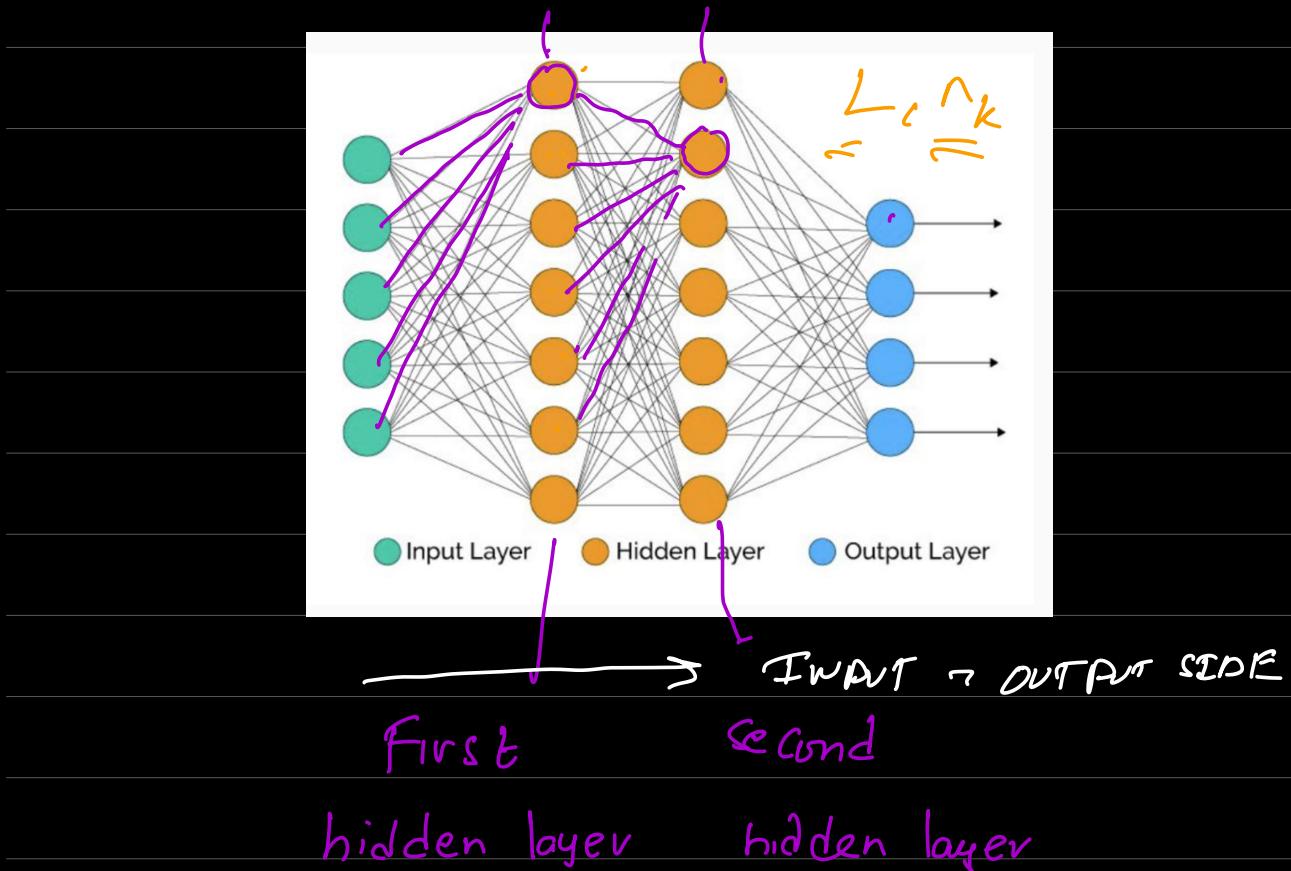


MLP → Multi Layer perceptrons

ANN → Artificial Neural Network

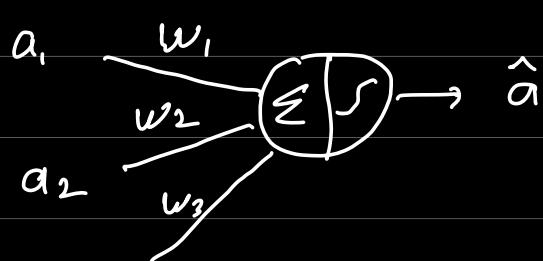
↳ Feed forward Neural Network

↳ Fully Connected feed forward Neural Network



Neuron - Computations

①



$$z = \sum w_i a_i$$

$$\hat{a} = g(z)$$

$a_3 \rightarrow$ activations

Sigmoid $\sim \sigma(z)$

ReLU $\sim \text{ReLU}(z)$

Forward pass

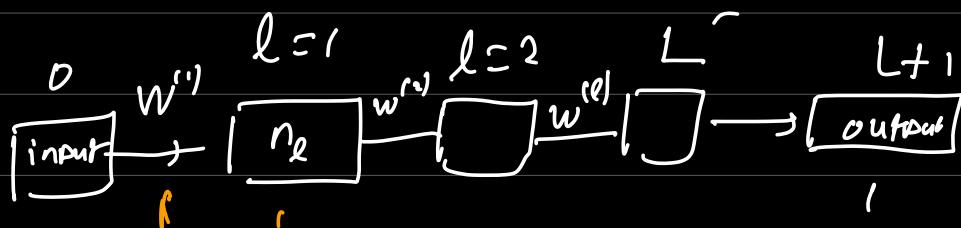
General case \rightarrow Arbitrary input vector

(Multiple rows in input)

