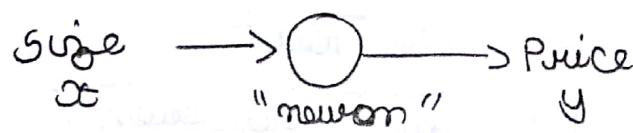
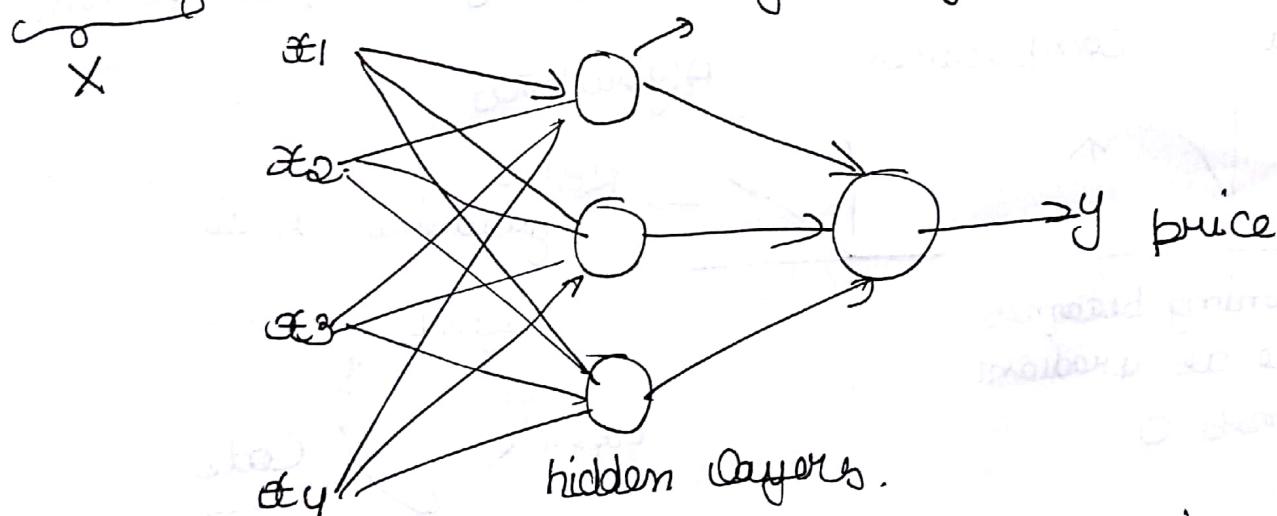
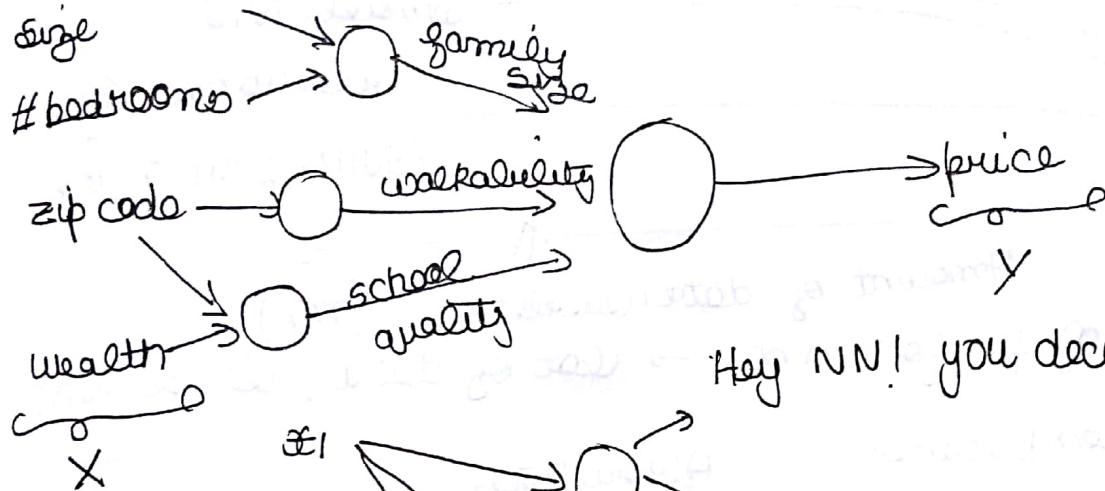


Deep Learning. ai

- AI is new electricity
- Deep learning highly sought after.
- Training Neural network.



$\text{ReLU} \rightarrow \text{Rectified Linear Unit}$.



Given (x, y) NNs can do well in mapping $x \rightarrow y$.

Supervised learning

Input (x)

standard Home features
NN of fd, user info

CNN of Image

RNN of Audio
(sequence)

Output (y)

Price

click on ad?

object (1..1000)

Text transcript

Application

Real Estate

Online Advert

Photo tagging

Speech

Structured Data

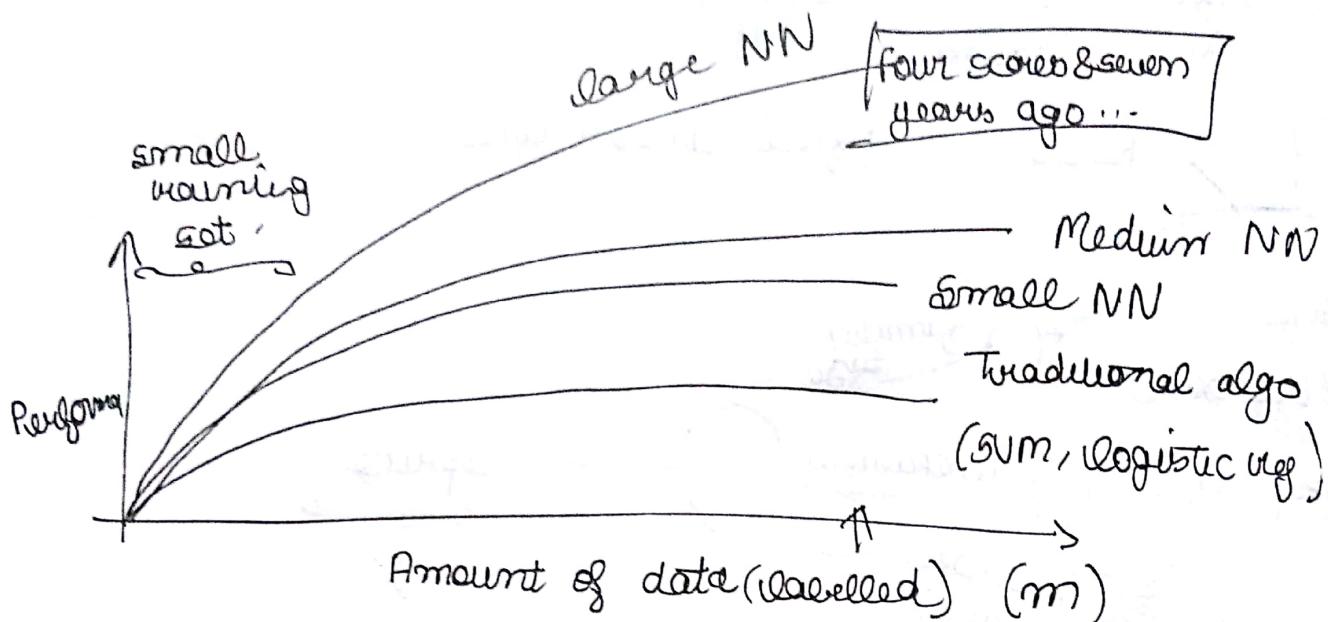
Size	#bedrooms	Price
small	1	1000
medium	2	2000
large	3	3000
huge	4	4000

Unstructured Data

Audio

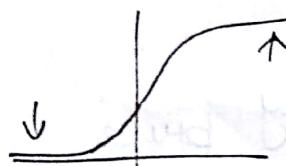
Image \rightarrow pixel values

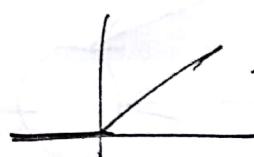
Treat



High level of performance \rightarrow lot of data, large NN.

