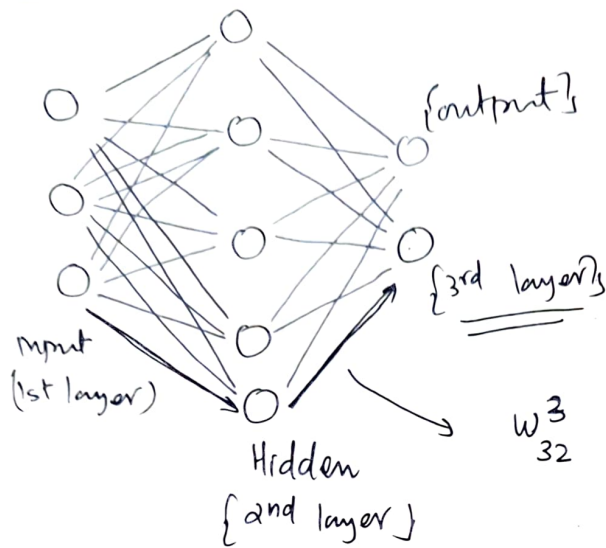


Notations :-



w_{jk}^z = z = final layer no.
 j = neuron no. in final layer
 k = neuron no. in initial layer.

b_j^l = bias in j -th neuron in l -th layer

a_j^l = activation of j -th neuron in l -th layer.

\Rightarrow Vectorizing activation

activation of $a_j^l = \sigma \left(\underbrace{\sum_k w_{jk}^l a_k^{l-1}}_{\text{dot product}} + b_j^l \right)$

this implies

$$\vec{a}^l = \sigma(\omega^l \vec{a}^{l-1} + \vec{b}^l)$$

Remember,

weight matrix is of form, w_{jk}^l $\begin{cases} j = \text{no. rows} \\ k = \text{no. columns} \end{cases}$

$$W = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1k} \\ A_{21} & A_{22} & \dots & A_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ A_{j1} & A_{j2} & \dots & A_{jk} \end{bmatrix}$$

$$\& \boxed{W \vec{a} + \vec{b}}$$

$$= \begin{bmatrix} A_{11} & \dots & A_{1k} \\ \vdots & \ddots & \vdots \\ A_{j1} & \dots & A_{jk} \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ \vdots \\ a_k \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_k \end{bmatrix}$$

Results in column matrix
 \Rightarrow a vector of activation.

Cost functions :

1. Quadratic cost :

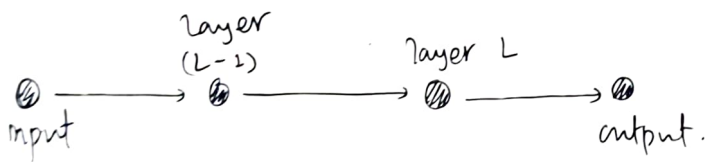
$$C = \frac{1}{2n} \sum_x \left\| \underbrace{y(x)}_{\text{"desired output"}} - \underbrace{a^L(x)}_{\text{"activation of last layer"}} \right\|^2 \quad \left. \vphantom{\sum_x} \right\} x = \text{no. of input.}$$

- defines how good/bad is our model performing with some weights & bias over a Example input.

• Other's later in hands on Experimentation.

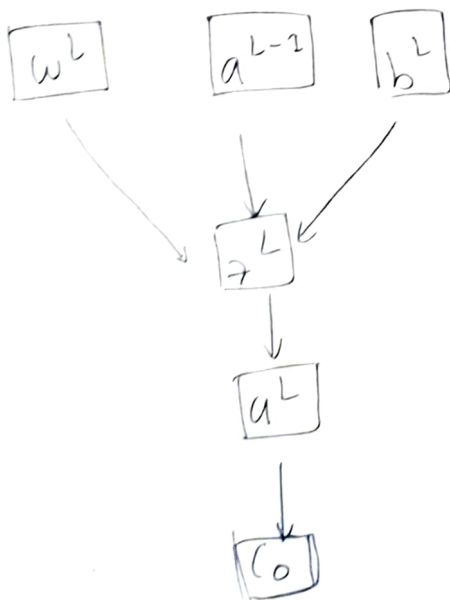
Forward pass & Back propagation :-

Imagine a simple network of neurons,



The cost of the final layer $C_0 = (a^L - y)^2 / 2$
activation $a^L = \sigma(z^L)$ $\left\{ \begin{array}{l} z^L = (w^L a^{L-1} + b^L) \end{array} \right.$

