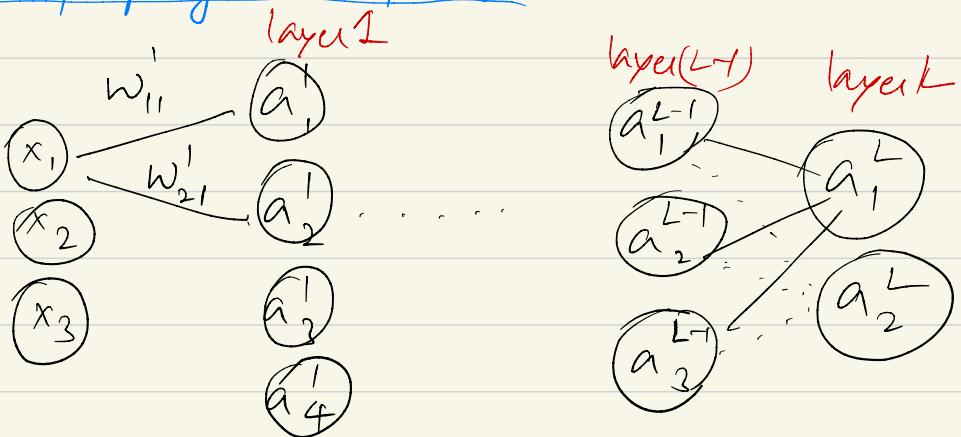


Back propagation proof (Intuition):



Example Network

what we know,

$$z_j^{[l]} = \sum_i w_{ji}^{[l]} \cdot a_i^{[l-1]} + b_j^{[l]}$$

$$\sigma(z_j^{[l]}) = a_j^{[l]}$$

$$(vectorized form) \quad z^{[l]} = \begin{matrix} [l] \\ \vdots \end{matrix} w^{[l]} \cdot \begin{matrix} [l-1] \\ \vdots \\ [1] \end{matrix} a^{[l-1]} + b^{[l]}$$

$$(n_e \times 1) \quad (n_e \times n_{e-1}) \quad (n_{e-1}) \quad (n_e \times 1)$$

$$A^{[l]} = \frac{\sigma(z^{[l]})}{(n_e \times 1)}$$

Cost Function : $C = \frac{1}{2} \sum (y_i^L - q_i^L)^2$

(Simple cost function is used for simpler illustration
of backprop Eq proof)

As we know,

$$\frac{\partial C}{\partial w_{ji}^L} = \frac{\partial C}{\partial z_j^L} \times \frac{\partial z_j^L}{\partial w_{ji}^L} \quad \text{and} \quad \frac{\partial C}{\partial b_j^L} = \frac{\partial C}{\partial z_j^L} \times \frac{\partial z_j^L}{\partial b_j^L}$$

so, once we find $\left(\frac{\partial C}{\partial z_j^L} \right)$ it will be piece of cake

we define, $\frac{\partial C}{\partial z_j^L} = \delta_j^L ; \frac{\partial C}{\partial z_j^l} = \delta_j^l$

what (δ_j^l) physically means ?

By how many times cost is changing, if we slightly change (z_j^l) .

$$\frac{\partial C}{\partial z_j^L} = \underbrace{\frac{\partial C}{\partial a_j^L}}_{\leftarrow} \cdot \underbrace{\frac{\partial a_j^L}{\partial z_j^L}}_{\downarrow}$$

$$= (a_j^L - y_j^L) \cdot \sigma'(z_j^L)$$

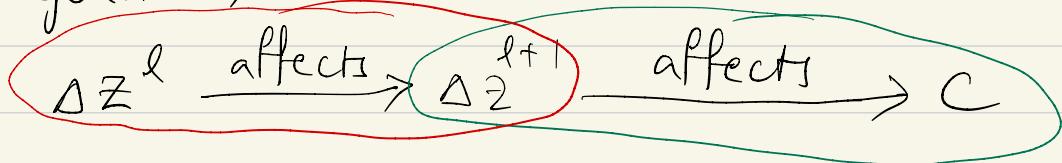
$$\boxed{\delta_j^L = (a_j^L - y_j^L) \cdot \sigma'(z_j^L)}$$

Now,

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}$$

As we know change in z_j^l affects all the neurons in $(l+1) \dots (L)$ we need to keep in mind that fact.