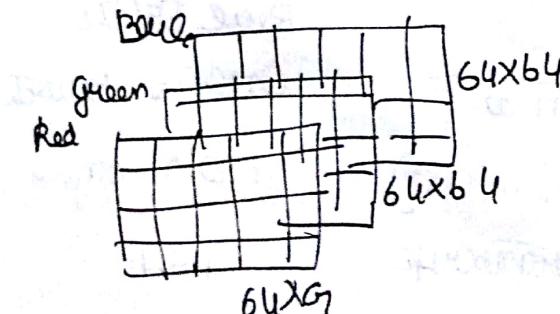
Data computation Algorithms


Learning becomes slow as gradient becomes 0,

 \rightarrow RELU gradient is 1.
Idea
Explore
Code

Week - 2

Logistic Regression - binary classification
 \rightarrow No iter loop. (1 (cat) v.s 0 (non cat))



$$\rightarrow X = \begin{bmatrix} 0.55 \\ 0.31 \\ \vdots \\ 0.56 \\ 0.34 \\ \vdots \end{bmatrix} \rightarrow 64 \times 64 \times 3$$

$$n_x = 10288$$

Notation

$(\mathbf{x}, y) \quad \mathbf{x} \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

m training examples: $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

$m = m_{\text{train}}$ $m_{\text{test}} \rightarrow$ test examples

$$\mathbf{X} = \begin{bmatrix} | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & | \end{bmatrix} \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

$\mathbf{x}^{(1)^T} \quad \mathbf{x}^{(2)^T} \quad \dots \quad \mathbf{x}^{(m)^T}$
difficult

$\mathbf{X} \cdot \text{shape} = (n_x, m)$

$$\mathbf{y} = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}] \quad \mathbf{y} \in \mathbb{R}^{1 \times m}$$

$\mathbf{X} \cdot \text{shape} = (1, m)$

Logistic Regression chance to tell it is a cat picture

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

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

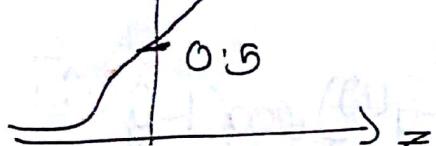
Parameters: $\mathbf{w} \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

$$\text{Output } \hat{y} = \underbrace{\mathbf{w}^T \mathbf{x} + b}_{\text{Linear regression}} = \sigma(\mathbf{w}^T \mathbf{x} + b) \quad \begin{matrix} \text{Logistic} \\ \text{Regression} \end{matrix}$$

but $(0 \& 1) \rightarrow$ we want

$$1 \rightarrow \sigma(z)$$

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



$$\text{if } z \text{ is very very big } \sigma(z) = \frac{1}{1+0} = 1$$

$$\text{If } z \text{ is very less } \sigma(z) = \frac{1}{1+\text{big number}} = 0$$

Learn w & b so that \hat{y} becomes 1.

w & b are separate.

In some conventions $x_0 = 1$, $\mathbf{x} \in \mathbb{R}^{n_{\text{in}}+1}$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_{n_{\text{in}}} \end{bmatrix} \quad \begin{array}{l} \text{more leditnes} \\ \text{y} \end{array}$$

Cost Function

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

Given $\{(\mathbf{x}^{(1)}, y^{(1)}) \dots (\mathbf{x}^{(m)}, y^{(m)})\}$, $\hat{y}^{(i)}$ vs $y^{(i)}$.

Loss (error function):

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \rightarrow \text{results in non convex. } \curvearrowleft \text{(gradient descent gets stuck at local)}$$

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

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

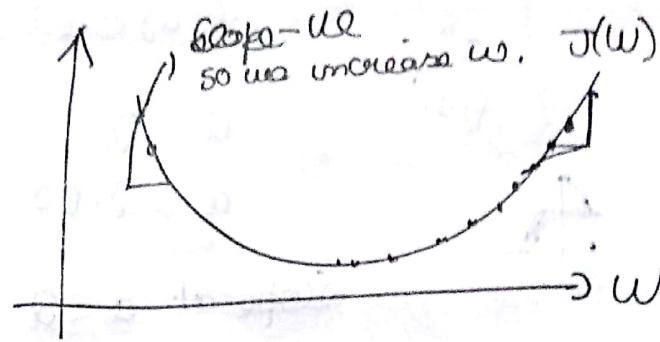
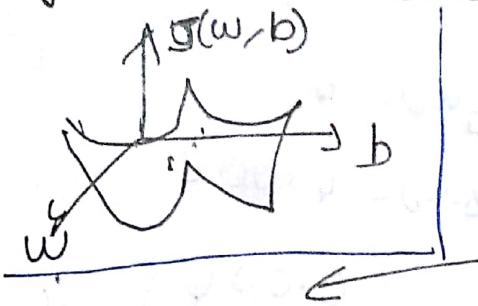
If $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{want } \log(1-\hat{y}) \text{ large, } \hat{y} \text{ small.}$

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



Repeat of

$$w := w - \alpha \frac{dJ(w)}{dw} \quad (\text{derivative is positive})$$

↑ $\frac{dJ(w)}{dw}$ → slope

g Learning Rate (control step) (want to find w, b that minimize $J(w, b)$)

→ Convex function: not

Initializing: w, b to zero (that's what ppe do instead of random initialisation)

Repeat of

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

g

Sometimes represented as

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

$$J(w, b)$$

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

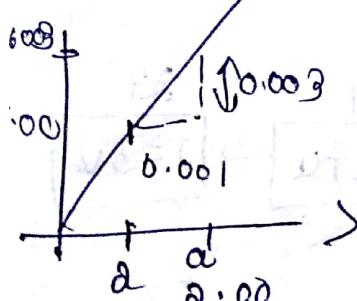
$$b := b - \alpha \frac{dJ(w, b)}{db}$$

Derivatives

$$f(a) = 3a \quad g(a) = 6 \quad a = 0.001 \quad g(a) = 6.003$$

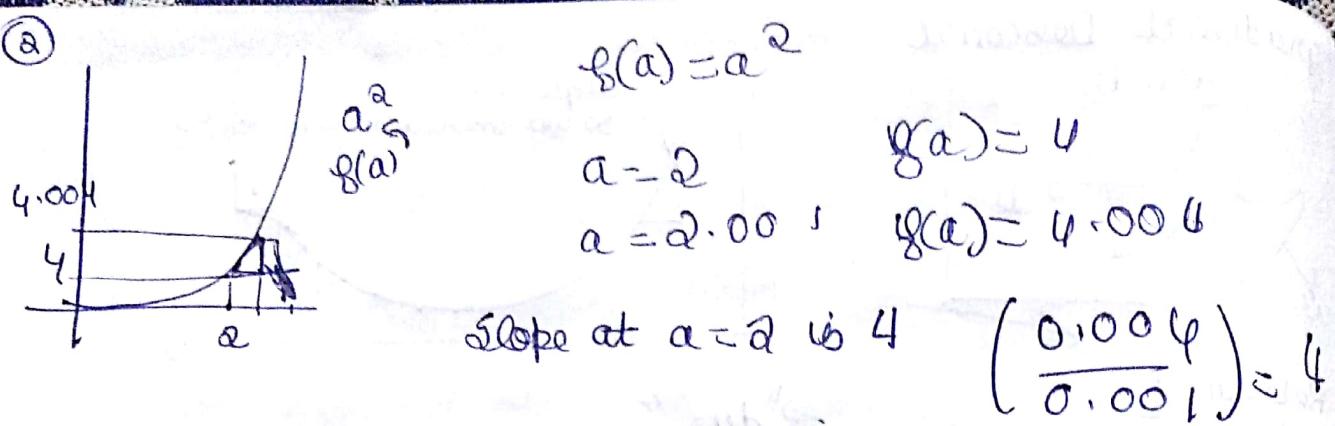
$$f'(a) = 3a$$

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



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

(At any point of function slope is same)



