

Neural Networks:-

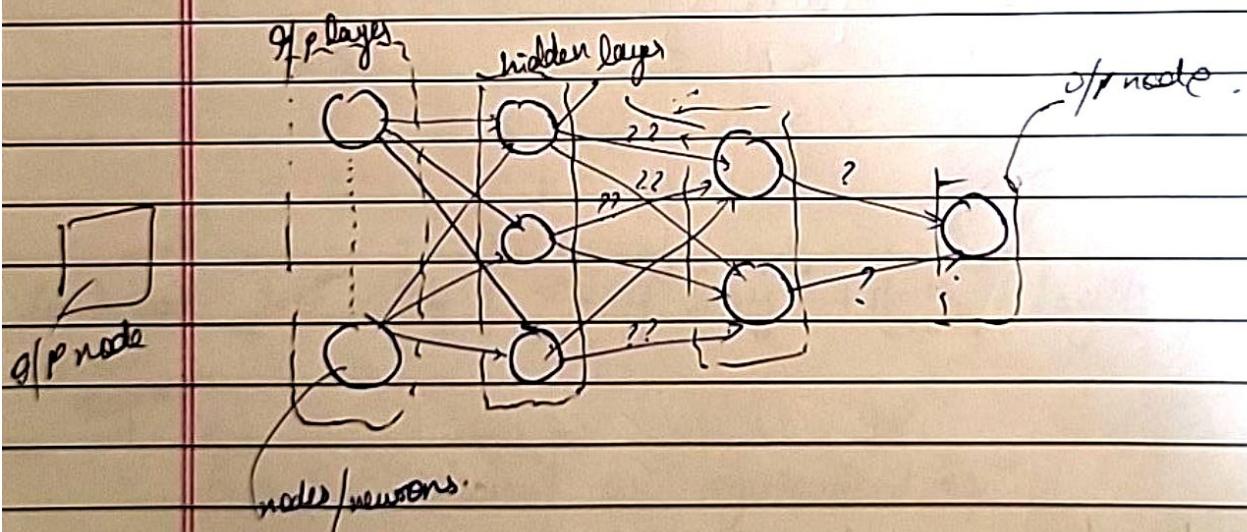
gradient descent here is termed as back propagation.

- every M.L model be it regression, classification, neural n/w etc will give probability not the exact value.

→ we can put as many hidden layers as we want in neural n/w.

→ neural n/w with 1 or 2 hidden layer is shallow neural n/w

→ neural n/w with more than 2 hidden layer is deep neural n/w



we can feed the i/p as a matrix converted into 1-d array (each row wise) at each node of the i/p layer.

$$\text{Size of i/p layer} = \text{Size of i/p (matrix)}$$

Computation Graph

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

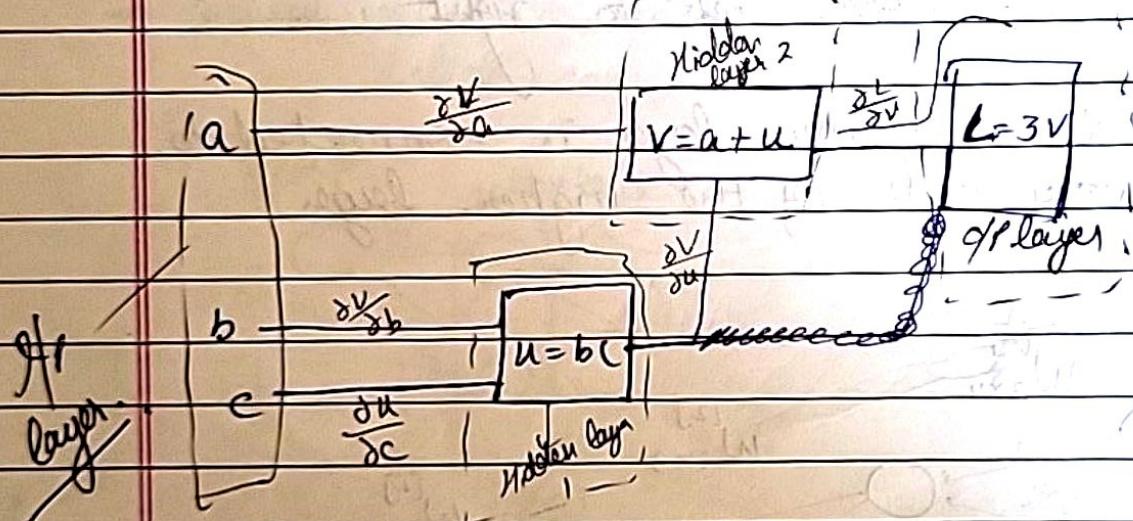
Calculate $\frac{\partial L}{\partial a}$, $\frac{\partial L}{\partial b}$, $\frac{\partial L}{\partial c}$??