\rightarrow Arbitrary no of layers

\rightarrow Arbitrary no of neurons
in each layer

\rightarrow No of output based
on problem



$n_e \rightarrow$ number of neurons in
each layer

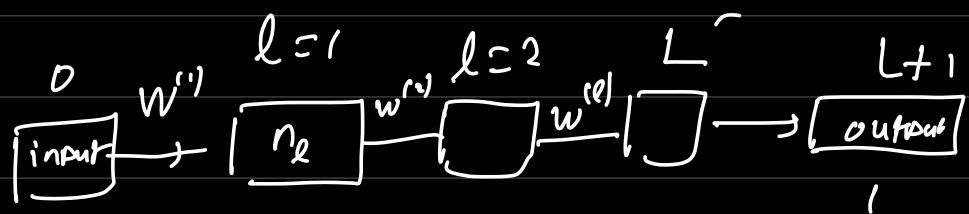
$n_e \sim$ is the number of input features

$m \sim$ Total number of data points

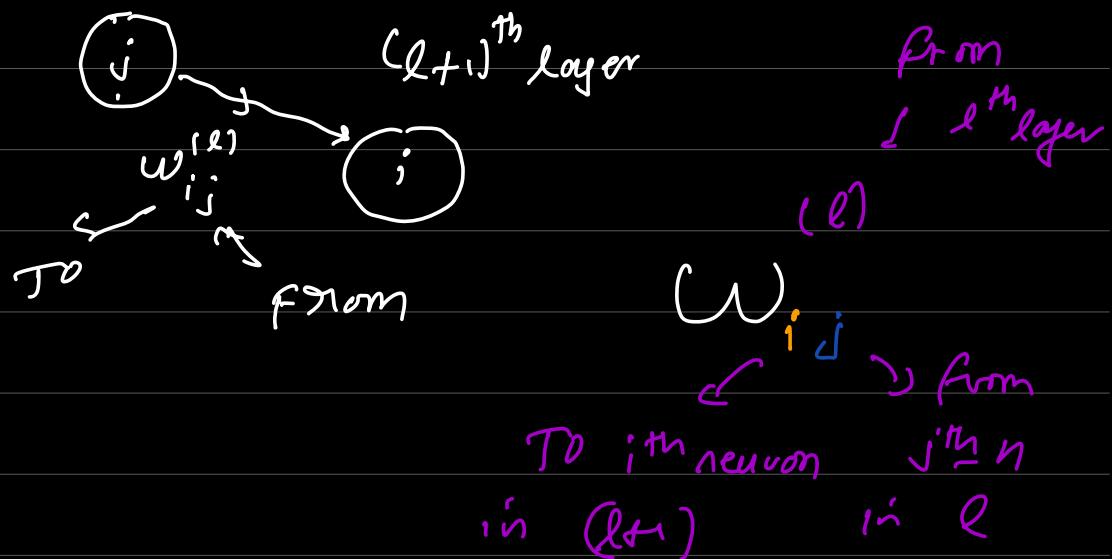
input $\rightarrow (x_1, x_2) \dots (x_l, \dots, x_n)$

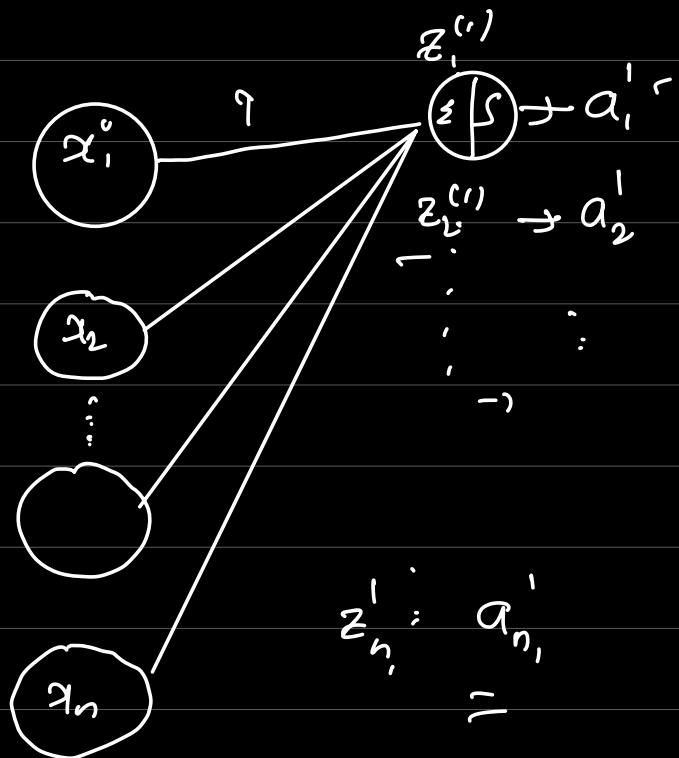
$$X = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \dots & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & \dots & x_n^{(2)} \\ \vdots & & & & \\ x_1^{(m)} & x_2^{(m)} & \dots & \dots & x_n^{(m)} \end{pmatrix}$$

$$m \times n - X_{m \times n}$$



l^{th} layer $1 \dots n_l$





$$a_1 = x_1, \quad a_2 = x_2, \dots, \quad a_n = x_n$$

for one row of X

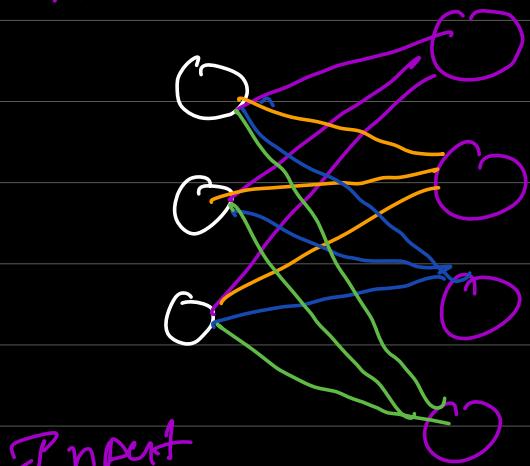
$$\left\{ \begin{array}{l} w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + \dots + w_{1n}x_n = z_1^{(1)} \\ \vdots \\ w_{11}a_1 + w_{12}a_2 + w_{13}a_3 + \dots + w_{1n}a_n = z_1^{(1)} \\ \\ w_{21}a_1 + w_{22}a_2 + w_{23}a_3 + \dots + w_{2n}a_n = z_2^{(1)} \\ \vdots \\ w_{n1}a_1 + w_{n2}a_2 + w_{n3}a_3 + \dots + w_{nn}a_n = z_n^{(1)} \end{array} \right.$$