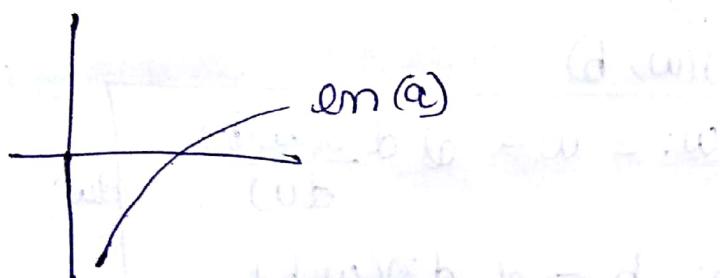
Now at $a=5$ $f(a)=25$, $a=5.001$, $f(a)=25.010$

$$\frac{d}{da} f(a) = 10 \text{ when } a=5, \boxed{\frac{d}{da} a^2 = 2a}$$

$$f(a) = \ln(a)$$

$$\frac{d}{da} f(a) = \frac{1}{a} \quad a=2; f(a)=0.693147$$

$$f(a) = 0.69365$$



$$\frac{d}{da} f(a) = 0.005$$

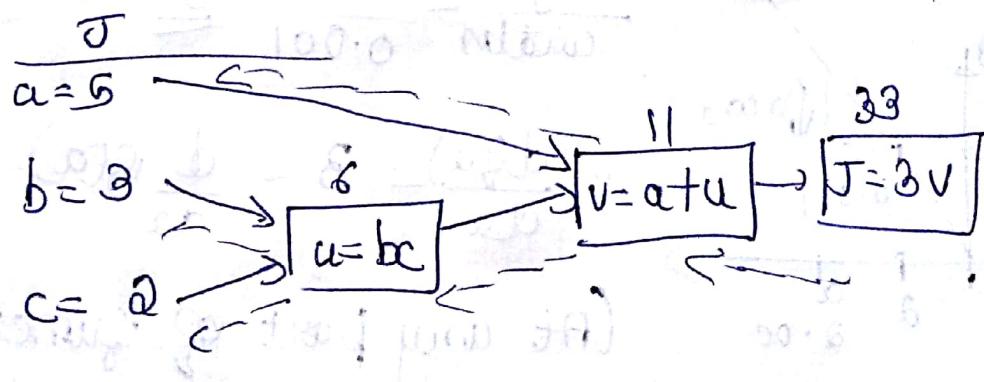
computation graph

$$J(a, b, c) = 3\left(a + \frac{bc}{u}\right) = 3(5 + 3 \cdot 2) = 33,$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



--> backward pass to
left computation of derivative

$$\frac{dJ}{dV} = ? = \underline{\underline{3}}$$

Now $J = 3V$
 $\& V = 11 \rightarrow V = 11.001$
 $J = 33.003$

$$g(a) = 3a$$

$$\frac{dg(a)}{da} = \frac{dg}{da} = 3 \quad \frac{dJ}{dV} = 3$$

Chain rule
 $a \rightarrow V$

$$\frac{dJ}{da} = \underline{\underline{3}}$$

$$a = b \rightarrow 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{da} = \frac{dV}{da} \cdot \frac{dJ}{dV}$$

$$\frac{dV}{da} = \frac{3}{1} = \underline{\underline{3}}$$

$J \rightarrow$ Final Output Variable

$$\frac{d \text{Final}/p}{d \text{var}} \rightarrow \text{intermediate}$$

$$\frac{d \text{var}}{d \text{var}}$$

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

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

$$\frac{dV}{db} = 6$$

$$\frac{dV}{dc} = 9$$

Logistic regression

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

$$y = a = \sigma(z)$$

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

$$w_1 \rightarrow \frac{dL}{dw_1} = "d_{w_1}" = \frac{d}{dz} z$$

$$x_1 \rightarrow \frac{dL}{dx_1} = "d_{x_1}" = \frac{d}{dz} (w_1 x_1 + w_2 x_2 + b)$$

$$w_2 \rightarrow \frac{dL}{dw_2} = "d_{w_2}" = \frac{d}{dz} (w_1 x_1 + w_2 x_2 + b)$$

$$b \rightarrow \frac{dL}{db} = "d_b" = \frac{d}{dz}$$

$$= a - y$$

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

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

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

$$\begin{aligned} w_1 &:= w_1 - \alpha d_{w_1} \\ w_2 &:= w_2 - \alpha d_{w_2} \\ b &:= b - \alpha d_b \end{aligned}$$

Gradient descent on m examples

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \ell(a^{(i)}, y^{(i)})$$

Compute $dw_1^{(i)}$ for $\underline{dw_1^{(i)}} \rightarrow$

all training examples & average it

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

for $i = 1$ to m

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

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

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

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

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

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

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

$$J / = m$$

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

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

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

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

Better to use Vectorization

What is vectorization?
 $z = w^T x + b$

(Shift Enter
to execute)

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

SIMD \rightarrow Single Instruction, Multiple Data Both CPU & GPU
Instructions. (Parallelisation)

$\rightarrow u = np.\text{abs}(v)$ $np.\log(v)$ $np.\text{abs}(v)$
 $np.\text{maximum}(v, 0)$ (element wise)

Vectorizing Logistic regression

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

$v * x_2$
square

v \rightarrow element
wise

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

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

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

\Downarrow auto converts b to row
(broadcasting)

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

Vectorizing gradient descent.
 $dz = [a - y]_{1 \times m}$

$$dZ = A - Y.$$

$$\begin{aligned} dw_t &= x^{(1)} dz^{(1)} & db &= 0 \\ &\vdots & db &+ = d(b^{(1)}) \\ dw &= x^m dz^{(2)} & & \\ dw &/= m & db &/= m \end{aligned}$$

$$\begin{aligned} db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} & \frac{1}{m} np.\text{sum}(dz) \\ dw &= \left[\frac{1}{m} X \cdot dz^T \right]_{n \times 1} & X \rightarrow n \times m \\ && dz \rightarrow m \times 1 \\ && dw \rightarrow n \times m \times m \times 1 \\ && = n \times 1 \end{aligned}$$

Put all boxes for multiple
iterations in a loop
for i in range(1000):