$$\text{Let } (a + bc) = v.$$

$$\text{Let } bc = u.$$

- Direct calc = $\frac{\partial L}{\partial a} = 3$, $\frac{\partial L}{\partial b} = 3c$, $\frac{\partial L}{\partial c} = 3b$

- Computational Graph method :-



$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial v} \times \frac{\partial v}{\partial a} = 3 \cdot 1 = 3$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial b} = 3 \cdot 1 \cdot c = 3c$$

$$\frac{\partial L}{\partial c} = \frac{\partial L}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial c} = 3 \cdot 1 \cdot b = 3b$$

$v + w$ uc ud

$$\text{eg } f = (a+b)c + (a+b)d$$

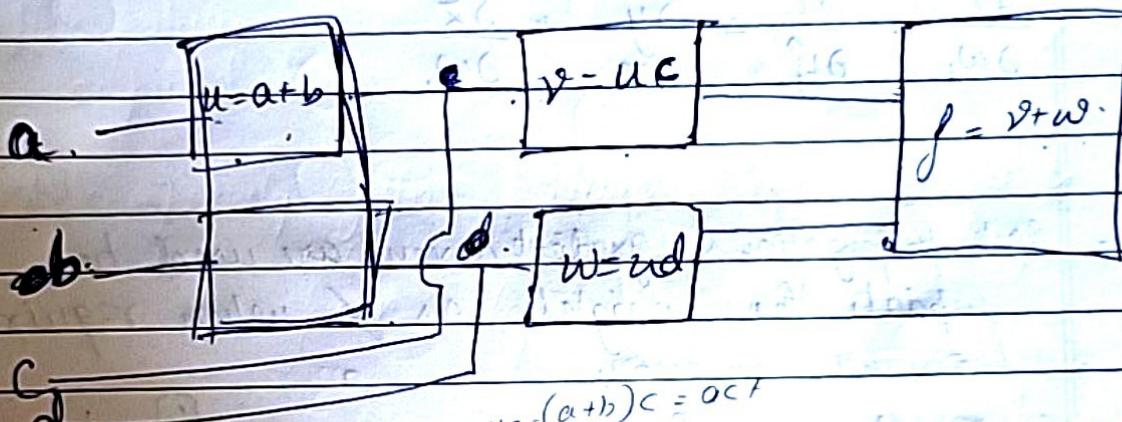
Calculate $\frac{\partial f}{\partial a}$, $\frac{\partial f}{\partial b}$, $\frac{\partial f}{\partial c}$, $\frac{\partial f}{\partial d}$ - ?

Let $a+b = u$.

Let $uc = v$

Let $ud = w$.

Let $v+w = x$.



$$uc = (a+b)c = \text{oct}$$

$$\begin{aligned} \frac{\partial f}{\partial a} &= \frac{\partial f}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial a} + \frac{\partial f}{\partial w} \times \frac{\partial w}{\partial u} \times \frac{\partial u}{\partial a} \\ &= 1 \times c * 1 + 1 \times 0 * 1 = c+d. \end{aligned}$$

$$\frac{\partial f}{\partial d} = \frac{\partial f}{\partial z} \times \frac{\partial z}{\partial d} = 1 \cdot v = a+b.$$

$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial u} \times \frac{\partial u}{\partial c} = 1 \cdot v = a+b$$

Acc to Sir :-

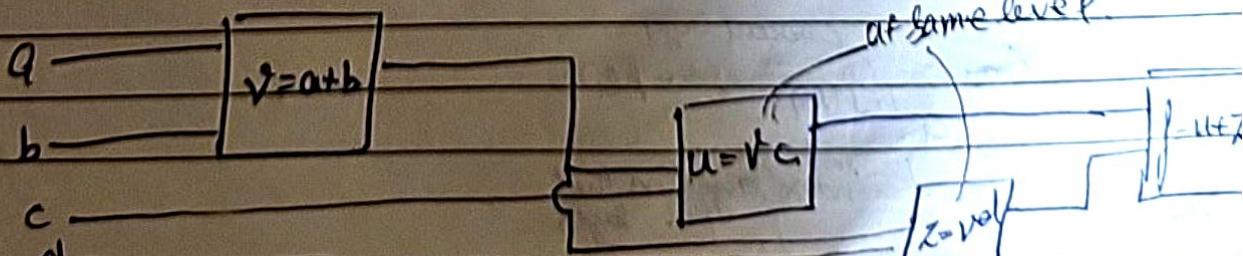
$$a+b = v.$$

$$vc = u.$$

$$vd = x.$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial u} \times \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial a} + \frac{\partial f}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial a} = c+d$$

