$$E_n \triangleq \frac{1}{2} \sum_{k=1}^{K} (y_k - t_k)^2$$

$$\delta_k \triangleq \frac{\partial E_n}{\partial a_k} = \frac{1}{2} \cdot 2 (y_k - t_k) \boxed{\frac{\partial y_k}{\partial a_k}} \rightarrow = 1 \text{ as}$$

output layer
activation

$$y_k = a_k$$
$$(\sigma(\cdot) = \text{unit fn.})$$

$$\delta_k = y_k - t_k \quad (\text{Else, according to the specific activation function } \sigma(\cdot) \text{ at the output layer)}$$

③ $\underline{\text{Backpropagate these to obtain } \delta_j\text{'s for}}$
$\underline{\text{the hidden layer units}}$

$$\boxed{\delta_j = (1 - z_j^2) \sum_{k=1}^{K} \omega_{kj} \boxed{\delta_k}}$$

What is this, and how?

→ previous step (step #2) [output layer]

[hidden layer]



$$\delta_j \triangleq \frac{\partial E_n}{\partial a_j} = \sum_{\forall k} \boxed{\frac{\partial E_n}{\partial a_k}} \boxed{\frac{\partial a_k}{\partial a_j}}$$

$\delta_k$

[step #2]

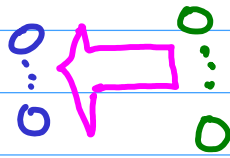$$a_k = \underline{w}^{(2)T} \underline{z}$$
$$= \sum_{j=0}^{M} \omega_{kj}^{(2)} z_j$$
$$= \sum_{j=0}^{M} \omega_{kj}^{(2)} h(a_j)$$

$$\Rightarrow \frac{\partial a_k}{\partial a_j} = \omega_{kj}^{(2)} \frac{\partial h(a_j)}{\partial a_j} = \omega_{kj}^{(2)} h'(a_j)$$
$$= \omega_{kj}^{(2)} (1 - z_j^2)$$

$$\delta_j = \sum_{\forall k} \delta_k \omega_{kj}^{(2)} (1 - z_j^2) = (1 - z_j^2) \sum_{\forall k} \omega_{kj}^{(2)} \delta_k$$

④ Use the chain rule to evaluate the gradient

$$\frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \delta_j x_i, \quad \frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \delta_k z_j$$

What is this, and how?



$$\bar{\nabla} E = \begin{bmatrix} \partial E_n / \partial \omega \end{bmatrix}$$

all weights (1st & 2nd layers of connections)

$$\frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \boxed{\frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial \omega_{ji}^{(1)}}} \quad ; \quad a_j = w_{\cdot j}^{(1)T} x$$

$$= \sum_{i=0}^{D} \omega_{ji}^{(1)} x_i$$

$\delta_j$ (previous step: step # 3)

$$\Rightarrow \frac{\partial a_j}{\partial \omega_{ji}^{(1)}} = x_i$$

$$\Rightarrow \frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \delta_j x_i$$

$$\frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \boxed{\frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial \omega_{kj}}} \quad ; \quad a_k = w_k^{(2)T} z = \sum_{j=0}^{M} \omega_{kj}^{(2)} z_j$$

$\delta_k$ (step # 2)

$$\Rightarrow \frac{\partial a_k}{\partial \omega_{kj}^{(2)}} = z_j$$

$$\Rightarrow \frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \delta_k z_j$$

# Side topic! Numerical Evaluation of the gradient

Empirically, all of these alternative methods are numerically not as good as Backpropagation

$$(*) \quad \frac{\partial E_n}{\partial \omega_{ji}} = \frac{E_n(\omega_{ji} + \epsilon) - E_n(\omega_{ji})}{\epsilon} + o(\epsilon)$$

$$(\omega_{ji} + \epsilon) - (\omega_{ji})$$

OR: symmetrical central differences

$$\frac{\partial E_n}{\partial \omega_{ji}} = \frac{E_n(\omega_{ji} + \epsilon) - E_n(\omega_{ji} - \epsilon)}{2\epsilon} + o(\epsilon)$$

why? The first direct formula?

$$E_n(\omega_{ji} + \epsilon) = E_n(\omega_{ji}) + \epsilon \frac{\partial E_n}{\partial \omega_{ji}} + \frac{1}{2} \epsilon H \epsilon + \text{higher order terms}$$

$$\Rightarrow \quad \frac{E_n(\omega_{ji} + \epsilon) - E_n(\omega_{ji})}{\epsilon} = \frac{\partial E_n}{\partial \omega_{ji}} + \frac{1}{2} \epsilon H$$

why? The symmetrical central differences formula:

$$E_n(\omega_{ji} + \epsilon) = E_n(\omega_{ji}) + \epsilon \frac{\partial E_n}{\partial \omega_{ji}} + \frac{1}{2} \epsilon H \epsilon + o(\epsilon^3)$$

$$E_n(\omega_{ji} - \epsilon) = E_n(\omega_{ji}) - \epsilon \frac{\partial E_n}{\partial \omega_{ji}} + \frac{1}{2} \epsilon H \epsilon - o(\epsilon^3)$$

$$E_n(\omega_{ji} + \epsilon) - E_n(\omega_{ji} - \epsilon) = 2\epsilon \frac{\partial E_n}{\partial \omega_{ji}} + 2 o(\epsilon^3)$$

$$\frac{\partial E_n}{\partial \omega_{ji}} = \frac{E_n(\omega_{ji} + \epsilon) - E_n(\omega_{ji} - \epsilon)}{2\epsilon} + O(\epsilon^2)$$

<u>Physical Significance</u> : NN-based solution vis-a-vis the linear (or restricted linear) method done before.

(*) NN: the weight parameters in the first layer are shared between the outputs
linear: each classification is performed independently.

(*) The first layer of the network can be viewed as performing a non-linear feature extraction, and sharing the features between different outputs can lead to savings on computation and also lead to improved generalisation.