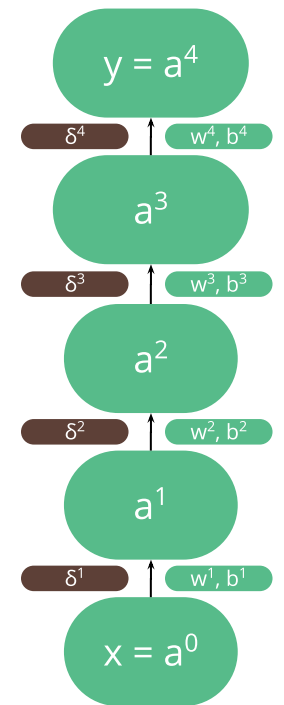


Notation

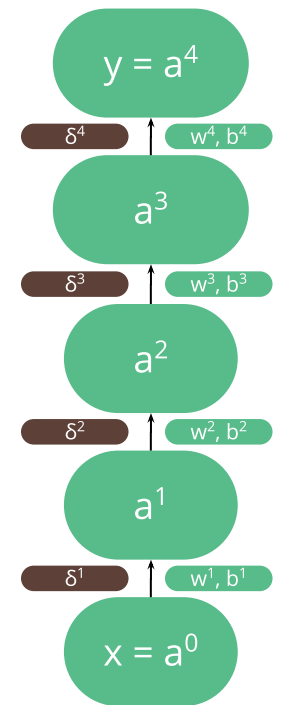
- t -> Ground Truth Output
- C -> Cost Function
- δ -> Gradient



The Cost Function

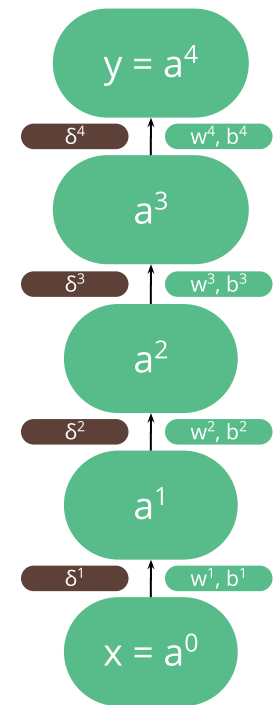
For a scalar output

- Mean Squared Error: $C = \frac{1}{2} * (y - t)^2$
- Cross Entropy: $C = t * \ln(y) + (1-t) * \ln(1-y)$



Backpropagation

- Goal: Compute $\partial C / \partial w$ and $\partial C / \partial b$
- Why: Use them for Stochastic Gradient Descent
- Define: $\delta^l = \partial C / \partial z^l$



Backward Pass

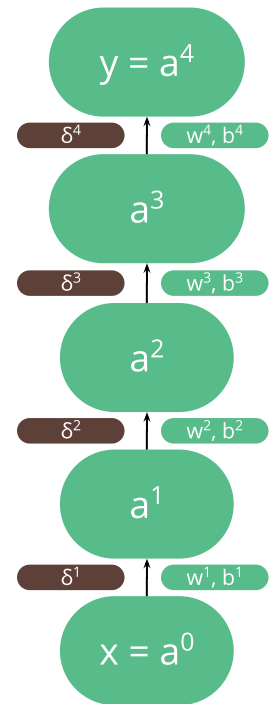
$$\delta^4 = \partial C / \partial z^4 = \partial C / \partial y * \partial y / \partial z^4$$

Now

- $\partial C / \partial y = (y - t)$
- $\partial y / \partial z^4 = \partial a^4 / \partial z^4 = f'(z^4)$

where $f'(\cdot)$ is derivative of $f(\cdot)$ w.r.t (\cdot)

$$\Rightarrow \delta^4 = (y - t) * f'(z^4)$$



Backward Pass

$$\delta^3 = \partial C / \partial z^3$$

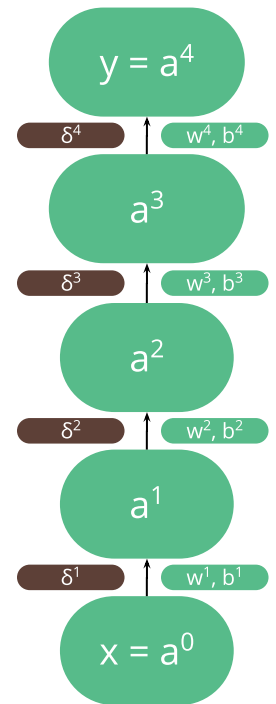
Now

- $z^4_1 = \dots + w^4_{1j} * f(z^3_j) + \dots$
- $z^4_k = \dots + w^4_{kj} * f(z^3_j) + \dots$

i.e. all elements of z^4 depend on z^3_j

Thus, by chain rule we can say that

$$\delta^3_j = \partial C / \partial z^3_j = \sum_k \partial C / \partial z^4_k * \partial z^4_k / \partial z^3_j$$