$$\begin{aligned} &\text{& then } w = w - ddw \\ &b = b - ddw. \end{aligned}$$

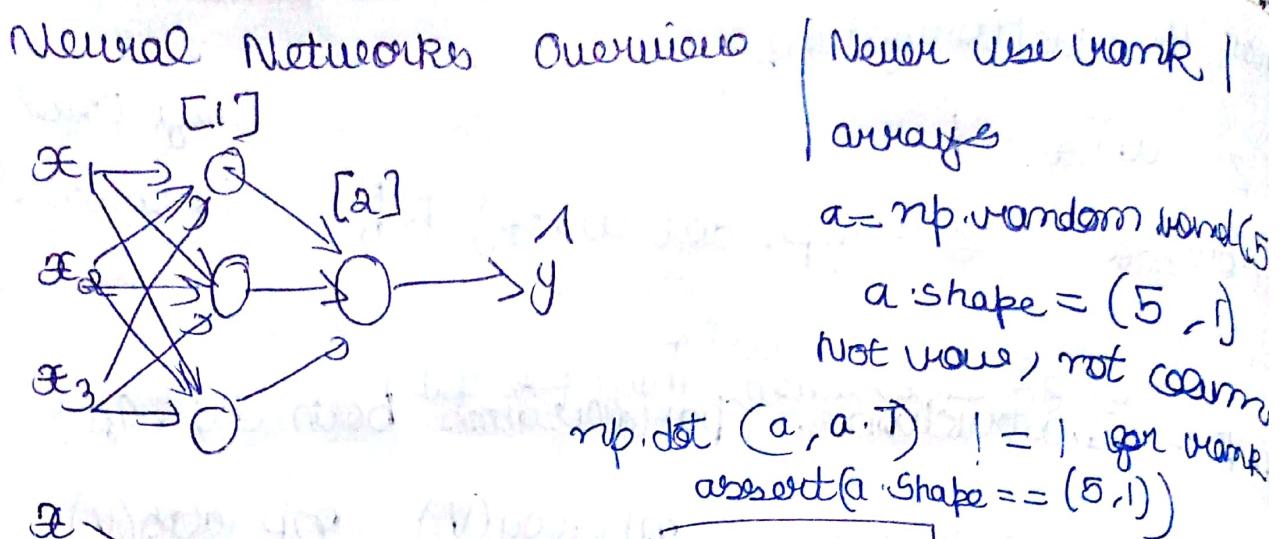
$$(m, n) \stackrel{?}{=} (1, m)$$

$\stackrel{?}{=}$
 $\stackrel{?}{=}$

$$(m, 1) \stackrel{?}{=} R \stackrel{?}{\rightarrow} (m, 1)$$

$\stackrel{?}{=}$
 $\stackrel{?}{\rightarrow}$

A.sum(axis=0)
sum vertically

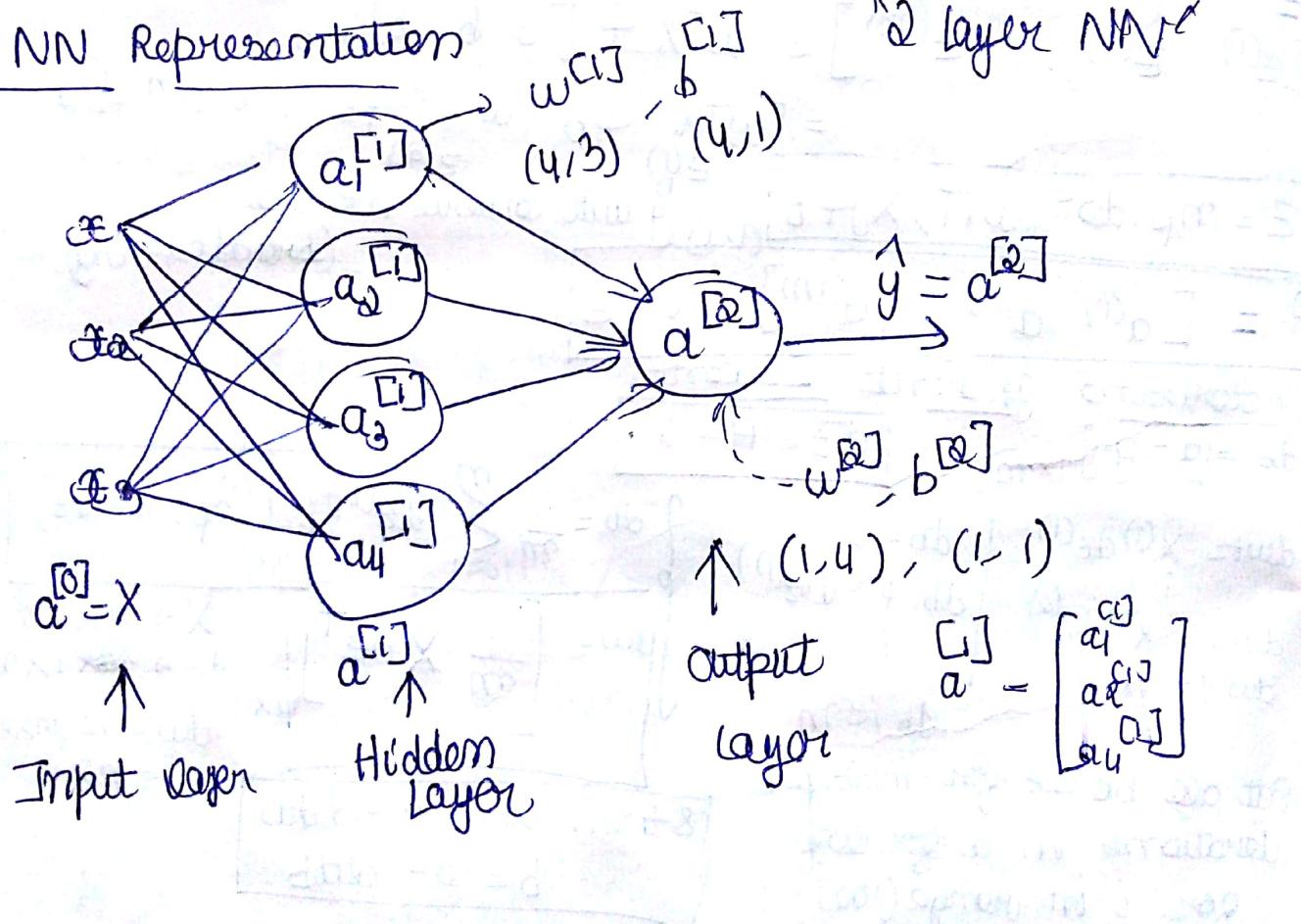


$$x \rightarrow z^{[1]} = w^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

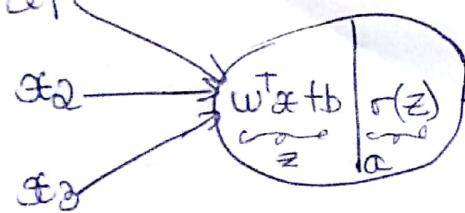
$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow L(a^{[2]}, y)$$

$$\frac{\partial L}{\partial a^{[2]}}$$

NN Representation



computing NN O/P



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

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

$$z_3^{[1]} = w_3^{[1]T} \cdot a + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} \cdot a + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

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

$$a = \sigma(z)$$

z as a vector

$$z = \begin{bmatrix} -(w_1)^T \\ -(w_2)^T \\ -(w_3)^T \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_1^T \cdot x + b \\ \vdots \\ w_4^T \cdot x + b \end{bmatrix}$$

(unvectorized)

given input x :

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

$$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$$

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

$$4 \times 1$$

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

$$(1,4) \quad (4,1) \quad (1,1)$$

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

$$1 \times 1$$

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

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

unvectorized

for $i=1$ to m ,

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

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

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

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

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

$$x \rightarrow a^{[0]} = \hat{y}$$

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

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

$$x^{(m)} \rightarrow a^{[1]}(m) = \hat{y}(m)$$

layer example

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

Vectorized Implementation

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

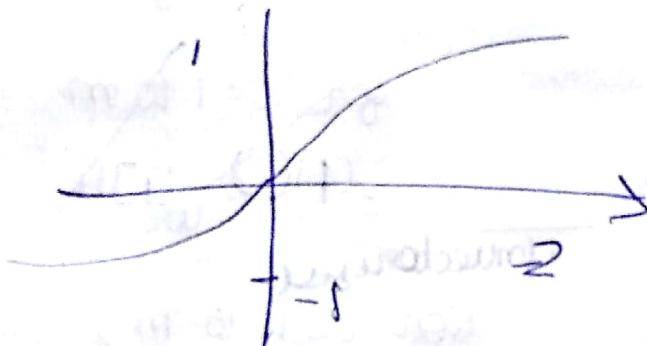
$$A^{[1]} = \sigma(z^{[1]})$$

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

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

Activation Functions

tanh function hyperbolic tan



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

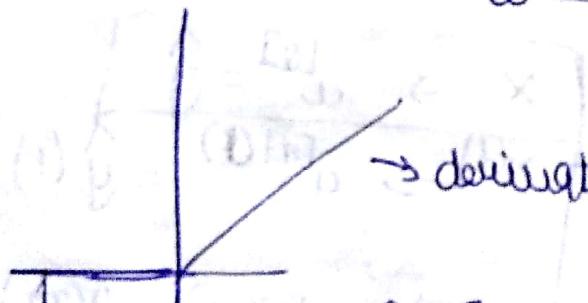
works better as
mean will $b = 0$

$\tanh h \rightarrow$ hidden layer (data is centred)
sigmoid \rightarrow output layer (as \hat{y} should be $0 \leq \hat{y} \leq 1$)