$$g \left(\begin{pmatrix} z_1^{(1)} \\ z_2^{(1)} \\ \vdots \\ z_{n_1}^{(1)} \end{pmatrix} \right) \rightarrow \begin{pmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_{n_1}^{(1)} \end{pmatrix}$$

$$\begin{pmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \dots & x_n^{(1)} \end{pmatrix} \xrightarrow{\quad \text{!} \quad} \begin{pmatrix} z_1 & z_2 & z_{n_1} \\ w_{1,1}^{(1)} & w_{1,2}^{(1)} & \vdots & w_{n_1,1}^{(1)} \\ w_{1,2}^{(1)} & w_{1,3}^{(1)} & \vdots & w_{n_1,2}^{(1)} \\ \vdots & \vdots & \vdots & \vdots \\ w_{1,n}^{(1)} & w_{2,n}^{(1)} & \dots & w_{n_1,n}^{(1)} \end{pmatrix}$$

$n=3$

$[n \times n_1]$



Input

Layer

First layer $n_1=3$

$$\left(\begin{matrix} x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots \\ x_1^{(n)} & \dots & x_n^{(n)} \end{matrix} \right) \left(\begin{matrix} w^{(1)} \\ \vdots \\ w^{(n)} \end{matrix} \right) = \left[\begin{matrix} z_1^{(1)} & \dots & z_{n_1}^{(1)} \\ \vdots & & \vdots \\ z_1^{(n)} & \dots & z_{n_1}^{(n)} \end{matrix} \right]$$

$[m \times n] \quad [n \times n_1] \rightarrow [m \times n_1]$

$\begin{matrix} [m \times n] \\ \in \end{matrix} \quad [n \times n_1] \rightarrow [m \times n_1]$,
 Z-Matrix

$$\rightarrow \left(\begin{matrix} x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & & \vdots \\ x_1^{(m)} & \dots & x_n^{(m)} \end{matrix} \right) \left(\begin{matrix} w^{(1)} \\ \vdots \\ w^{(n)} \end{matrix} \right) = \left(\begin{matrix} z_1^{(1)} & \dots & z_{n_1}^{(1)} \\ \vdots & & \vdots \\ z_1^{(m)} & \dots & z_{n_1}^{(m)} \end{matrix} \right)$$

X
—

$$X \cdot W = Z$$

$\rightarrow m \times n \quad n \times n_1 \quad \rightarrow m \times n_1$

No of data points h number of neurons in the hidden layer

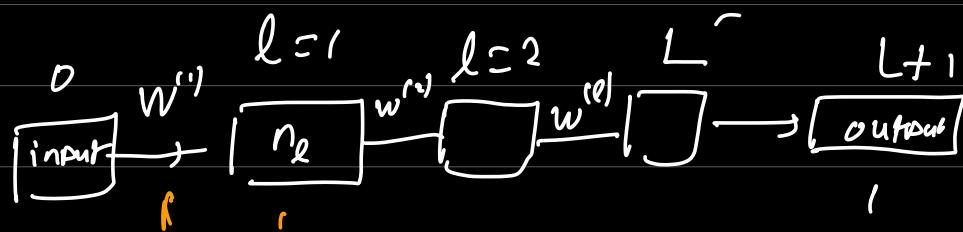
$$X \cdot W = Z \cdot g(Z)$$

↓
Pointwise

$$A = g(Z) \quad | \times n_1$$

↓

$M \times n_1$



$$X \rightarrow X \cdot W^{(1)} \xrightarrow{g} A^{(1)} \rightarrow A^{(1)} \cdot W^{(2)}$$

... $\leftarrow A^{(2)} \xleftarrow{g}$

$\underbrace{\qquad\qquad\qquad}_{A^{(L)} \cdot W^{(L+1)}} \rightarrow \text{Output}$

$$x \rightarrow [w^{(1)}]^T x \xrightarrow{g} A^{(1)} \rightarrow [w^{(m)}]^T A^{(1)}$$

-- ↘

$$\left[\begin{array}{c} \vdots \\ [w^{(L+1)}]^T A^{(L)} \end{array} \right] \rightarrow \text{output}$$

Assuming you know the weights

$$\left(\begin{array}{c} x \rightarrow xw \xrightarrow{g} A^T \end{array} \right) \rightarrow \left(\begin{array}{c} x \\ A \end{array} \right)$$