Backward Pass

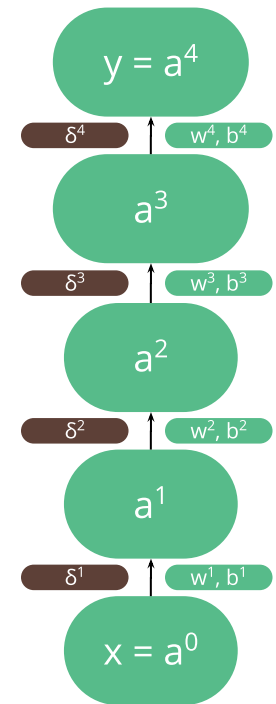
$$\delta_j^3 = \partial C / \partial z_j^3 = \sum_k \partial C / \partial z_k^4 * \partial z_k^4 / \partial z_j^3$$

$$\Rightarrow \delta_j^3 = \sum_k \partial C / \partial z_k^4 * \partial z_k^4 / \partial a_j^3 * \partial a_j^3 / \partial z_j^3$$

Now

- $\partial C / \partial z_k^4 = \delta_k^4$
- $\partial z_k^4 / \partial a_j^3 = w_{kj}^4$ [As $z_k^4 = \dots + w_{kj}^4 * a_j^3 + \dots$]
- $\partial a_j^3 / \partial z_j^3 = f'(z_j^3)$

$$\Rightarrow \delta_j^3 = (\sum_k \delta_k^4 * w_{kj}^4) * f'(z_j^3)$$

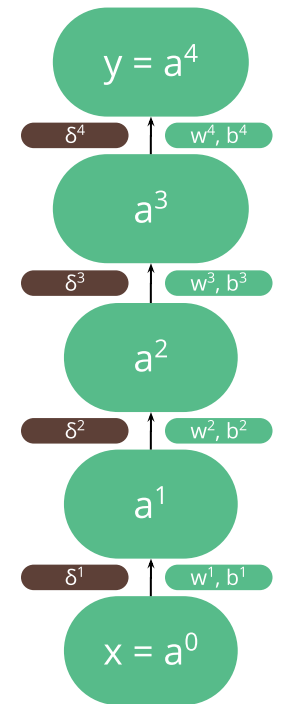
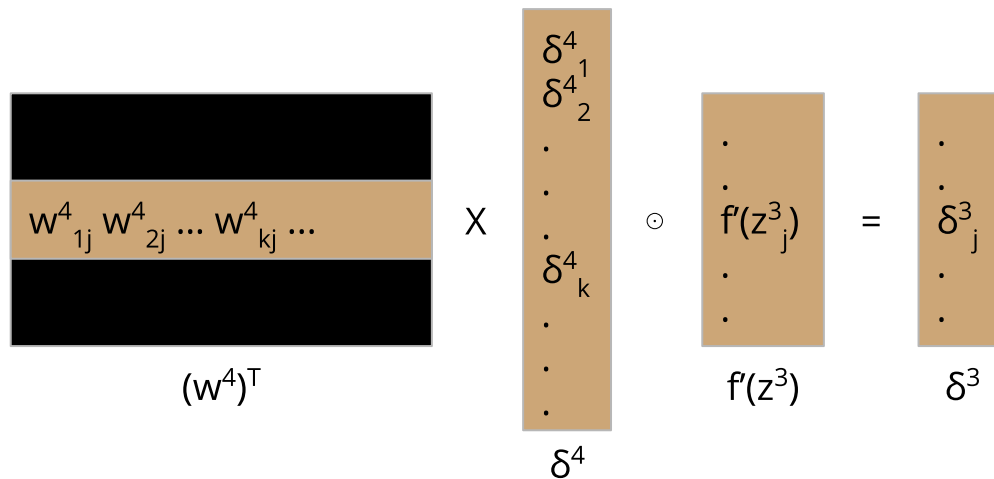