the computational graph is,



input with weight of $L-1$ -th layer gives a^{L-1} and that activation with subsequent weights & biases gives us z^L weighted input from last layer which again with activation function generates output.

- Cost C_0 is something we want to minimize,
and w 's & b 's are knobs we can turn & twist
- Lets look at how changes in w_L changes C_0

$$\frac{\partial C_0}{\partial w_L} = \frac{\partial}{\partial w_L} \left[\frac{(\sigma(z^L) - y)^2}{2} \right]$$

$$\frac{\partial C_0}{\partial w_L} = \frac{\partial z^L}{\partial w^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} \quad \left. \vphantom{\frac{\partial C_0}{\partial w_L}} \right\} \text{chain rule.}$$

$$\boxed{\frac{\partial C_0}{\partial w_L} = a^{L-1} \cdot \sigma'(z^L) \cdot (a^L - y)}$$

Similarly,

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial z^L}{\partial b^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L}$$

$$\boxed{\frac{\partial C_0}{\partial b^L} = w^L \sigma'(z^L) (a^L - y)}$$

C_0 depends on w_L & $b \Rightarrow \nabla C_0 = \left(\frac{\partial C}{\partial w^L}, \frac{\partial C}{\partial b} \right)^T$

gradient decent
step!!

For our simple network,

Total cost $\left| \frac{\partial C}{\partial w^L} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^L} \right|$

→ we can make
the errors smaller
& make correct-er prod.

Formal - 4 Equations of back propagation

definition : δ_j^L - "Error" is,

$$\delta_j^L \equiv \frac{\partial C}{\partial z_j^L} \quad \left. \vphantom{\frac{\partial C}{\partial z_j^L}} \right\} \text{Error Eq.}$$

Eq 1. : "Error in output layer"

$$\delta_j^L = \underbrace{\frac{\partial C}{\partial a_j^L}}_{\text{change of } C \text{ w.r.t activations}} \underbrace{\sigma'(z_j^L)}_{\text{change in } \sigma \text{ w.r.t } z_j^L}$$

These are the result of chain rule in
Error Eq.

→ converting it to matrix form,

$$\boxed{\delta_j^L = (a^L - y) \odot \sigma'(z^L)}$$

Hadamard
product.

Eq 2. Error δ^L as $\delta^L(\delta^{L+1})$:

$$\boxed{\delta^L = ((w^{L+1})^T \delta^{L+1}) \odot \sigma'(z^L)}$$

propagating error backwards.

Eq 3. Rate of change of cost wrt bias:

$$\boxed{\frac{\partial C}{\partial b_j^L} = \delta_j^L}$$

Eq 4. Rate of change of cost wrt weight:

$$\boxed{\frac{\partial C}{\partial w_{jk}^L} = a_k^{L-1} \delta_j^L}$$

↓
 $\{a \text{ in } \delta_{out}\}$

Back propagation Algorithm

- ① Input n : ^{get} activations in a^1
- ② Feed forward: for each layer compute
$$z^L = w^L a^{L-1} + b^L \quad \& \quad a^L = \sigma(z^L)$$
- ③ Error: $\delta^L = (a^L - y) \odot \sigma'(z^L) \rightarrow E_{n1}$.
- ④ Back propagate errors: $\delta^L = ((w^{L+1})^T \delta^{L+1}) \odot \sigma'(z^L)$
- ⑤ adjust weights & Biases:

Update Eq.

$$\boxed{\begin{aligned} w^L &\rightarrow w^L - \eta \delta^L (a^{L-1})^T \\ b^L &\rightarrow b^L - \eta \delta^L \end{aligned}}$$