one layer
calculation

Each row corresponds
to a data point

$$x \quad w^{(1)} \quad w^{(2)} \dots \quad w^{(L)}$$

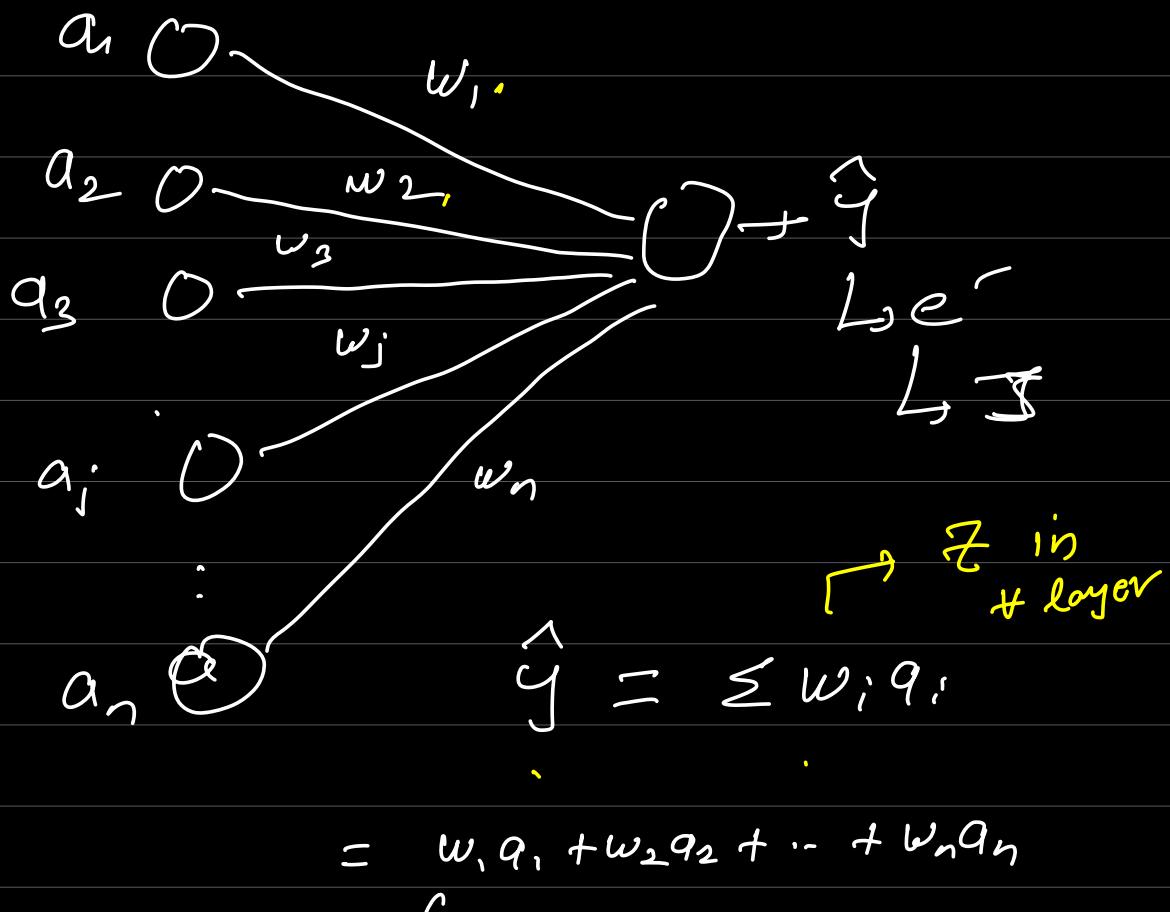
↳ Another linear regression
Model

Back propagation Algorithm

No hidden layer Model

With and without non-linearities

THE DELTA RULE



$$e \leftarrow \underline{y - \hat{y}}, \quad y \rightarrow \text{ground truth}$$

label

$$J(y, \hat{y}) = \frac{e^2}{2} = J(e),$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_1}$$

$=$

$$= e^{(-1)} a_1 = \boxed{-a_1, e}$$

$\Delta w_1 = - \frac{\partial J}{\partial w_1} = a_1 e$

$$w_j \leftarrow w_j + \alpha \Delta w_j$$

$$\hookrightarrow w_j - \alpha \frac{\partial J}{\partial w_j}$$

$$\text{update term} = \Delta w_j = a_j e$$

$a_j \rightarrow$ input to the weight

$e \rightarrow$ error in the output

This is for a linear neuron

$$z = \sum w_i a_i; \quad \hat{y} = g(z)$$

$$\begin{aligned} -\frac{\partial I}{\partial w_j} &= -\frac{\partial I}{\partial e} \underbrace{\frac{\partial e}{\partial \hat{y}}}_{F^{-1}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_j} \\ &= e \cdot g'(z) a_j \end{aligned}$$

$g'(z) \rightarrow$ Derivative of g w.r.t z , evaluated
at particular z

$$\Delta w_j = g'(z) a_j e \rightarrow$$

Sigmoid is bad

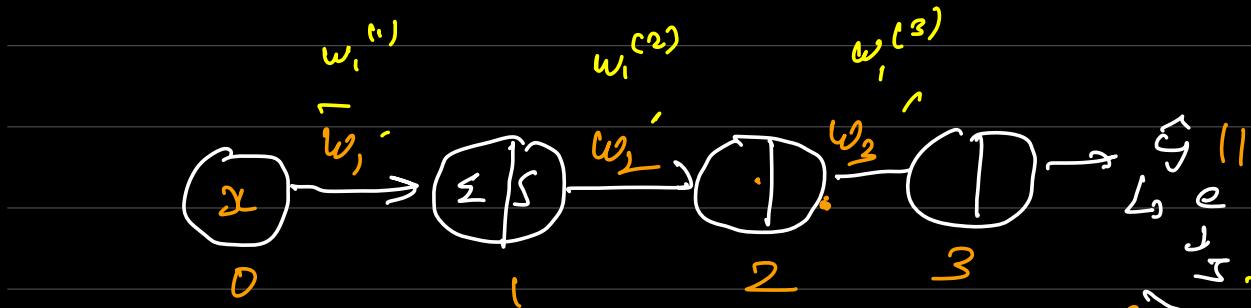
non-linearly

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

for $g(z) = \sigma(z)$

$$X = A^0$$

Multilayer Single Neuron Case



$$\begin{aligned} \underline{q}_0 &= x \\ \text{Input} & \quad \quad \quad \underline{z}_1 = w_1 q_0 \\ & \quad \quad \quad \underline{z}_2 = \underline{w}_2 \underline{q}_1 \\ & \quad \quad \quad \underline{z}_3 = \underline{w}_3 \underline{q}_2 \\ & \quad \quad \quad \underline{q}_1 = g(z_1) \\ & \quad \quad \quad \underline{q}_2 = g(z_2) \\ & \quad \quad \quad \underline{q}_3 = g(z_3) \\ & \quad \quad \quad e_3 = (y - \hat{y}) \end{aligned}$$

$\boxed{\hat{y} = q_3}$

what is the objective?

$$-\frac{\partial J}{\partial w_1}, \quad -\frac{\partial J}{\partial w_2}, \quad -\frac{\partial J}{\partial w_3} \quad ?$$