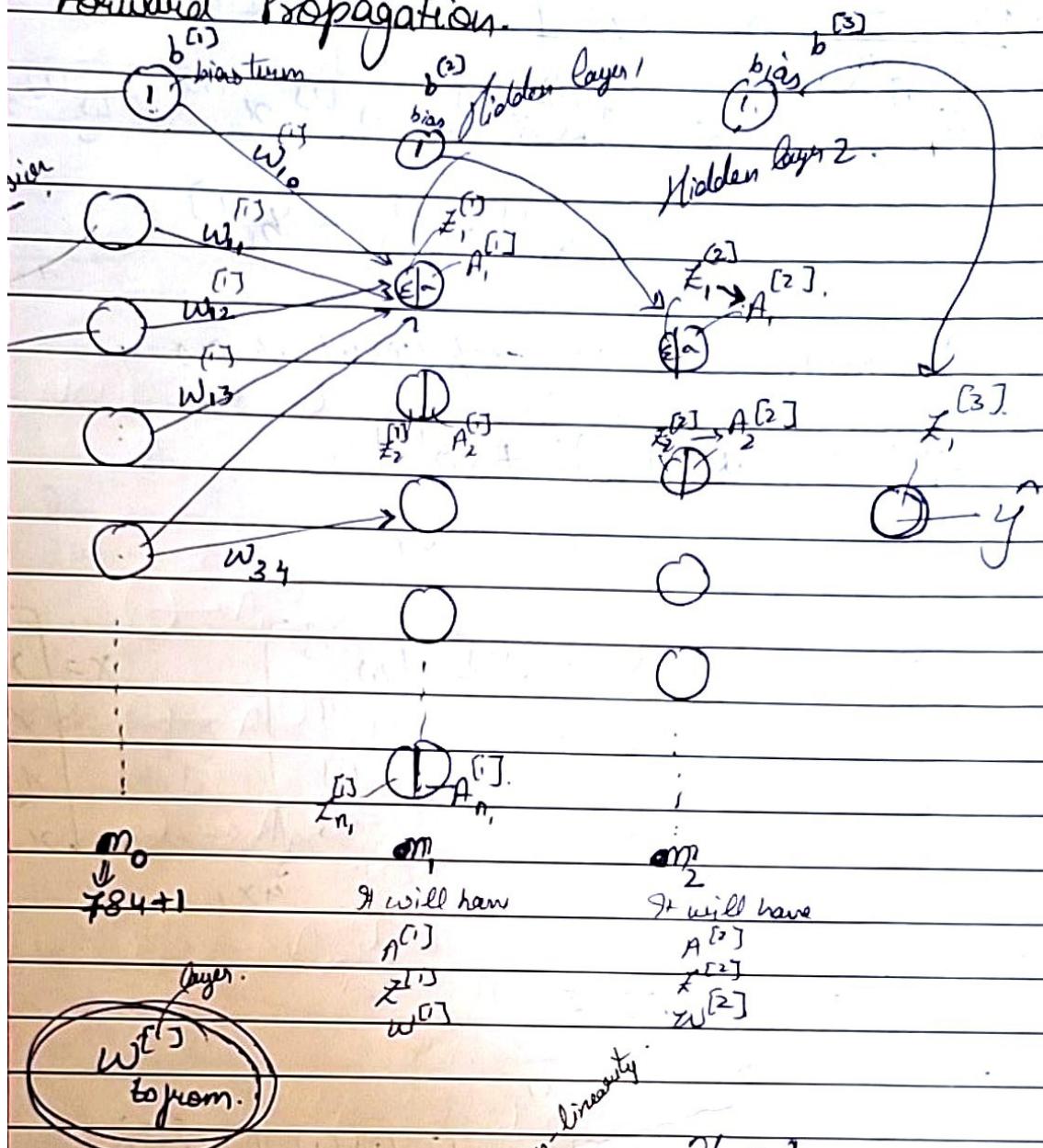
$$\begin{aligned} \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial z} \times \frac{\partial z}{\partial v} \times \frac{\partial v}{\partial b} + \frac{\partial f}{\partial v} \times \frac{\partial v}{\partial u} \times \frac{\partial u}{\partial b} \\ &= 1 \times d * 1 + 1 \times c * 1 = c+d \end{aligned}$$



$$28 \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \xrightarrow{1-D} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

RANKA
DATE / /
PAGE

Forward Propagation.



It will have $A^{(1)}$

$Z^{(1)}$

$w^{(1)}$

It will have $A^{(2)}$

$Z^{(2)}$

$w^{(2)}$

m_0

$\downarrow 784 + 1$

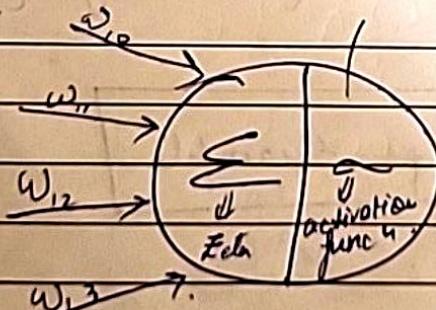
Layer.
 $w^{(1)}$
to from.

