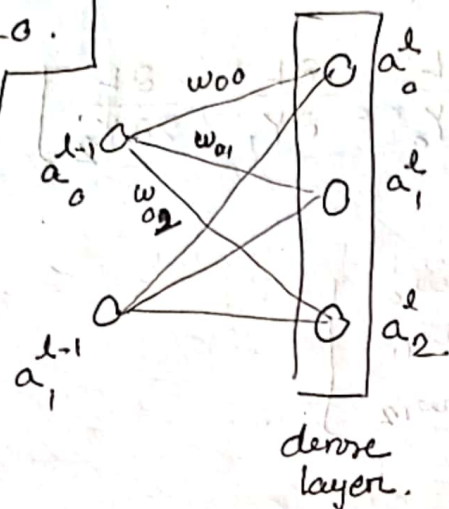$a^l_0$ = activation for layer $l$ and node $0$.

$w_{ij}$ = weight of edge between $i$th node of previous layer and $j$th node of current layer.


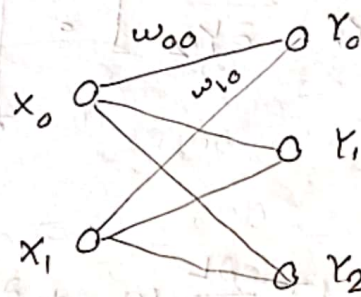
denre layer.

output mize = 3
input mize = 2

weight matrix = $\begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix}$

bias matrix = $\begin{bmatrix} b_0 & b_1 & b_2 \end{bmatrix}$

we can interprete like this:



$Y_0 = w_{00} x_0 + w_{10} x_1 + b_0$

$Y_1 = w_{01} x_0 + w_{11} x_1 + b_1$

$Y_2 = w_{02} x_0 + w_{12} x_1 + b_2$

$$[Y_0\ Y_1\ Y_2] = [x_0\ x_1] \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix} + [b_0\ b_1\ b_2]$$

$$\frac{dY}{dX} = w^T$$

$$\text{grad}_{\partial}\text{ output} = \frac{\partial L}{\partial Y} = \left[\frac{\partial L}{\partial Y_o} \, \& \, \frac{\partial L}{\partial Y_1} \quad \frac{\partial L}{\partial Y_2}\right]$$

$$W^T = \begin{bmatrix} W_{00} & W_{10} \\ W_{01} & W_{11} \\ W_{02} & W_{12} \end{bmatrix}$$

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \cdot W^T$$

$$= \left[\frac{\partial L}{\partial Y_o} \quad \frac{\partial L}{\partial Y_1} \quad \frac{\partial L}{\partial Y_2}\right] \cdot \begin{bmatrix} W_{00} & W_{10} \\ W_{01} & W_{11} \\ W_{02} & W_{12} \end{bmatrix}$$

$$= \left[\frac{\partial L}{\partial Y_o} \times W_{00} + \frac{\partial L}{\partial Y_1} W_{01} + \frac{\partial L}{\partial Y_2} W_{02} \quad \frac{\partial L}{\partial Y_o} \times W_{10} \right]$$

$$= \left[\frac{\partial L}{\partial X_o} \quad \frac{\partial L}{\partial X_1}\right]$$

So, grad input $= \dfrac{\partial L}{\partial X_o} = \dfrac{\partial L}{\partial Y_o} \times \left(\dfrac{\partial Y_o}{\partial X_o}\right) W_{00} + \dfrac{\partial L}{\partial Y_1}\left(\dfrac{\partial Y_1}{\partial X_o}\right) W_{01}$

$+ \dfrac{\partial L}{\partial Y_2}\left(\dfrac{\partial Y_2}{\partial X_o}\right) W_{02}$

thus, this is correct. $\dfrac{\partial L}{\partial Y_o} \cdot W_{00}$

[effect of everyone will be added]

$$\therefore \text{ grad input} = \left[ \frac{\partial L}{\partial x_0} \quad \frac{\partial L}{\partial x_1} \right]$$

this value will be
further back propagated.

in the same way,

$$\text{grad weight}, \quad \frac{\partial L}{\partial w} = X^T \frac{\partial L}{\partial r}$$

$$= \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} \begin{bmatrix} \dfrac{\partial L}{\partial r_0} & \dfrac{\partial L}{\partial r_1} & \dfrac{\partial L}{\partial r_2} \end{bmatrix}$$

$$= \begin{bmatrix} x_0 \dfrac{\partial L}{\partial r_0} & x_0 \dfrac{\partial L}{\partial r_1} & x_0 \dfrac{\partial L}{\partial r_2} \\[2mm] x_1 \dfrac{\partial L}{\partial r_0} & x_1 \dfrac{\partial L}{\partial r_1} & x_1 \dfrac{\partial L}{\partial r_2} \end{bmatrix}$$

for the first element
of the matrix

$$r_0 = w_{00} x_0 + b.$$

$$\therefore \frac{\partial r_0}{\partial w_{00}} = x_0.$$

$$\therefore \boxed{\frac{\partial L}{\partial w_{00}}} = \frac{\partial r_0}{\partial w_{00}} \cdot \frac{\partial L}{\partial r_0}$$

$$= \boxed{x_0 \frac{\partial L}{\partial r_0}}$$

$$\frac{\partial L}{\partial b} = \left[ \frac{\partial L}{\partial b_0} \quad \frac{\partial L}{\partial b_1} \quad \frac{\partial L}{\partial b_2} \right]$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial Y} \cdot \frac{\partial Y}{\partial b} = \frac{\partial L}{\partial Y}.$$

two
∴ for one out samples

$$\frac{\partial L}{\partial Y} = \begin{bmatrix} \frac{\partial L}{\partial Y_{00}} & \frac{\partial L}{\partial Y_{10}} & \frac{\partial L}{\partial Y_{20}} \\ \frac{\partial L}{\partial Y_{01}} & \frac{\partial L}{\partial Y_{11}} & \frac{\partial L}{\partial Y_{21}} \end{bmatrix}$$

$Y_{i3}$ = output
of $i^{th}$
node
for
$3^{th}$
sample

$$\frac{\partial L}{\partial b_0} = \frac{\partial L}{\partial Y_{00}} + \frac{\partial L}{\partial Y_{01}} \qquad \text{sum over all samples.}$$