$\Delta w_1, \quad \Delta w_2, \quad \Delta w_3$

$$w_1^{\text{new}} \leftarrow w_1^{\text{old}} + \Delta w_1 (\alpha)$$

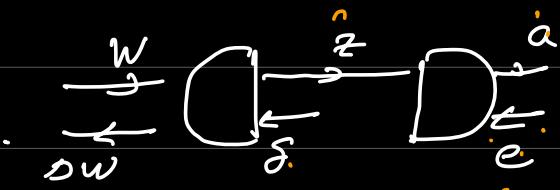
$$w_2^{\text{new}} \leftarrow w_2^{\text{old}} + \Delta w_2 (\alpha)$$

Δ

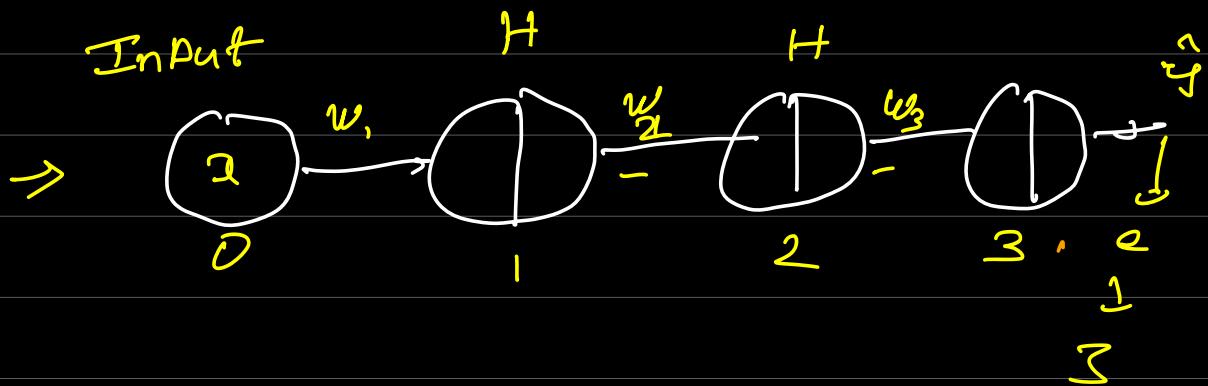
$$\boxed{w^{\text{new}} = w^{\text{old}} + \alpha \Delta w}$$

Requisites calculation of intermediate quantities

$\delta_1, \delta_2, \delta_3$
 $\frac{\partial J}{\partial z_1}, \frac{\partial J}{\partial z_2}, \frac{\partial J}{\partial z_3} \rightarrow \delta_1, \delta_2, \delta_3$
 'linear' terms



e_1, e_2, e_3
 $\frac{\partial J}{\partial q_1}, \frac{\partial J}{\partial q_2}, \frac{\partial J}{\partial q_3} \rightarrow e_1, e_2, e_3$
 'non-linear' terms



$$\frac{\partial J}{\partial w_1} \rightarrow \left[\frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial q_3} \cdot \frac{\partial q_3}{\partial q_2} \cdot \frac{\partial q_2}{\partial q_1} \right] \frac{\partial z_1}{\partial w_1}$$

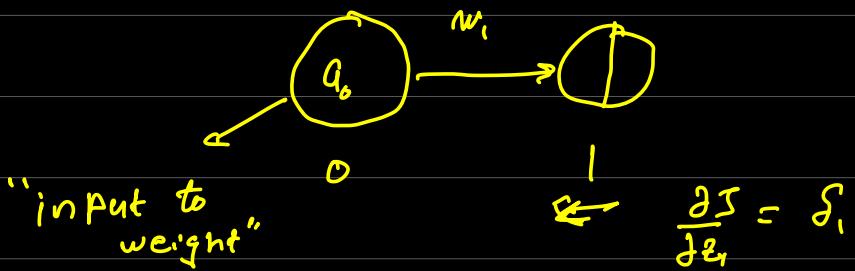
\nearrow

$$[\tau_1 = a_0 w_1]$$

$$= \frac{\partial J}{\partial z_1} \cdot q_o = \delta_1 \cdot q_o$$

\hookrightarrow input

$$\begin{aligned}\frac{\partial J}{\partial w_i} &= \left\{ \frac{\partial J}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_i} \right\} \\ &= \delta_1 \cdot q_o\end{aligned}$$



$$\frac{\partial J}{\partial w_i} = \left(\frac{\partial J}{\partial y} \frac{\partial y}{\partial a_3} \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \right) \frac{\partial z_1}{\partial w_i}$$

$$\frac{\partial J}{\partial z_1} q_o = \delta_1 q_o \quad \text{"delta rule"}$$

We don't know δ_1

$$\delta_1 = \left\{ \begin{pmatrix} \frac{\partial J}{\partial q_1} & \frac{\partial J}{\partial q_2} & \frac{\partial J}{\partial q_3} & \frac{\partial J}{\partial z_1} & \frac{\partial J}{\partial z_2} & \frac{\partial J}{\partial z_3} \end{pmatrix} \middle| \begin{array}{l} \frac{\partial q_1}{\partial a_1}, \\ \frac{\partial q_1}{\partial z_1} \end{array} \right\} = \frac{\partial J}{\partial z_1}$$

$$\frac{\partial J}{\partial a_1}, \quad \frac{\partial J}{\partial z_1}.$$

$$\frac{\partial a_1}{\partial z_1} = ? \quad a_1 = g(z_1) \\ \frac{\partial a_1}{\partial z_1} = g'(z_1)$$