Non-linearity

$w_0 = 1$

So $w_0 \times w_0$ = bias term.

for each layer (i.e.)
a bias term w_0 will
be added
where the weight
will change
(The bias term
will be used in the
next layer)



$$\Sigma = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + w_0 (x_0)$$

$$\alpha \left(\sum_{j=0}^{m_0} w_j x_j \right)$$

DATE / / PAGE / /

$$w \in R^{[1]} \times [m_1 \times (m_0 + 1)]$$

$$\rightarrow \begin{bmatrix} w_{1,0}^{[1]} & w_{1,1}^{[1]} & \dots & w_{1,m_0}^{[1]} \\ w_{2,0}^{[1]} & w_{2,1}^{[1]} & \dots & w_{2,m_0}^{[1]} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m_1,0}^{[1]} & w_{m_1,1}^{[1]} & \dots & w_{m_1,m_0}^{[1]} \end{bmatrix} = w^{[1]}$$

neurons
of first hidden layer

$z_i^{[1]} = 1^{\text{st}} \text{ neuron of } 1^{\text{st}} \text{ hidden layer ke saare parameters}$

$$= x_1 \cdot w_{1,1}^{[1]} + x_2 \cdot w_{1,2}^{[1]} + x_3 \cdot w_{1,3}^{[1]} + \dots + x_{m_0} \cdot w_{1,m_0}^{[1]}$$

$$z_i^{[1]} = \sum_{j=1}^{m_0} w_{1,j}^{[1]} x_j + b_1^{[1]}$$

$$z_3^{[1]} = \sum_{j=1}^{m_0} w_{3,j}^{[1]} x_j + b_3^{[1]}$$

So generalized k^{th} neuron of 1^{st} hidden layer ke saare parameters

$$z_k^{[1]} = \sum_{j=1}^{m_0} x_j \cdot w_{k,j}^{[1]} + b_k^{[1]}$$

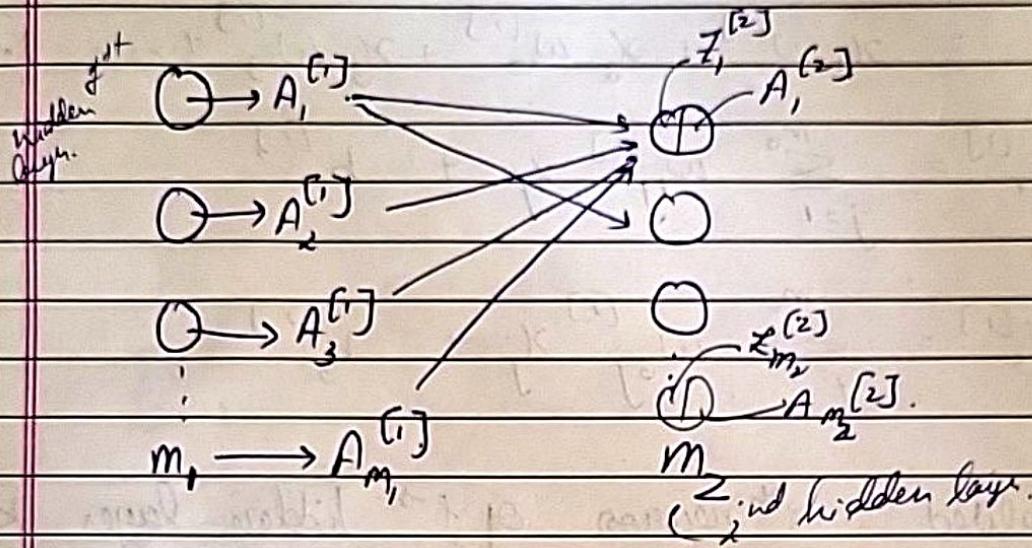
For multiple training samples or image.

$$1^{\text{st}} \text{ image} = z_i^{1} = \sum_{j=1}^{m_0} x_j^{(1)} w_{j,i}^{[1]} - b_i^{[1]}$$

$$2^{\text{nd}} \text{ image} = z_i^{[1](2)} = \sum_{j=1}^{m_0} x_j^{(2)} w_{j,i}^{[1]} + b_i^{[1]}$$

so $z_k^{[1](i)} = 1^{\text{st}} \text{ hidden layer ke } k^{\text{th}} \text{ neuron of } i^{\text{th}} \text{ image/example}$

$$= \sum_{j=1}^{m_0} x_j^{(i)} w_{k,j}^{[1]} + b_k^{[1]}$$



$$\text{So } Z_1^{[2]} = w_{11}^{[2]} A_1^{[1]} + w_{12}^{[2]} A_2^{[1]} + w_{13}^{[2]} A_3^{[1]} + \dots + w_{1m_1}^{[2]} A_{m_1}^{[1]} + b^{[2]}.$$

$$Z_2^{[2]} = w_{21}^{[2]} A_1^{[1]} + w_{22}^{[2]} A_2^{[1]} + \dots + w_{2m_1}^{[2]} A_{m_1}^{[1]} + b^{[2]}.$$

$$\text{So generalized } Z_K^{[2]} = \sum_{j=1}^{m_1} w_{kj}^{[2]} A_j^{[1]} + b_K^{[2]}$$

$$\text{or multiply in matrix form } Z_K^{(2)(i)} = \sum_{j=1}^{m_1} w_{kj}^{[2]} A_j^{(1)(i)} + b_K^{[2]}$$