Backward Pass

$$\delta_j^3 = \left(\sum_k \delta_k^4 * w_{kj}^4 \right) * f'(z_j^3)$$

$$\Rightarrow \delta^3 = (w^4)^T * \delta^4 \odot f'(z^3)$$

where \odot = Element-wise product



Backward Pass

Hence we have

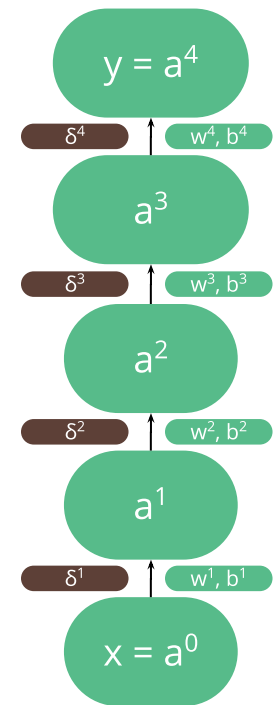
- $\delta^4 = (y - t) * f'(z^4)$
- $\delta^3 = (w^4)^T * \delta^4 \odot f'(z^3)$
- $\delta^2 = (w^3)^T * \delta^3 \odot f'(z^2)$
- $\delta^1 = (w^2)^T * \delta^2 \odot f'(z^1)$

Or in general

$$\delta^l = (w^{l+1})^T * \delta^{l+1} \odot f'(z^l) \quad \text{for } l = 1, 2, \dots, L-1$$

$$\delta^L = \nabla_y C \odot f'(z^L)$$

where $\nabla_y C$ is derivative of cost wrt output



Backward Pass

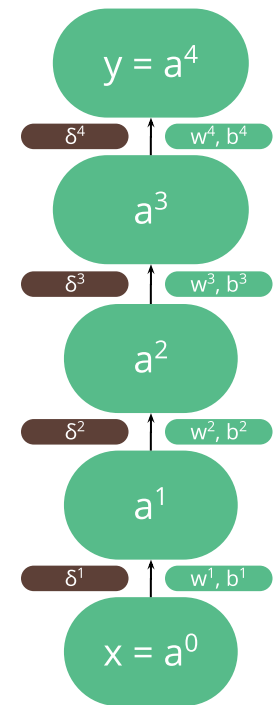
Now for our main objectives: $\partial C / \partial w_{jk}^l$ and $\partial C / \partial b_j^l$

$$\partial C / \partial w_{jk}^l = \partial C / \partial z_j^l * \partial z_j^l / \partial w_{jk}^l$$

Since

- $\partial C / \partial z_j^l = \delta_j^l$
 - $\partial z_j^l / \partial w_{jk}^l = a_k^{l-1}$
- [As $z_j^l = \dots + w_{jk}^l * a_k^{l-1} + \dots$]

$$\Rightarrow \partial C / \partial w_{jk}^l = \delta_j^l a_k^{l-1}$$



Backward Pass

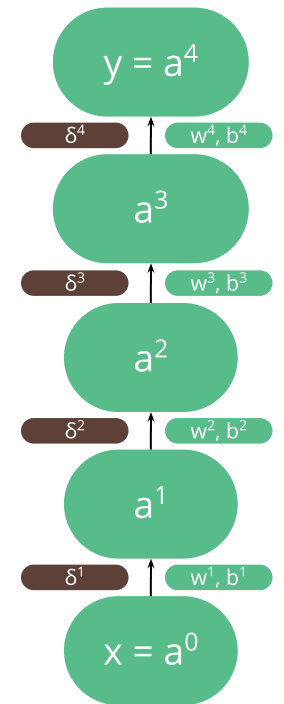
$$\partial C / \partial w_{jk}^l = \delta_j^l a_{jk}^{l-1}$$