$$\delta_1 = g'(z_1) e_1$$

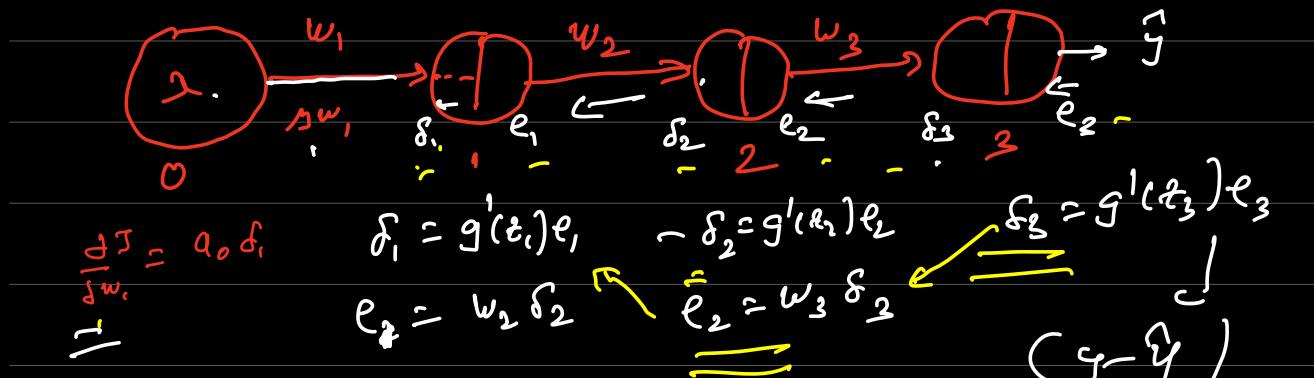
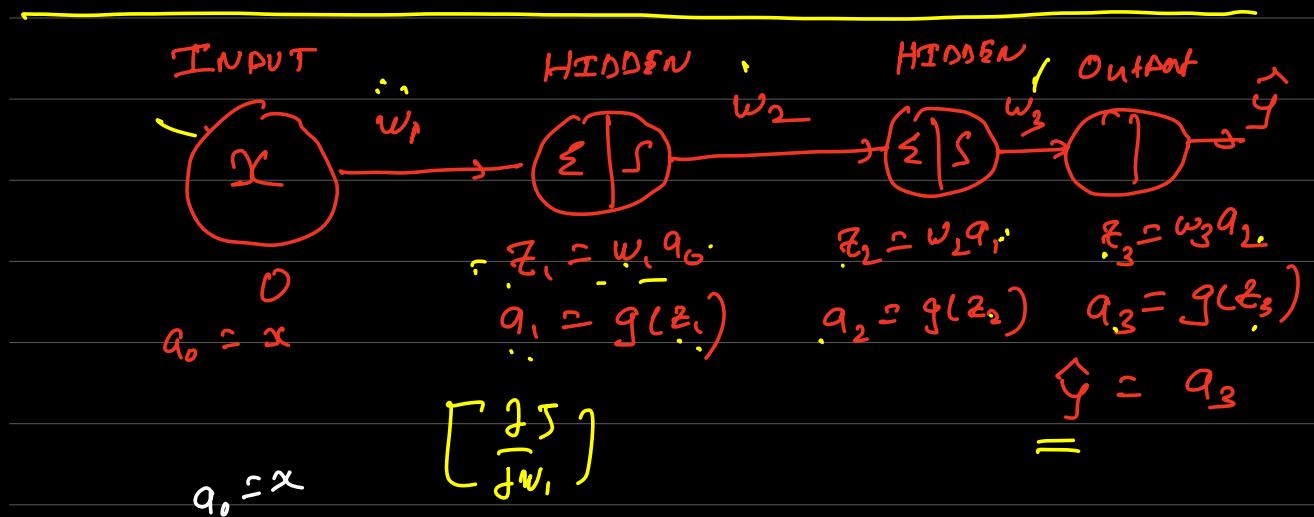
what is e_1 ?

$$\left\{ \begin{pmatrix} \frac{\partial J}{\partial q_1} & \frac{\partial J}{\partial q_2} & \frac{\partial J}{\partial q_3} & \frac{\partial J}{\partial z_1} & \frac{\partial J}{\partial z_2} & \frac{\partial J}{\partial z_3} \end{pmatrix} \middle| \begin{array}{l} \frac{\partial q_1}{\partial a_1}, \\ \frac{\partial q_1}{\partial z_1} \end{array} \right\} = \frac{\partial J}{\partial a_1}$$

$$= \frac{\partial J}{\partial z_2}, \quad \frac{\partial z_2}{\partial a_1}, \quad z_2 = w_2 a_1, \\ \frac{\partial z_2}{\partial a_1} = w_2$$

$$\frac{\partial J}{\partial a_1} = e_1 = \frac{\partial J}{\partial z_2} \cdot w_2$$

$$= e_1 = \delta_2 w_2$$



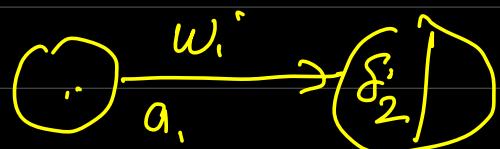
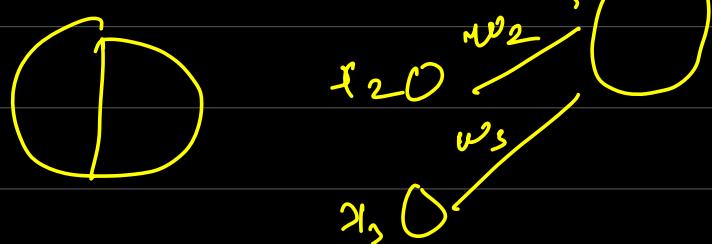
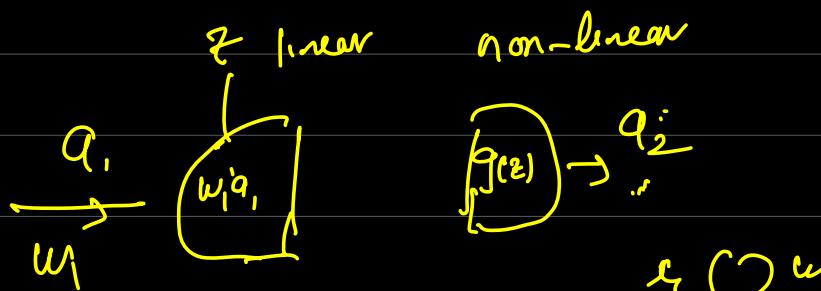
$$\frac{\partial J}{\partial w_2} = a_1 \delta_2$$