ReLU

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

→ derivative 1 everywhere

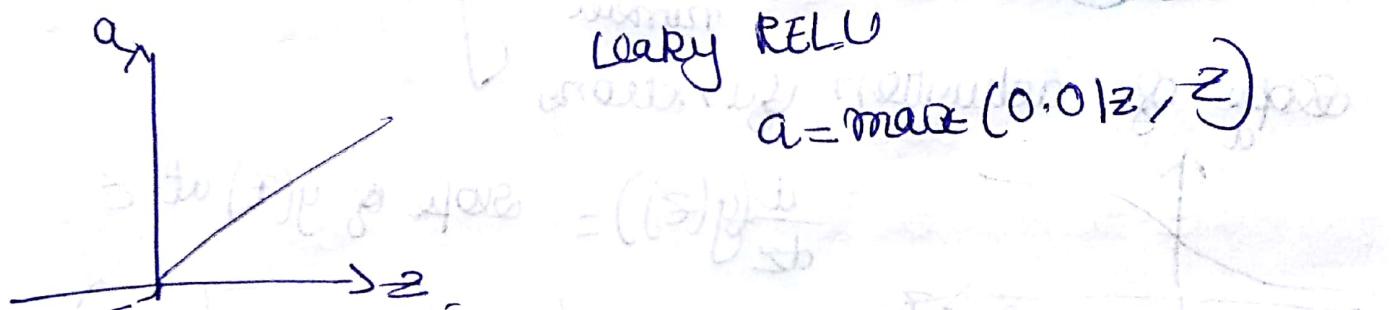
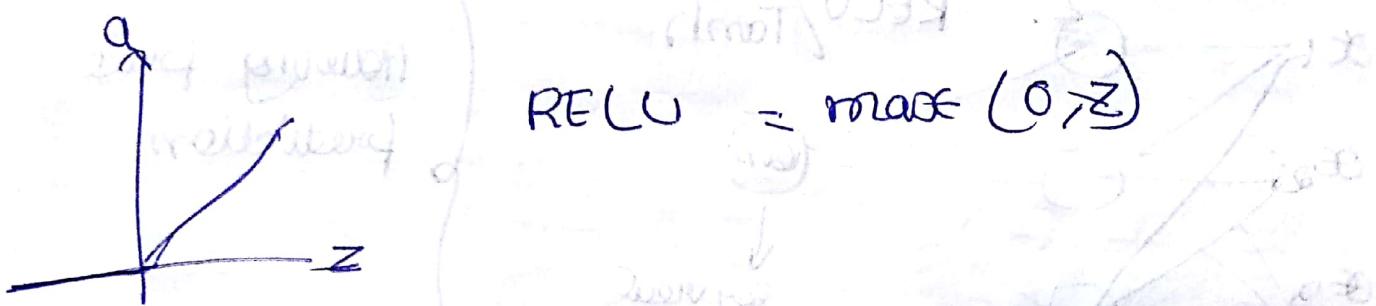
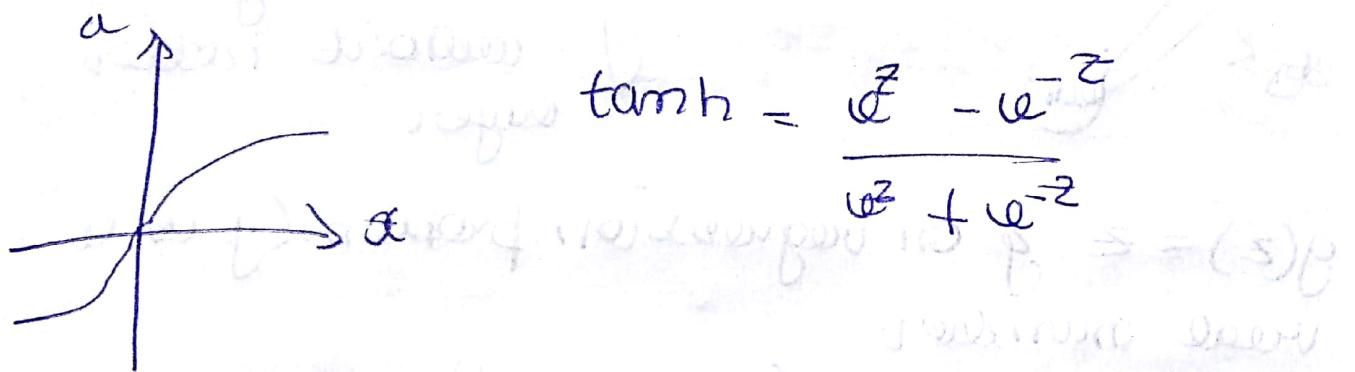
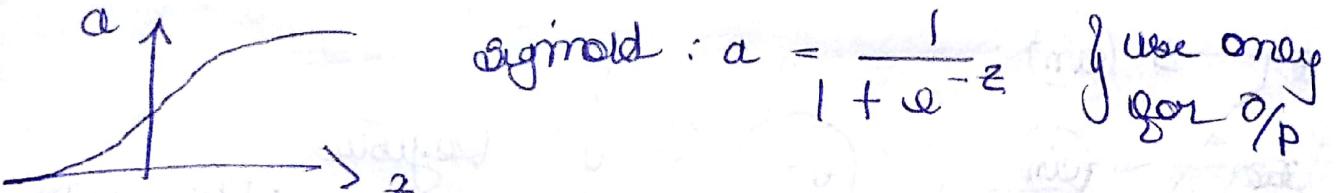


slope 0 (steery ROL better)

ReLU (tanh sometimes)

default for
all hidden
layers

Sigmoid \rightarrow O/P 0.5



Why use activation functions?

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

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

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

"linear activation function"

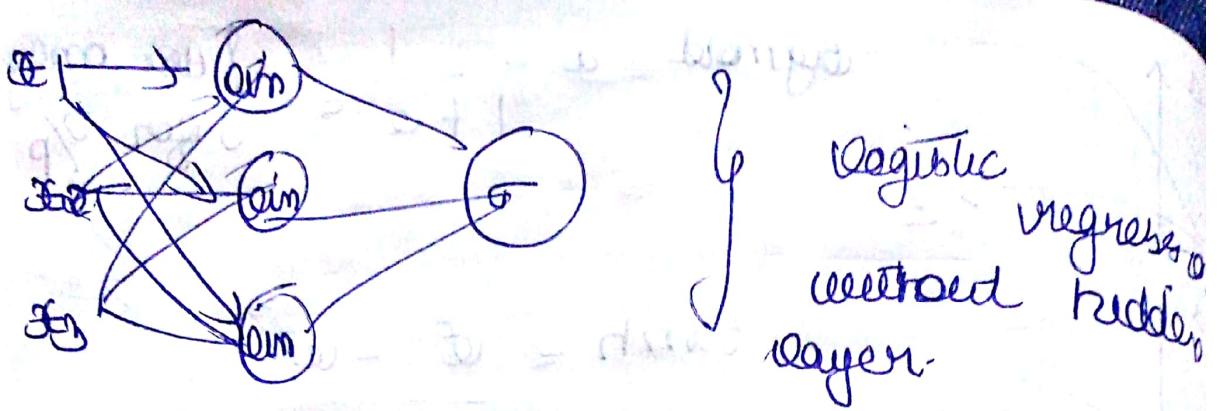
$$a^{[l]} = g(z^{[l]})$$

$$\Rightarrow a^{[l]} = w^{[l]} (w^{[l-1]} x + b^{[l-1]}) + b^{[l]}$$

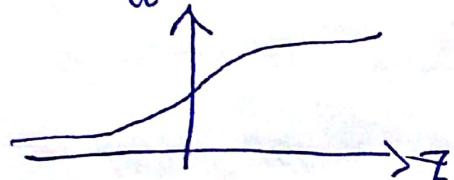
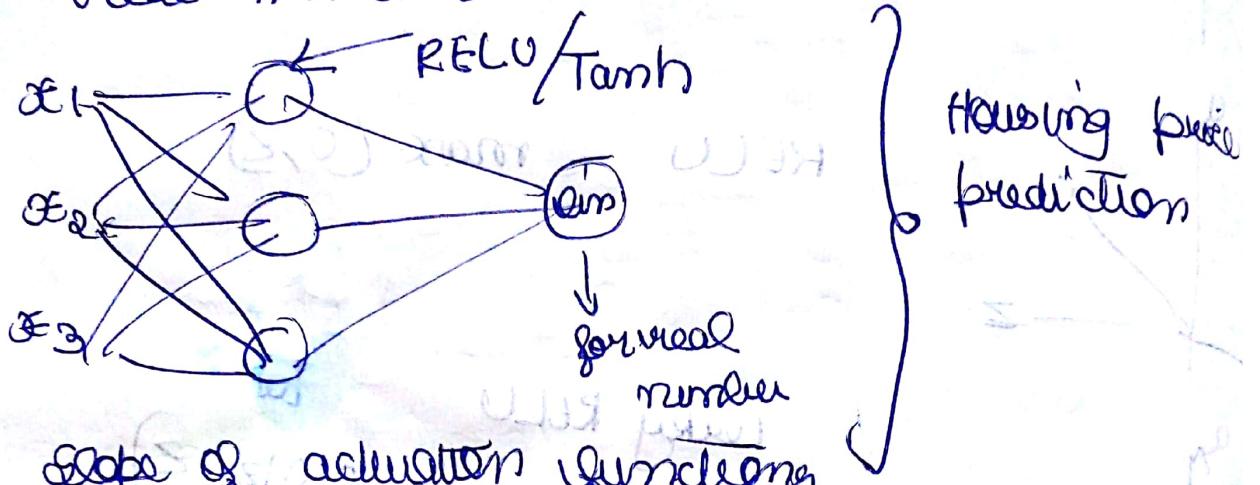
$$= (\underbrace{w^{[l]} w^{[l-1]}}_{w^l}) x + (\underbrace{w^{[l]} b^{[l-1]} + b^{[l]}}_{b^l})$$

$$= w^l x + b^l$$

} linear regression



$g(z) = z$ on regression problem (y is a real number)



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

$$z = 10, g(z) = \frac{1}{1+e^{-10}} = 1 - \frac{1}{1+e^{-10}} = \frac{1}{1+e^{-10}} \left(1 - \frac{1}{1+e^{-10}} \right)$$

$$\begin{aligned} \frac{d(g(z))}{dz} &= 1(1-1) \\ &= \underline{0} \end{aligned}$$

$$z = -10 \quad g(z) = \frac{1}{1+e^{10}}$$

$$\begin{aligned} \frac{d(g(z))}{dz} &= 0(1-0) \\ &= \underline{0} \end{aligned}$$