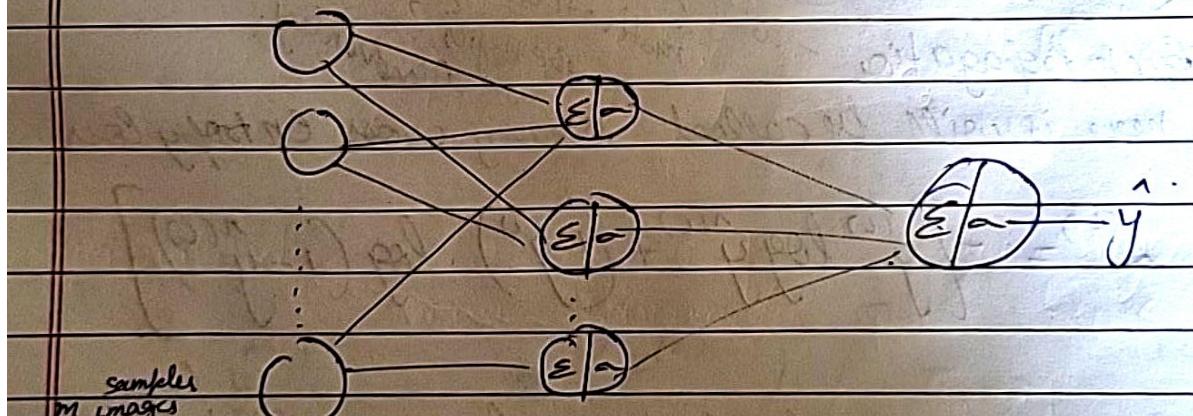
$$A_k^{[2](i)} = \alpha(Z_k^{[2](i)})$$

for Output or 3rd Layer

$$Z_p^{[3](i)} = \sum_{j=1}^m w_{pj}^{[3]} A_j^{[2](i)} + b_p^{[3]}$$

$$A_p^{[3](i)} = \alpha(Z_p^{[3](i)}) = \hat{y}$$

unvectorised (1 eg by 1 eg.)



$$\text{samples } m \text{ images} \quad x^{(i)} \in \mathbb{R}^{n^{[0]} \times 1}$$

$$w^{[1]} \in \mathbb{R}^{n^{[0]} \times n^{[1]}}$$

$$b^{[1]} \in \mathbb{R}^{n^{[1]} \times 1}$$

$$\text{back in range}(1, n): w^{[2]} \in \mathbb{R}^{n^{[1]} \times n^{[2]}}$$

$$b^{[2]} \in \mathbb{R}^{n^{[2]} \times 1}$$

$i \in \text{range}(1, m)$:

$$z^{[1](i)} = \underbrace{(w^{[1]} x^{(i)})}_{(n^{[0]} \times n^{[1]}) (n^{[0]} \times 1)} + \underbrace{b^{[1]}}_{n^{[1]} \times 1}$$

any non linear activation func.

$$A = \begin{bmatrix} A^{(1)(i)} \\ A^{(2)(i)} \\ \vdots \\ A^{(m)(i)} \end{bmatrix}_{n^{(1)} \times 1} \xrightarrow{\text{f}(z(i))} g(z^{(1)(i)})_{n^{(2)} \times 1}$$

So now $z^{(2)(i)} = w^{[2]} A^{[1](i)} + b^{[2]}_{1 \times 1}$

$$A^{[2](i)}_{1 \times 1} = g(z^{[2](i)}) = y^{(i)}$$

Back propagation - why? motivation??
benefit? disadvantage?
here it will be called Binary cross entropy loss

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

here we are calculating loss image by image.

$$L = \frac{1}{m} \sum_{i=1}^m \lambda^{(i)}$$

So here we have to calculate.

$$w^{[1]} = w^{[1]} - \alpha \frac{\partial L}{\partial w^{[1]}}$$

$$b^{[1]} = \frac{\partial L}{\partial b^{[1]}}$$

$$w^{[2]} = w^{[2]} - \alpha \frac{\partial L}{\partial w^{[2]}}$$

$$b^{[2]} = \frac{\partial L}{\partial b^{[2]}}$$

$w^{[2]}$ & $b^{[2]}$ are nearest to the o/p so we will calculate them first.

