$$

$$dA^{[L]} = a^{[L]} = (n^{[L]}, m)$$

$$b^{[L]} = (n^{[L]}, 1) \xrightarrow{\text{broadcasting}} (n^{[L]}, m)$$

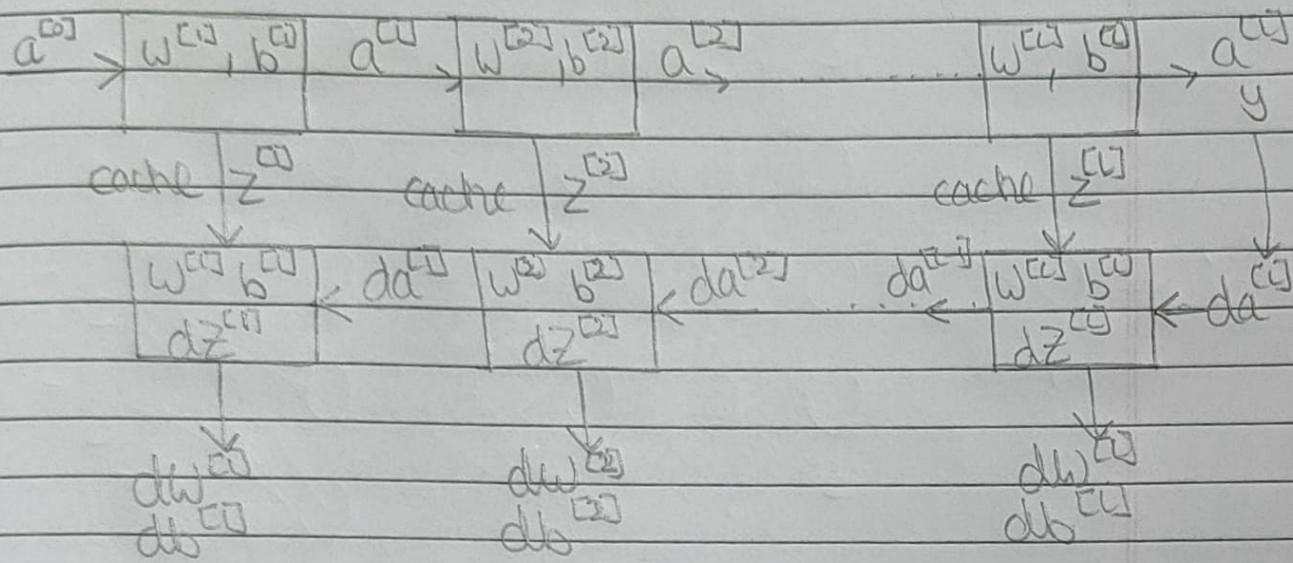
$$dW^{[L]} = (n^{[L]}, n^{[L-1]})$$

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

generalized
form

vectorized
implementation

- ★ Deep learning starts with simpler function and slowly making it more and more complex and hence deeper NN start performing better and better than the shallower ones.



$$W^{[l]} := W^{[l]} - \alpha dw^{[l]}$$

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

★ Hyperparameters :

α (learning rate)

no of iterations

no of hidden layers

no of nodes in a layer

activation functions