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

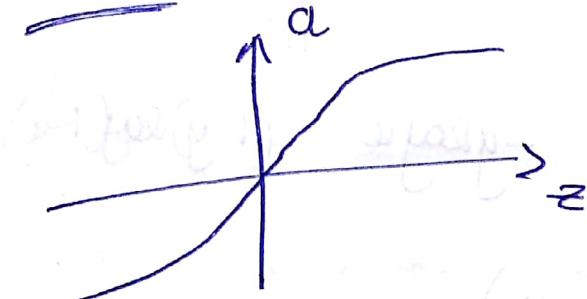
$$z = 0 \quad g(z) = \frac{1}{2}$$

$$\begin{aligned} \frac{d(g(z))}{dz} &= \frac{1}{2} \left(1 - \frac{1}{2} \right) \\ &= \underline{\frac{1}{4}} \end{aligned}$$

Im N'Ns

$$g'(z) = a(1-a) \quad \text{as } a^{[0]} = x \\ a^{[1]} \text{ next layer etc}$$

tanh h



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

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

$$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ = 1 - (\tanh(z))^2$$

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

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

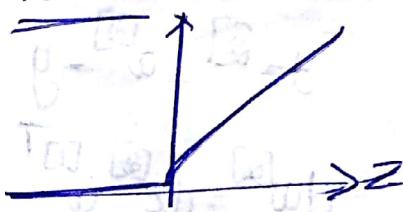
$$g'(z) = 1 - 1^2 = 0$$

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

$$g'(z) = 1$$

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

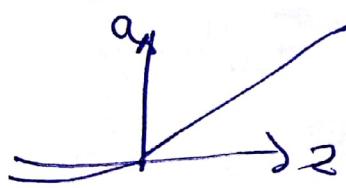
ReLU



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

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

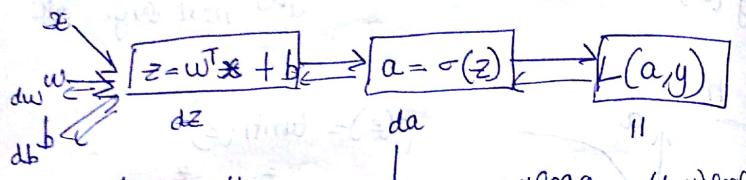
leaky ReLU



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

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

Logistic gradient



$$\begin{aligned} dz &= a - y \\ dz &= da \cdot g'(z) \\ g(z) &= \sigma(z) \\ \frac{\partial L}{\partial z} &= \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \end{aligned}$$

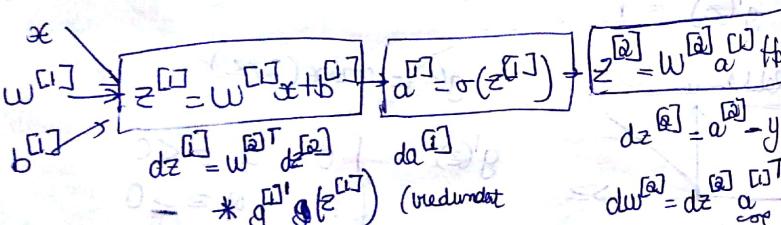
$= \frac{-y}{a} + \frac{1-y}{1-a}$

charmille

$$dz = da \cdot g'(z) \quad \frac{\partial g(z)}{\partial z} = g'(z)$$

$$dw = dz \cdot x \quad db = dz$$

In NNs we do this twice



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

(redundant)

$$dw^{[l]} = dz^{[l]} \cdot x^T$$

$$db^{[l]} = dz^{[l]} \cdot 1$$

Diagram showing a neural network structure with inputs x_1, x_2, x_3 and weights $w^{[1]}, w^{[2]}$. The output is y .

$$n^{[0]} = 3, n^{[1]} = 2, n^{[2]} = 1$$

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

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

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

$$dz^{[2]} = w^{[2]T} dz^{[1]} * g^{[2]}(z^{[2]})$$

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

$$a^{[2]} = \sigma(z^{[2]}) \rightarrow L(a^{[2]}, y)$$

$$da^{[2]} = dz^{[2]}$$

(redundant)

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

plays role of x

$$dw = dz \cdot x$$

$$w^{[1]} = \dots$$

w^{[1]} = w^{[1]} + dw^{[1]}

w^{[1]} = w^{[1]} - dw^{[1]}

Vectorized implementation

$$dz^{[2]} = A^{[2]} - y \quad \text{if } y \text{ has sigmoid}$$

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

To prevent

rank 1 always

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

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

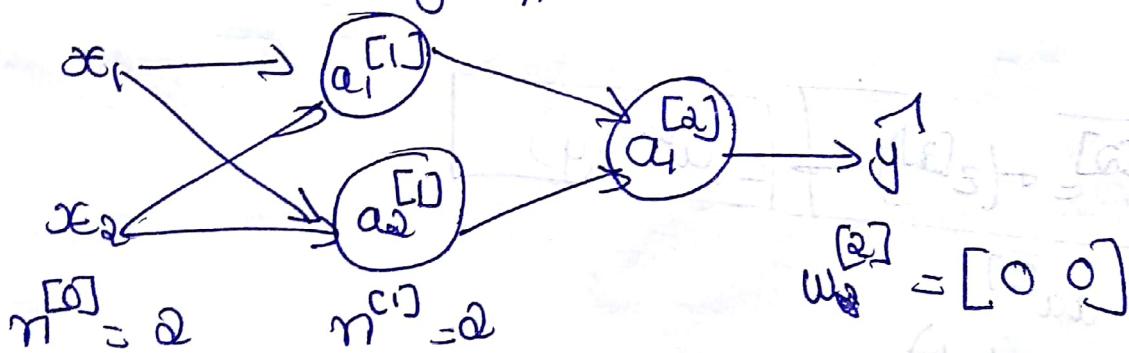
$$n^{[1]}, m \quad n^{[2]}, m$$

$n^{[1]}, m$
derivative of activation

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

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

Random initialization



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow a_1^{[1]} = a_2^{[2]}$$

In backprop $dz_1^{[1]} = dz_2^{[1]}$ [Symmetric]

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

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

$w^{[1]}$ 1st row = 2nd row

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

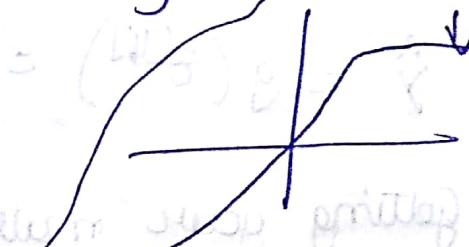
Hidden units calculate same thing in all iterations

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

$$b^{[1]} = \text{np.zeros}((2, 1)) \quad \begin{matrix} \text{to make} \\ \text{weights} \end{matrix}$$

$$w^{[2]} = \text{np.random.randn}(2, 1) * 0.01 \quad \begin{matrix} \text{small} \\ \text{weights} \end{matrix}$$

$$f^{[2]} = \text{np.zeros}$$

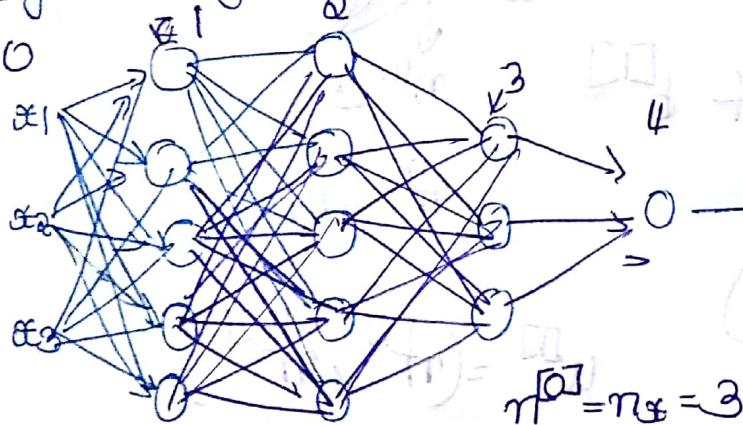


Week - 4 ~~Implementation with numpy~~

Deep Neural Networks

Logistic regression \rightarrow "shallow" model

\uparrow
to not end up
here (gradient will
be ~~zero~~ slow)



$$n^{[0]} = 5$$

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

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

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

$$L = 4$$

$$n^{[0]} = \# \text{units in layer 0} \quad w^x = a^{[0]}$$

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

$$z^{[l]} = g^{[l]}(z^{[l]}) \quad w^{[l]} = \text{weights for } z^{[l]}$$

$$z^{[0]} = w^{[0]} x + b^{[0]}$$

$$a^{[0]} = g^{[0]}(z^{[0]})$$

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

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

$$\boxed{\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}}$$

$$z^{[0]} = w^{[0]} A^{[0]} + b^{[0]}$$

$$A^{[0]} = g^{[0]}(z^{[0]})$$

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

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

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

to 4

Getting your matrix dimensions right.