$$J = \frac{1}{2} (\hat{y} - y)^2$$

$$J = \frac{1}{2} (y - a_3)^2$$

$$e_2 = \frac{\partial J}{\partial a_2} = (\hat{y} - a_3) = (y - \hat{y})$$

$$\begin{pmatrix} \frac{\partial J}{\partial q_1} & \frac{\partial J}{\partial q_2} & \frac{\partial J}{\partial q_3} \\ \frac{\partial J}{\partial z_1} & \frac{\partial J}{\partial z_2} & \frac{\partial J}{\partial z_3} \end{pmatrix}$$



$$z = \underline{w_1 a_1 + w_2 a_2 + w_3 a_3}$$

$$\leftarrow \delta$$

$$\frac{\partial J}{\partial w_2} \doteq \delta x_3$$

Full algorithm (scalar prod)

i) Do the forward pass

$$q_0 = x; \rightarrow z_1, q_1, , z_2, q_2, , z_3, q_3 .. z_i, q_i$$

$$\hat{y} \rightarrow e = (y - \hat{y}) = e_3$$

Back prop

$$e^{c(k-1)} = w_k \underline{\delta_k},$$

$$\rightarrow \Delta w = \frac{q_i \cdot \delta_i}{e^{ch})}$$

$$D \xrightarrow{w^{ch}} D \leftarrow s^{ch}$$

3) Continue to calculate $s^{(k)}$ and $e^{(k)}$ till you reach the first layer

$$4) \quad \Delta w_i = q_i \delta_i \quad \frac{\partial J}{\partial w_i} =$$

$\frac{\partial J}{\partial w_i} \rightarrow$ gradient to update weights

why backpropagation?

If you use finite difference

$$\frac{\partial J}{\partial w_i} = \frac{J(w_1 + \Delta w_i, w_2, w_3) - J(w_1, w_2, w_3)}{\Delta w_i}$$

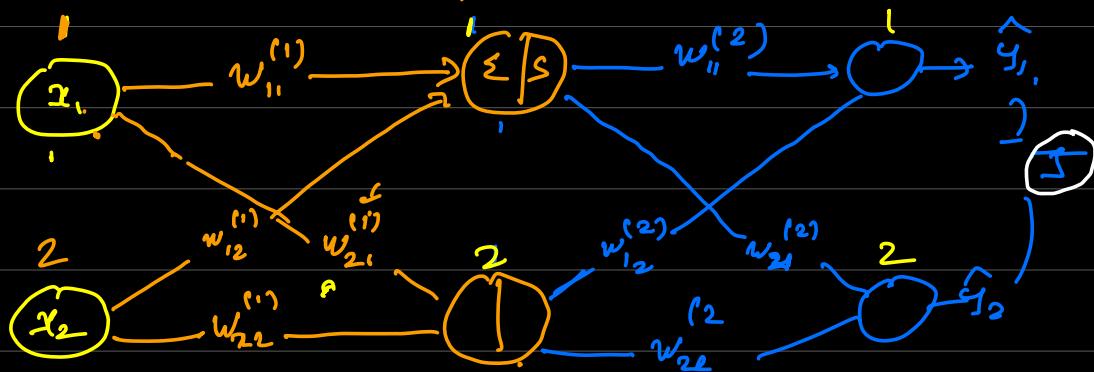
$\xrightarrow{\text{one forward}}$ $\xrightarrow{\text{one backward}}$

if passes forward-passes

MULTI-LAYER

MULTI-NEURON CASE

HIDDEN



(l-1) layer

$w_{i,j}$
to
from

$\Sigma \rightarrow$ scalar
 $(q_1 - \hat{q}_1)^2 + (q_2 - \hat{q}_2)^2$

FORWARD PASS: $\vec{z}_1^{(1)}$
 $\vec{z}_2^{(1)}$

$$\vec{z}_1^{(1)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2$$

$$\vec{z}_2^{(1)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2$$

$$\begin{bmatrix} \vec{z}_1^{(1)} \\ \vec{z}_2^{(1)} \end{bmatrix} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \left| \begin{array}{l} a_1^{(1)} = g(\vec{z}_1^{(1)}) \\ a_2^{(1)} = g(\vec{z}_2^{(1)}) \\ a = g(z) \\ \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = g\left(\begin{bmatrix} \vec{z}_1^{(1)} \\ \vec{z}_2^{(1)} \end{bmatrix}\right) \end{array} \right.$$

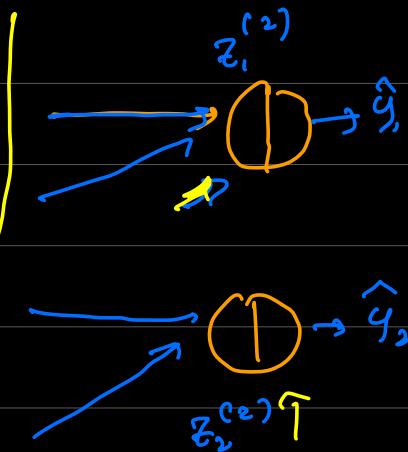
$$(\vec{z}^{(1)} = W^{(1)} x)$$

$$2 \times 1 \quad 2 \times 2 \quad 2 \times 1$$

Final layer

$$\frac{\partial J}{\partial \hat{y}_1} = e_1^{(2)} = y - \hat{y}_1$$

$$\frac{\partial J}{\partial \hat{y}_2} = e_2^{(2)} = y - \hat{y}_2$$



$$\delta_1^{(2)} = \frac{\partial J}{\partial z_1^{(2)}} = \frac{\partial J}{\partial a_1} \frac{\partial a_1}{\partial z_1} = g'(z_1) \cdot e_1$$

$$\delta_2^{(2)} = \frac{\partial J}{\partial z_2^{(2)}} = \frac{\partial J}{\partial a_2} \frac{\partial a_2}{\partial z_2} = g'(z_2) \cdot e_2$$