In general,



So, we can use chain Rule here

$$\delta_j^l = \sum_k \frac{\partial C}{\partial z_K^{l+1}} \times$$

$$\frac{\partial z_K^{l+1}}{\partial z_j^l}$$

$$z_j^l = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \times \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$= \sum_k \delta_k^{l+1} \times \left(w_{kj}^{l+1} \sigma'(z_j^l) \right)$$

$$\overbrace{\frac{\partial z_k^{l+1}}{\partial z_j^l}} \Rightarrow$$

$$z_k^{l+1} = \sum_i w_{ki}^{l+1} a_i^l + b_k^{l+1}$$

$$a_i^l = \sigma(z_i^l)$$

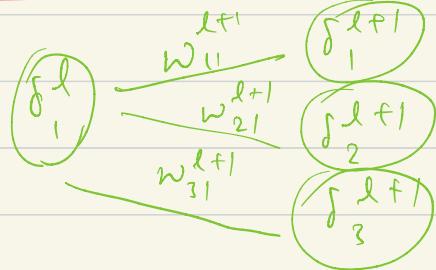
$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = \frac{\partial}{\partial z_j^l} \left[w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1} \right]$$

$$= w_{kj}^{l+1} \sigma'(z_j^l)$$

$$\delta_j^l = \sum_k \delta_k^{l+1} \cdot w_{kj}^{l+1} \sigma'(z_j^l)$$

$$\delta_j^l = \left(\sum_k w_{kj}^{l+1} \delta_k^{l+1} \right) \cdot \sigma'(z_j^l)$$

vectorization :



$$\delta^l_{(n_l \times 1)} = \begin{pmatrix} w^{l+1}^T & \delta^{l+1}_{(n_{l+1} \times 1)} \end{pmatrix} \odot \sigma'(z^l)_{(n_l \times 1)}$$

Element wise
multiplication

$$\delta^L_{(n_L \times 1)} = (A^L - y^L) \odot \sigma'(z^L)_{(n_L \times 1)}$$

Finally ,

$$\frac{\partial C}{\partial w_{ij}^l} = \left(\frac{\partial C}{\partial z_i^l} \right) * \frac{\partial z_i^l}{\partial w_{ij}^l}$$

↓

$$= (\delta_i^l) * a_j^{l-1}$$

↗

$$z_i^l = \sum_j a_j^{l-1} w_{ij}^l + b_i^l$$

For a particular (j) $\in (l-1)^{th}$ layer.

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = \frac{d}{d w_{ij}^l} \left[a_j^{l-1} w_{ij}^l + b_i^l \right]$$

~~a_j^{l-1}~~ ~~w_{ij}^l~~ ~~b_i^l~~

$$= a_j^{l-1}$$

Rectification:

$$\boxed{\frac{dw^l}{(n_l \times n_{l-1})} = \delta^l_{(n_l \times 1)} \cdot a^{l-1 \top}_{(1 \times n_{l-1})}}$$

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial b_j^l}$$

For proof, look
 prev slide

$$= \delta_j^l \cdot (1)$$

Vectorize:

$$\boxed{db^l = \delta^l \quad (n_e \times 1)}$$

All Vectorized eq at once:

$$\delta^L = (A^L - y^L) \odot \sigma'(z^L) \quad (n_e \times 1)$$

$$\delta^l = \left(\omega^{l+1}^T \cdot \delta^{l+1} \right) \odot \sigma'(z^l) \quad (n_e \times 1)$$

$$\frac{\partial C}{\partial b^l} = \delta^l$$

$$\frac{\partial C}{\partial \omega^l} = \delta^l \cdot a^{l-1}^T \quad (1 \times n_{l-1})$$

* Previous backprop eq are same for any activation function but not for any cost function. So, Let's create backprop eq for any cost function.

Let, $C \Rightarrow C(y^L, a^L)$

$$\text{Ex1: } C(y^L, a^L) = \sum \frac{1}{2} (y^L - a^L)^2$$

$$\text{Ex2: } C(y^L, a^L) = -\frac{1}{m} \sum (y \log a + (1-y) \log(1-a))$$

In this case everything remains same except $\frac{\partial C}{\partial z^L}$.

$$\delta^L = \frac{\partial C}{\partial z^L} = \frac{\partial C}{\partial a^L} \times \frac{\partial a^L}{\partial z^L}$$

$$\boxed{\delta^L = \frac{\partial C}{\partial a^L} \times \sigma'(z^L)}$$

Vectorized form:

$$\delta^L = \frac{da^L}{(n_L \times 1)} \odot \sigma'(z^L) \quad (n_L \times 1)$$

Let's evaluate (δ^L) for log-likelihood cost func:

$$C(y^L, a^L) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)}) \right)$$

$$\frac{\partial C}{\partial a^L} \Big|_{\text{for 1 example}} = - \left(\frac{y}{a} + \frac{(1-y)}{(1-a)} \right)$$

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

$$\boxed{\frac{\partial C}{\partial a^L} = \frac{a-y}{a(1-a)}}$$

If $\sigma(z)$ is sigmoid function, then

$$\sigma'(z) = \sigma(z)(1-\sigma(z)) = a(1-a)$$

Finally,

$$\delta^L = \frac{\partial C}{\partial a^L} \times \sigma'(z^L)$$

$$= \frac{a-y}{a(1-a)} \times a(1-a) = \boxed{a-y}$$