$$L = 5$$

forward propagation

$$z^{[0]} = w^{[0]} x + b^{[0]}$$

(3, 1)

(3, 2) (2, 1)

[3 hidden
layers]

(n^[0], n^[1]) (n^[0], 1)

$$\begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix}$$

(n^[0], 1)

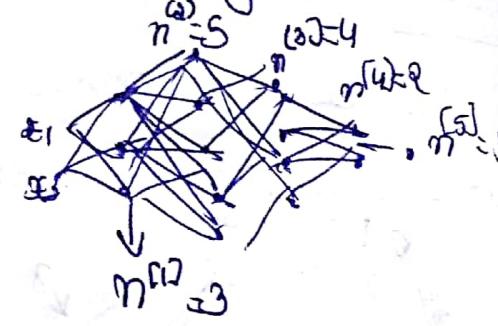
$$z^{[0]} = w^{[0]} \cdot a^{[0]} + b^{[0]}$$

(5, 1) (5, 3) (3, 1)

$$w^{[0]} = (4, 5)$$

$$w^{[1]} = (2, 4)$$

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



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

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

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

$$(5, 3)$$

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

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

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

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

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

$z^{[l]}, a^{[l]} = (n^{[l]}, 1)$

to 4
Tot. sparsity
Same dimensions

Now for multiple examples.

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

$\begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_n \end{bmatrix} \quad \begin{matrix} m \\ \vdots \\ \vdots \\ n^{[1]}, n^{[0]} \end{matrix}$

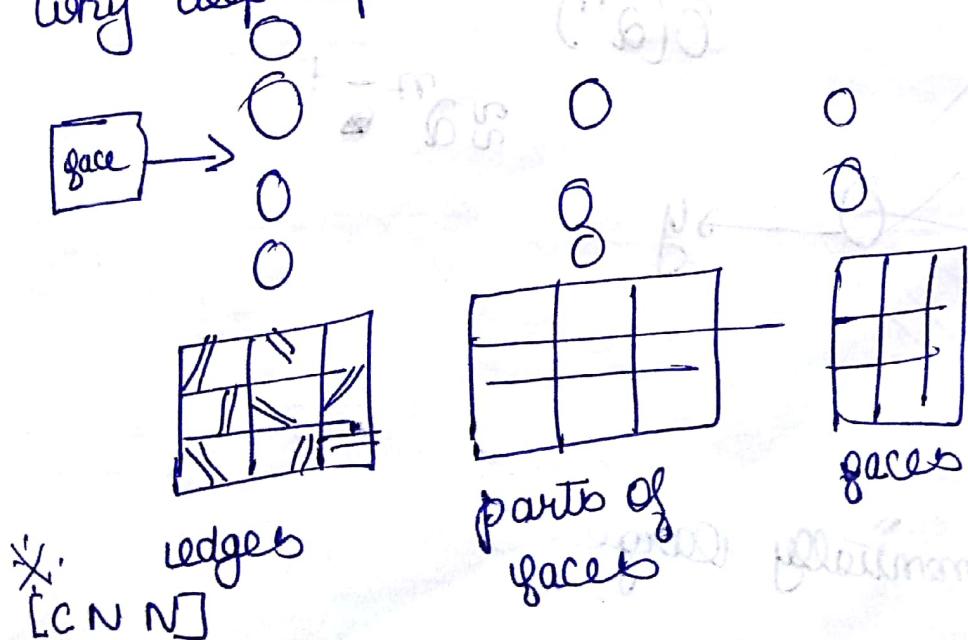
$\begin{bmatrix} z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \end{bmatrix} \quad \begin{matrix} n^{[1]}, m \\ \vdots \\ \vdots \\ n^{[0]}, m \end{matrix}$