$$\frac{\partial L^{(l)}}{\partial w^{[2]}} = \frac{\partial L^{(l)}}{\partial \hat{y}^{(i)}} \cdot \frac{\partial \hat{y}^{(i)}}{\partial z^{[2](i)}} \cdot \frac{\partial z^{[2](i)}}{\partial w^{[2]}}$$

$$\text{So } \frac{\partial L^{(l)}}{\partial \hat{y}^{(i)}} = - \left[\hat{y}^{(i)} \times \frac{1}{\hat{y}^{(i)}} + (1 - \hat{y}^{(i)}) \times \frac{1}{1 - \hat{y}^{(i)}} (-1) \right]$$

$$= \left[\frac{-\hat{y}^{(i)}}{\hat{y}^{(i)}} + \frac{1 - \hat{y}^{(i)}}{1 - \hat{y}^{(i)}} \right]$$

$$= -\hat{y}^{(i)} + \frac{\hat{y}^{(i)} \hat{y}^{(i)}}{\hat{y}^{(i)} (1 - \hat{y}^{(i)})} + \hat{y}^{(i)} - \frac{\hat{y}^{(i)} \hat{y}^{(i)}}{\hat{y}^{(i)} (1 - \hat{y}^{(i)})}$$

$$= \cancel{\frac{\hat{y}^{(i)}}{\hat{y}^{(i)}}} - \frac{\hat{y}^{(i)}}{1 - \hat{y}^{(i)}}$$

$$\frac{\partial \hat{y}^{(i)}}{\partial z^{[2](i)}} = \alpha(z^{[2](i)}) \cdot \left[1 - \alpha(z^{[2](i)}) \right]$$

$$\frac{\partial z^{[2](i)}}{\partial w^{[2]}} = A^{(1)(i)}$$

$$\begin{aligned} \frac{\partial L^{(l)}}{\partial z^{[2](i)}} &= \frac{\partial L^{(l)}}{\partial \hat{y}^{(i)}} \times \frac{\partial \hat{y}^{(i)}}{\partial z^{[2](i)}} \\ &= c \cancel{A^{(1)(i)}} (\text{notatio}) \end{aligned}$$

$$\text{So. } \frac{\partial L^{(l)}}{\partial w^{[2]}} = \frac{\hat{y}^{(i)} - \hat{y}^{(i)}}{\hat{y}^{(i)} (1 - \hat{y}^{(i)})} * \cancel{\hat{y}^{(i)} (1 - \hat{y}^{(i)})} * A^{(1)(i)}$$

$$= \cancel{\left(\frac{\hat{y}^{(i)} - \hat{y}^{(i)}}{\hat{y}^{(i)} (1 - \hat{y}^{(i)})} \right)} * A^{(1)(i)}$$

but $w^{[2]} \in \mathbb{R}^{1 \times n^{(i)}}$ so we need to take $(A^{(1)(i)})^T$

$$\frac{\partial L^{(i)}}{\partial w^{(2)}} = \left(\hat{y}^{(i)} - y^{(i)} \right) \left(A^{(1)(i)} \right)^T$$

$$w^{(2)(\text{new})} = w^{(2)(\text{old})} - \alpha \left(\hat{y}^{(i)} - y^{(i)} \right) \left(A^{(1)(i)} \right)^T$$

$$\Rightarrow \frac{\partial L^{(i)}}{\partial b^{(2)}} = \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} * \frac{\partial \hat{y}^{(i)}}{\partial z^{(2)(i)}} * \frac{\partial z^{(2)(i)}}{\partial b^{(2)}}$$

$$\frac{\partial b^{(2)}}{\partial L^{(i)}} = \frac{\partial L^{(i)}}{\partial b^{(2)}} = d(\cancel{L^{(2)}})$$

$$\frac{\partial w^{(2)}}{\partial L^{(i)}} = \frac{\partial L^{(i)}}{\partial w^{(1)}} * \frac{\partial \hat{y}^{(i)}}{\partial z^{(2)(i)}} * \frac{\partial z^{(2)(i)}}{\partial A^{(1)(i)}} * \frac{\partial A^{(1)(i)}}{\partial z^{(1)(i)}} * \frac{\partial z^{(1)(i)}}{\partial w^{(1)}}$$

$$\frac{\partial Z^{[2](c)}}{\partial A^{[1](c)}} = \omega^{[2]}$$

$$\frac{\partial A^{[1](c)}}{\partial Z^{[1](c)}} = g'(Z^{[1](c)})$$

$$\frac{\partial Z^{[1](c)}}{\partial \omega^{[1]}} \leftarrow x^{(c)}$$

$$\delta_0 \frac{\partial L^{(i)}}{\partial w^{[1]}} = \underbrace{\begin{pmatrix} \hat{y}^{(i)} - y^{(i)} \\ \hat{y}^{(i)}(1 - \hat{y}^{(i)}) \end{pmatrix}}_{1 \times n^{[1]}} * \underbrace{\hat{y}^{(i)}(1 - \hat{y}^{(i)})}_{1 \times n^{[1]}} * \underbrace{w^{[2]} * g'(z^{[1](i)})}_{n^{[1]} \times 1} * \underbrace{x_i}_{1 \times 1}$$

we need to
take transpose

$$\delta_0 \frac{\partial L^{(i)}}{\partial w^{[1]}} = \left(\hat{y}^{(i)} - y^{(i)} \right) * w^{[2]} * g'(z^{[1](i)}) * \left(x^{(i)} \right)^T$$