Or in general

$$\partial C / \partial w^l = \delta^l * (a^{l-1})^T \text{ for } l = 1, \dots, L$$

$$\begin{array}{|c|} \hline \delta_1^l a^{l-1}_1 \dots \delta_1^l a^{l-1}_m \\ \hline \cdot \qquad \qquad \cdot \\ \hline \cdot \qquad \qquad \cdot \\ \hline \delta_n^l a^{l-1}_1 \dots \delta_n^l a^{l-1}_m \\ \hline \end{array} = \begin{array}{|c|} \hline \delta_1^l \\ \hline \cdot \\ \hline \cdot \\ \hline \delta_n^l \\ \hline \end{array} \times \begin{array}{|c|} \hline a^{l-1}_1 \dots a^{l-1}_m \\ \hline \end{array}$$

$\partial C / \partial w^l$
 δ^l
 $(a^{l-1})^T$



Backward Pass

Also

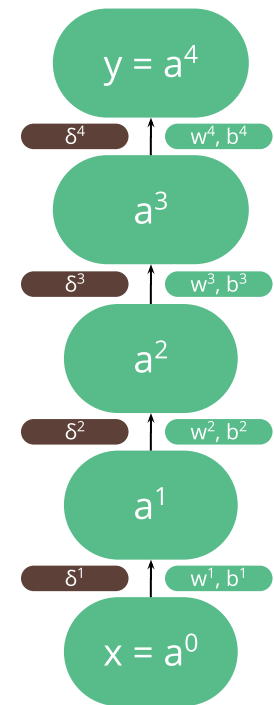
$$\partial C / \partial b_j^l = \partial C / \partial z_j^l * \partial z_j^l / \partial b_j^l$$

Since

- $\partial C / \partial z_j^l = \delta_j^l$
- $\partial z_j^l / \partial b_j^l = 1$ [As $z_j^l = \dots + b_j^l$]

$$\Rightarrow \partial C / \partial b_j^l = \delta_j^l$$

Or in general $\partial C / \partial b^l = \delta^l$ for $l = 1, \dots, L$



Backward Pass

In general:

$\delta^L = \nabla_y C \odot f'(z^L)$ where $\nabla_y C$ is derivative of cost wrt output

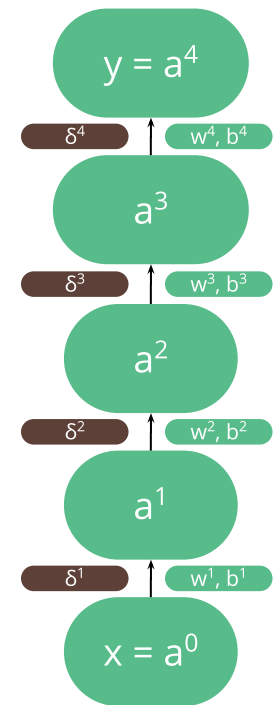
Then for $l = 1, 2, \dots, L-1$

$\delta^l = (w^{l+1})^T * \delta^{l+1} \odot f'(z^l)$ where \odot stands for element wise product

Finally for $l = 1, \dots, L$

$$\partial C / \partial w^l = \delta^l * (a^{l-1})^T$$

$$\partial C / \partial b^l = \delta^l$$



Summary

Forward Pass

The Input

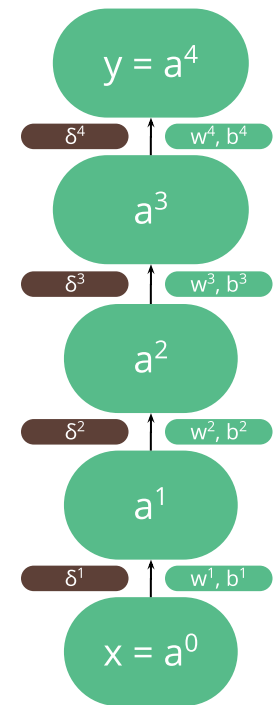
- $a^0 = x$

For $l = 1, \dots, L$ layers

- $z^l = w^l * a^{l-1} + b^l$
- $a^l = f(z^l)$

Finally

- $y = a^L$



Backward Pass

$$\delta^L = \nabla_y C \odot f'(z^L)$$

where $\nabla_y C$ is derivative of cost wrt output

Then for $l = 1, 2, \dots, L-1$

$$\delta^l = (w^{l+1})^T * \delta^{l+1} \odot f'(z^l)$$

where \odot stands for element wise product

Finally for $l = 1, \dots, L$

$$\partial C / \partial w^l = \delta^l * (a^{l-1})^T$$

$$\partial C / \partial b^l = \delta^l$$