$(n^{[1]}, 1)$
broadcasting

$n \times m$

$w^{[1]} = 0 \quad A^{[0]} = X = (n^{[0]}, m)$

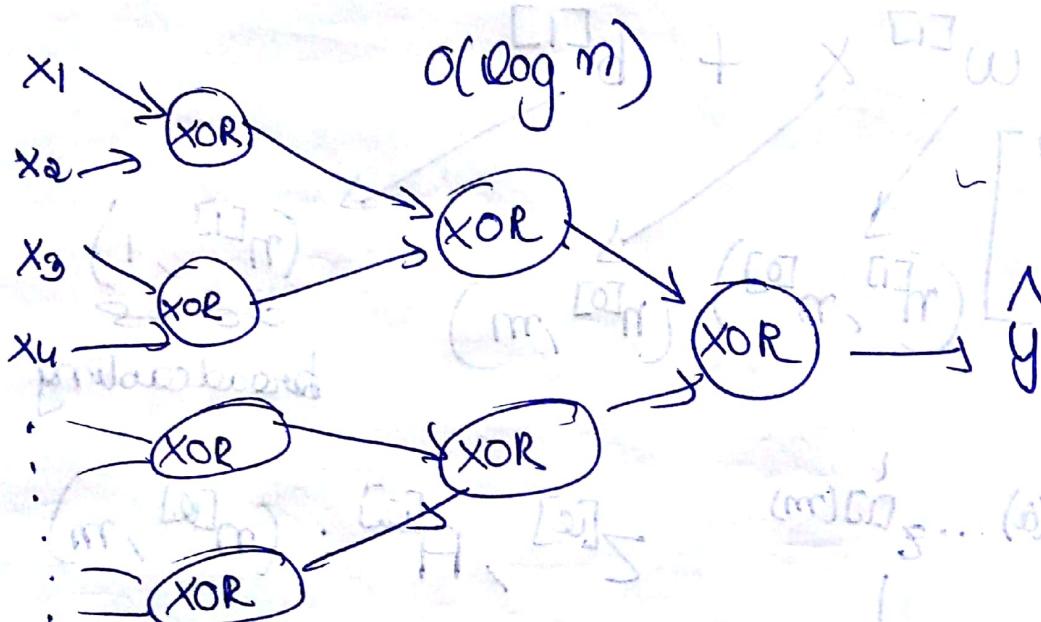
why deep representation



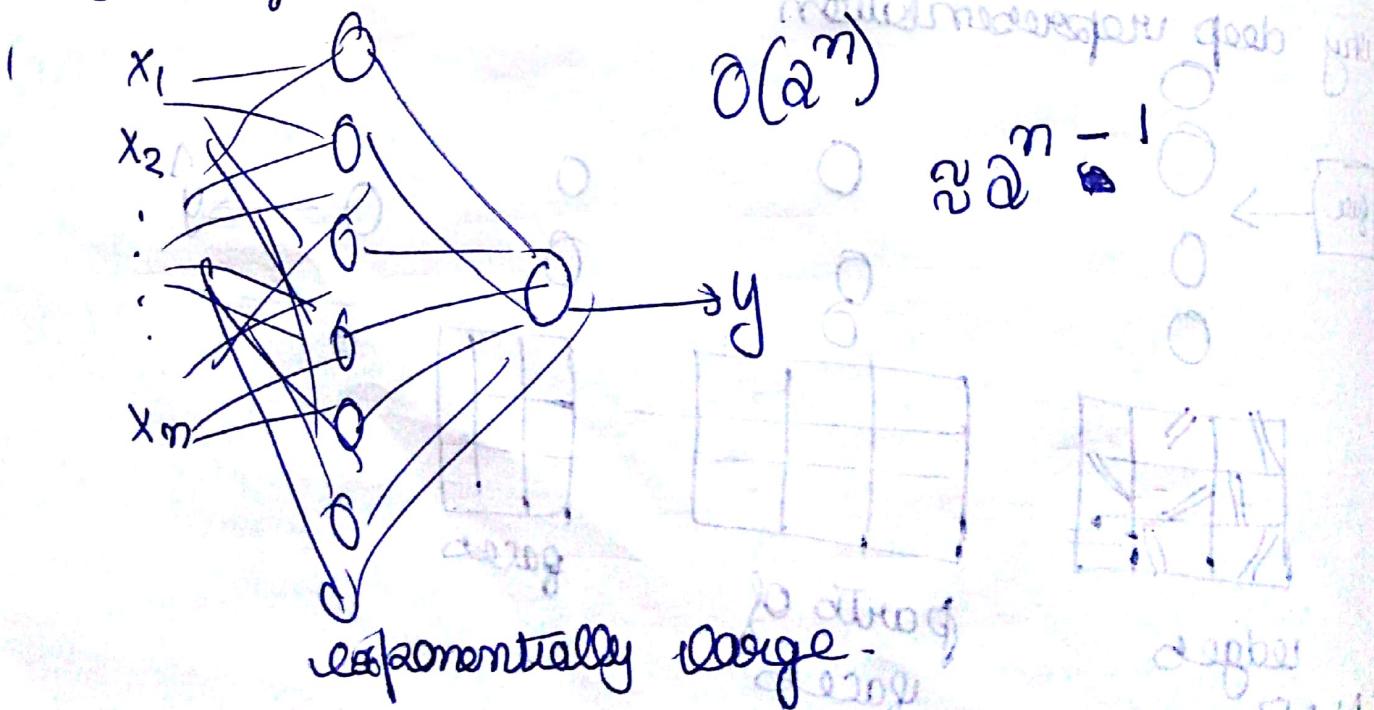


Circuit theory

There are functions you can compute with a "small" L-layer deep neural network. Shallower networks require exponentially more hidden units to compute.



But if we're not allowed L-layer



Building blocks of deep neural nets

Layer ℓ : $w^{[\ell]}, b^{[\ell]}$

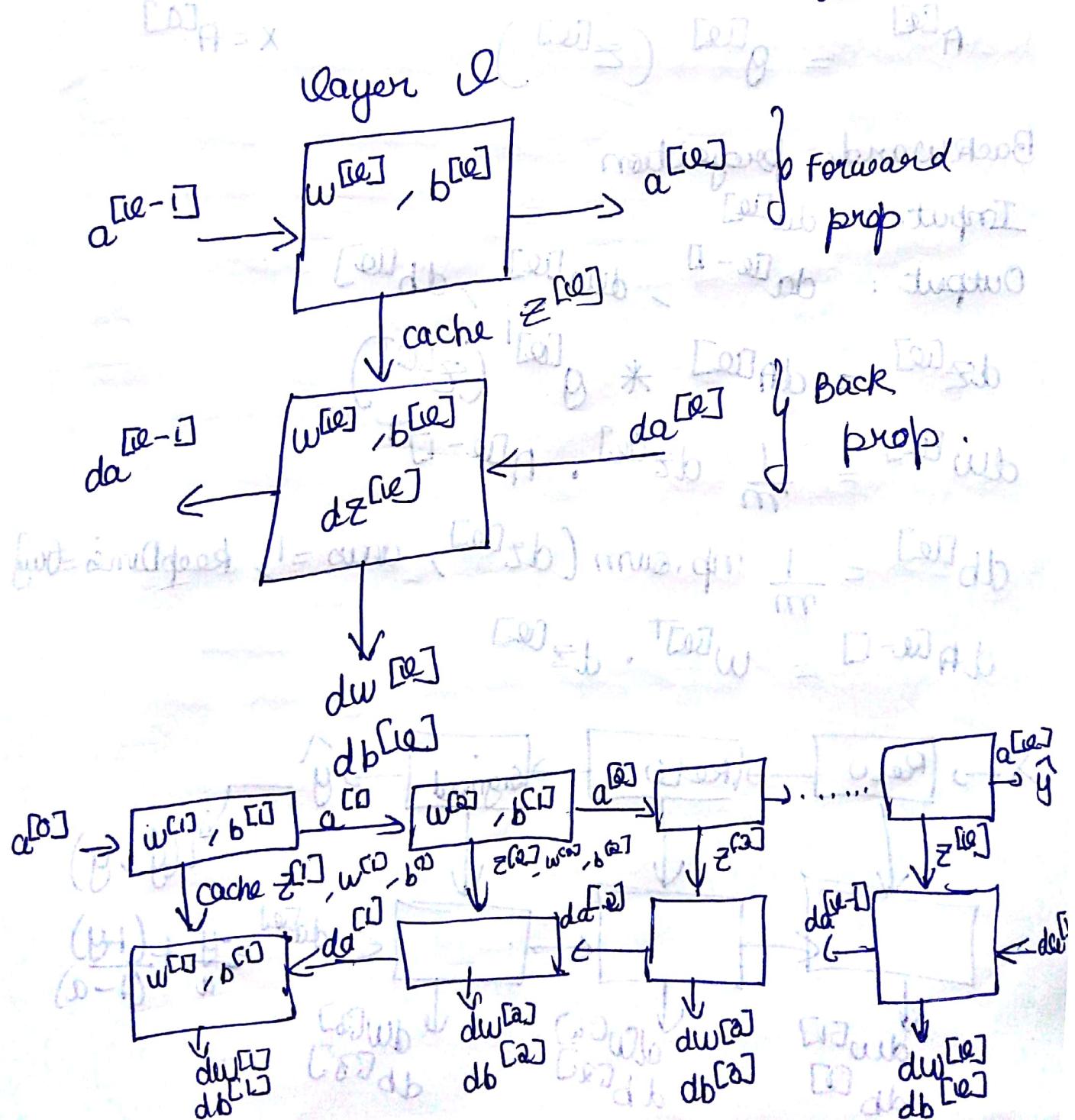
Format Input: $a^{[\ell-1]}$, output $a^{[\ell]}$.

$$z^{[\ell]} = w^{[\ell]} a^{[\ell-1]} + b^{[\ell]} \quad [\text{cache } z^{[\ell]}]$$

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

Backward: Input $da^{[\ell]}$, output $da^{[\ell-1]}$
([\text{cache } z^{[\ell]}]).

$$\begin{aligned} dw^{[\ell]} \\ db^{[\ell]} \end{aligned}$$



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

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

Forward propagation.

Input : $a^{[l-1]}$

Output $a^{[l]}$, cache ($z^{[l]}$), $w^{[l]}$, $b^{[l]}$

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

$$A^{[l]} = g^{[l]}(z^{[l]}) \quad x = A^{[0]}$$

Backward propagation

Input : $da^{[l]}$

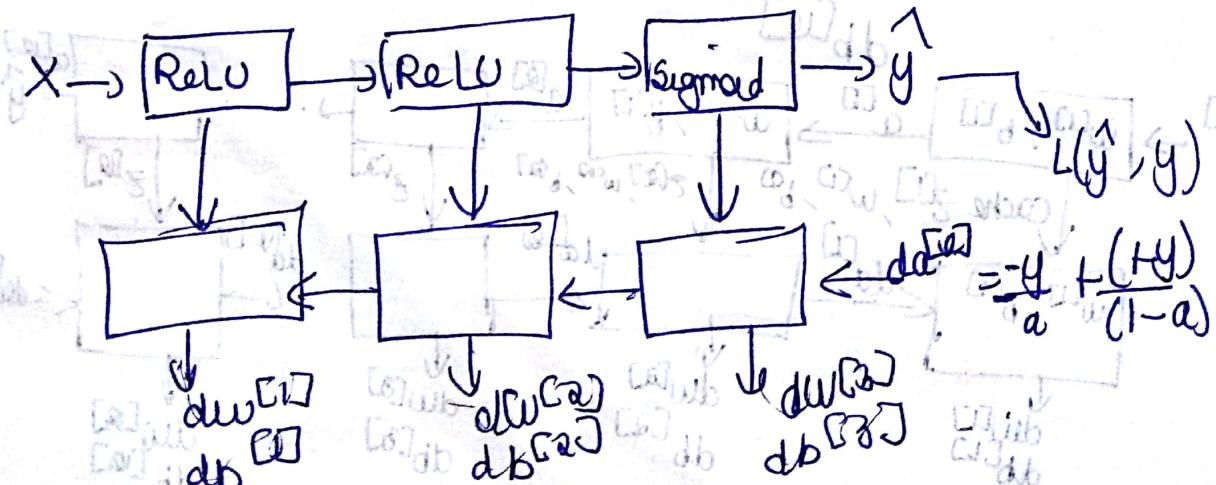
Output : $da^{[l-1]}$, $dw^{[l]}$, $db^{[l]}$

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

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

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

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



Hyperparameters

Learning rate α

iterations

hidden layer L

hidden units $n^{[l]}, n^{[0]}$

Choice of activation functions

"It's like the brain"