$$\frac{\partial L^{(i)}}{\partial w^{[1]}} = \underbrace{\left(w^{[2]} \right)^T}_{n^{[2]} \times 1} \underbrace{\left(\hat{y}^{(i)} - y^{(i)} \right)}_{1 \times 1} \odot \underbrace{g'(z^{[1](i)})}_{\substack{\text{element} \\ \text{by} \\ \text{element} \\ \text{multiplication}}} \underbrace{\left(x^{(i)} \right)^T}_{1 \times n^{[0]}}$$

of the 2nd matrix

$$\Rightarrow \frac{\partial L^{(i)}}{\partial b^{[1]}} = \frac{\partial L^{(i)}}{\partial \hat{y}^{(i)}} * \frac{\partial \hat{y}^{(i)}}{\partial z^{[2](i)}} * \frac{\partial z^{[2](i)}}{\partial A^{[1][i]}} * \frac{\partial A^{[1][i]}}{\partial z^{[1](i)}} * \frac{\partial z^{[1](i)}}{\partial b^{[1]}}$$

$$\frac{\partial b^{[1]}}{\partial b^{[1]}} = \left(\hat{y}^{(i)} - y^{(i)} \right) * w^{[2]} * g'(z^{[1](i)}) * \underline{\underline{1}}$$

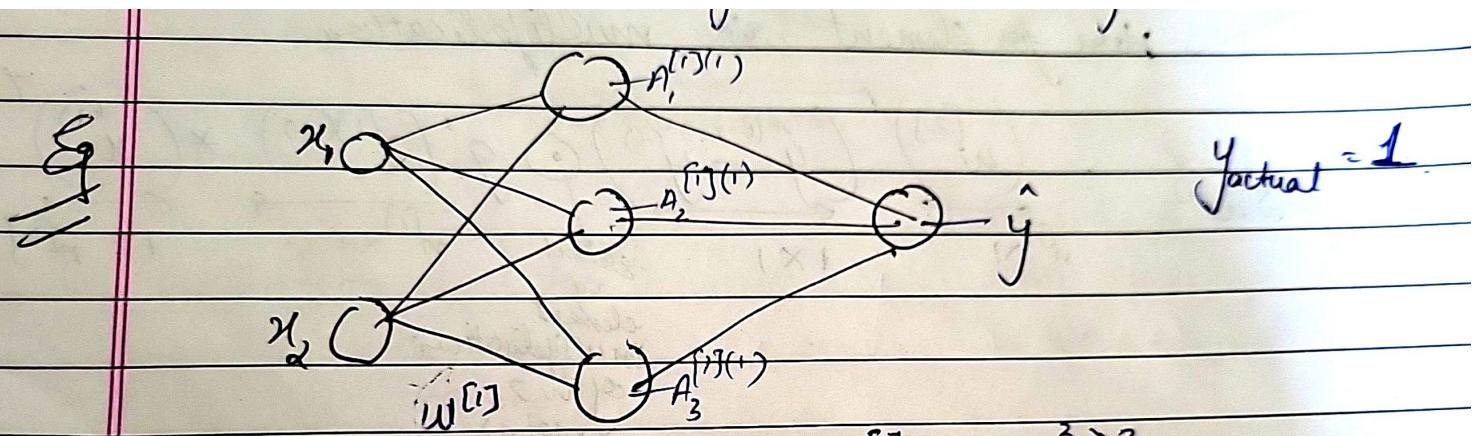
So Gradient update of the parameters in the neural network

$$w^{[2]} = w^{[2]} - \alpha d w^{[2]}$$

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

$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

$$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$$



$$x = \begin{bmatrix} 0.5 \\ 0.3 \end{bmatrix}$$

$$w^{[1]} \in \mathbb{R}^{3 \times 2}$$

$$b^{[1]} \in \mathbb{R}^{3 \times 1}$$

$$w^{[2]} \in \mathbb{R}^{1 \times 3}$$

$$b^{[2]} \in \mathbb{R}^{1 \times 1}$$

$$w^{[1]} = \begin{bmatrix} 0.2 & 0.4 \\ 0.1 & 0.3 \\ 0.5 & -0.2 \end{bmatrix} \text{ will be given in exam.}$$

$$b^{[1]} = \begin{bmatrix} 0.1 \\ -0.1 \\ 0.05 \end{bmatrix} \quad w^{[2]} = \begin{bmatrix} 0.3 & -0.2 & 0.4 \end{bmatrix}$$

$$b^{[2]} = \begin{bmatrix} -0.2 \end{bmatrix}$$

Target \rightarrow get new value of w from old values.