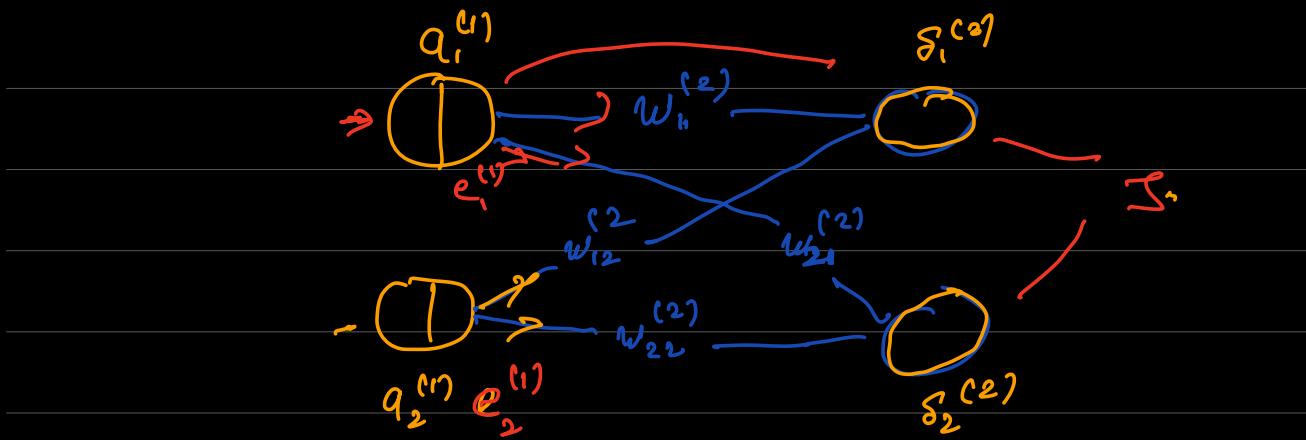
$$\delta^{(2)} = g'(z^{(2)}) \odot e$$

\odot → hadamard product

$$\begin{pmatrix} \delta_1^{(2)} \\ \delta_2^{(2)} \end{pmatrix} = \underbrace{\begin{pmatrix} g'(z_1^{(2)}) \\ g'(z_2^{(2)}) \end{pmatrix}}_{\sim} \cdot \begin{pmatrix} e_1^{(2)} \\ e_2^{(2)} \end{pmatrix}$$

$$c = a \odot b$$

$$c_i = a_i \cdot b_i$$



$$e_1^{(1)} = \frac{\partial J}{\partial q_1^{(1)}} = \underbrace{\frac{\partial J}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial q_1^{(1)}}}_{\text{derivative of loss}} + \frac{\partial J}{\partial z_2^{(2)}} \cdot \frac{\partial z_2^{(2)}}{\partial q_1^{(1)}}$$

$$e_1^{(1)} = \delta_1^{(2)} w_{11}^{(2)} + \delta_2^{(2)} w_{21}^{(2)}$$

$$e_2^{(1)} = \frac{\partial J}{\partial q_2^{(1)}} = \frac{\partial J}{\partial z_1} \frac{\partial z_1}{\partial q_2} + \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial q_2}$$

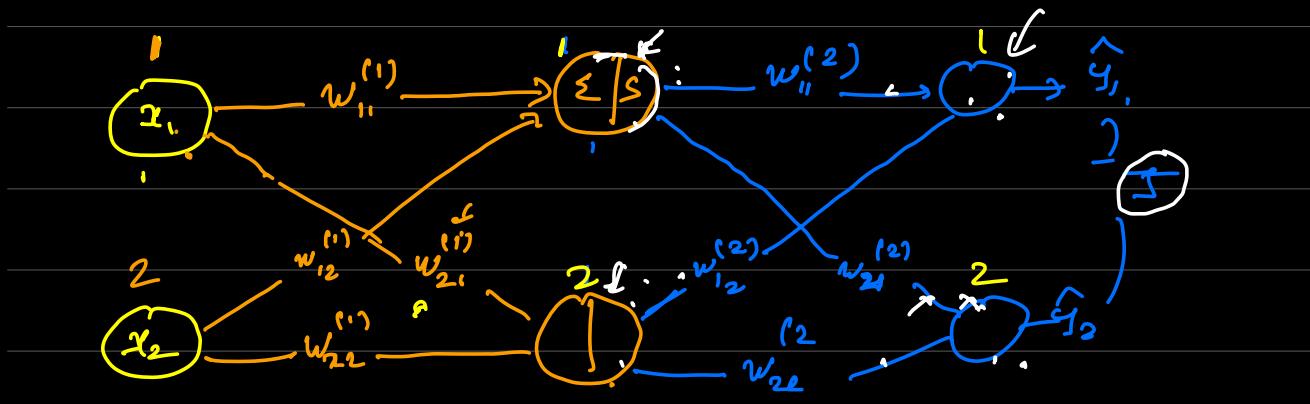
$$e_2^{(1)} = \delta_1^{(2)} w_{12}^{(2)} + \delta_2^{(2)} w_{22}^{(2)}$$

$$\begin{bmatrix} e_1^{(1)} \\ e_2^{(1)} \end{bmatrix} = \begin{bmatrix} w_{11}^{(2)} & w_{21}^{(2)} \\ w_{12}^{(2)} & w_{22}^{(2)} \end{bmatrix} \begin{bmatrix} \delta_1^{(2)} \\ \delta_2^{(2)} \end{bmatrix}$$

$$e^{(1)} = W^T (W^{(2)})^T \cdot \delta^{(2)}$$

$$e^{(1)} = (W^{(2)})^T \delta^{(2)}$$

$$f^{(2)} = g'(z^{(2)}) \odot e^{(2)}$$



from hidden to weight update



① Forward Pass $z^{k+1} = (w^{(k)} a^{(k)})$
 $a^{(k+1)} = g(z^{(k)})$

② Calculate error vector at final layer $\rightarrow e$

③ $\delta^{(k)} = g'(z^{(k)}) \odot e^{(k)}$
 $e^{(k+1)} = (w^{(k)})^T \delta^{(k)}$

Continue to first layer

④ update weights

$$\Delta w_{ij} = \delta_i a_j \quad (\text{Automatic Differentiation})$$