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

$$= \begin{bmatrix} 0.2 & 0.4 \\ 0.1 & 0.3 \\ 0.5 & -0.2 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.1 \\ 0.05 \end{bmatrix}$$

$$= \begin{bmatrix} 0.22 \\ 0.14 \\ 0.19 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.1 \\ 0.05 \end{bmatrix} = \begin{bmatrix} 0.32 \\ 0.04 \\ 0.24 \end{bmatrix}$$

$$z^{1} = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-0.32}} = 0.579$$

$$= \frac{1}{1+e^{-(0.04)}} = 0.510$$

$$= \frac{1}{1+e^{-(0.24)}} = 0.559 = 0.560$$

$$A^{1} = a(z^{1}) = a\begin{bmatrix} 0.32 \\ 0.04 \\ 0.24 \end{bmatrix} = \begin{bmatrix} a(0.32) \\ a(0.04) \\ a(0.24) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{1+e^{(-0.32)}} \\ \frac{1}{1+e^{-(0.04)}} \\ \frac{1}{1+e^{-(0.24)}} \end{bmatrix} = \begin{bmatrix} 0.579 \\ 0.510 \\ 0.559 = 0.560 \end{bmatrix}$$

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

$$= \begin{bmatrix} 0.3 & -0.2 & 0.4 \end{bmatrix}_{1 \times 3} \begin{bmatrix} 0.579 \\ 0.510 \\ 0.560 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} -0.2 \end{bmatrix}$$

$$= [1.737 + (-0.102) + 0.224]$$

$$= [0.2957] + [-0.2]$$

$$= 0.0957$$

$$A^{[2](1)} = \alpha (z^{[2](1)})$$

$$= \alpha(0.0957) = \frac{1}{1 + e^{-(0.0957)}}$$

$$\hat{y} = 0.0524 \quad (\text{it is not close to 1 (our goal)})$$

loss =

So now

$$\begin{aligned} dz^{[2](1)} &= A^{[2](1)} - y \\ &= 0.524 - 1 = (0.4476) \end{aligned}$$

$$\begin{aligned} dw^{[2]} &= d(z^{[2](1)}) (A^{1})^T \\ &= -0.4476 [0.579 \quad 0.510 \quad 0.560] \end{aligned}$$

$$\frac{\partial J^{(1)}}{\partial w^{[2]}} = dw^{[2]} = \begin{bmatrix} -0.275 & -0.243 & -0.267 \end{bmatrix}$$

$$d\tilde{w}^{(2)} = d\tilde{x}^{(2)(1)} = (-0.476)$$

$$d\tilde{w}^{(2)} = (\tilde{w}^{(2)})^T \cdot d\tilde{x}^{(2)(1)} \odot a'(x^{(1)(1)})(x^{(1)})^T$$

$$= \begin{bmatrix} 0.3 \\ -0.2 \\ 0.4 \end{bmatrix} (-0.476) \odot \begin{bmatrix} a(0.32)(1-a(0.32)) \\ a(0.04)(1-a(0.04)) \\ a(0.24)(1-a(0.24)) \end{bmatrix} \quad [0.5 \ 0.3]$$

$$= \begin{bmatrix} -0.01428 \\ 0.00952 \\ -0.01904 \end{bmatrix} \odot \begin{bmatrix} \frac{1}{1+e^{-(0.32)}} (1 - \frac{1}{1+e^{-(0.32)}}) \\ \frac{1}{1+e^{-(0.04)}} (1 - \frac{1}{1+e^{-(0.04)}}) \\ \frac{1}{1+e^{-(0.24)}} (1 - \frac{1}{1+e^{-(0.24)}}) \end{bmatrix} \quad [0.5 \ 0.3]$$

$$= \begin{bmatrix} -0.01428 \\ 0.00952 \\ -0.01904 \end{bmatrix} \odot \begin{bmatrix} 0.579 (1-0.579) \\ 0.510 (1-0.510) \\ 0.560 (1-0.560) \end{bmatrix} \quad [0.5 \ 0.3]$$

$$= \begin{bmatrix} -0.01428 \\ 0.00952 \\ -0.01904 \end{bmatrix} \odot \begin{bmatrix} 0.244 \\ 0.250 \\ 0.246 \end{bmatrix} \quad [0.5 \ 0.3]_{1 \times 2}$$

$$= \begin{bmatrix} -0.00348 \\ 0.0238 \\ -0.0468 \end{bmatrix} \quad [0.5 \ 0.3]_{1 \times 2}$$

$$= \begin{bmatrix} -0.0175 & 0.0105 \\ 0.0120 & 0.0072 \\ -0.0235 & 0.0141 \end{bmatrix}$$

$$\delta b^{[1]} = (\omega^{[2]})^T \left(\hat{y}^{[2]} - y^{[1]} \right) \odot g'(\boldsymbol{x}^{[1]})$$

$$= \begin{bmatrix} 0.3 \\ -0.2 \\ 0.4 \end{bmatrix} (-0.0476) \odot \alpha(\boldsymbol{x}^{[1]})$$

$$= \begin{bmatrix} -0.0343 \\ 0.0238 \\ -0.0468 \end{bmatrix}$$

$$\omega^{[1]} = \omega^{[1]} - \alpha \delta \omega^{[1]}$$