

# DEFINITIONS - Colored Orange

## Network

We will fix a single input sample  $\vec{x}$  which expects output  $\vec{t}$

Layer index  $k = 0, \dots, L$  (i.e.: there are  $L+1$  layers)

Input Neuron index in layer  $k$   $i = 0, \dots, \#k-1$

Neuron index in layer  $k$   $j = 0, \dots, \#k$

Output neuron index  $q = 0, \dots, \#L$

$\#k$  Size of layer  $k$  ( $\#0 = 784$   $\#L = 10$ )

$\#k = (\#k) - 1$  last index of layer  $k$  (without bias neuron)

$$a_j^k = g(h_j^k) \text{ (Neuron } j \text{ of layer } k)$$

$$a_j^0 = x_j \text{ (input } j) \quad a_{\#k}^k = 1 \text{ (bias extension)}$$

$$h_j^k = \sum_{i=0}^{\#k-1} a_i^{k-1} w_{ij}^{k-1} \quad g(x): \text{ activation function}$$

$w_{ij}^k$ : weight from  $a_i^{k-1}$  to  $a_j^k$

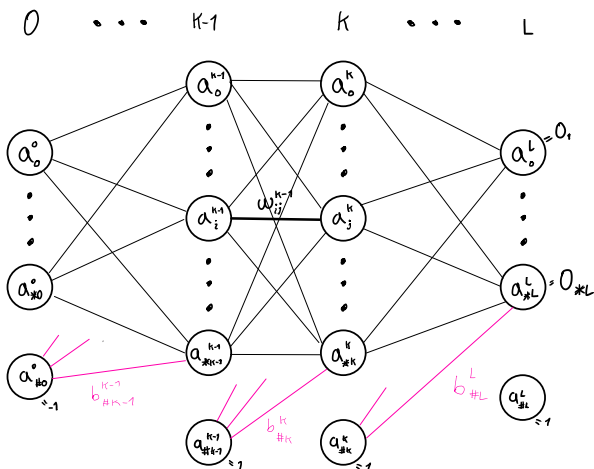
$b_j^k = w_{\#k-1, j}^{k-1}$  bias of neuron  $a_j^k$

$$O_q = a_q^L \text{ (Output neuron } q) \text{ (we want } \vec{O} = \vec{t})$$

$$E(\vec{x}, \vec{t}) = \sum_{q=0}^{\#L} (O_q - t_q)^2 / \#L$$

$$\nabla E_{ij}^k = \frac{\partial E(\vec{x}, \vec{t})}{\partial w_{ij}^k} \text{ (influence of } w_{ij}^k \text{ in } E)$$

## NETWORK DIAGRAM



# EQUATIONS DEVELOPMENT

## Gradient

We would like to decrease the error  $E$  by tweaking  $w_{ij}^k$

$$\nabla E_{ij}^k = \frac{\partial E(\bar{y}, I)}{\partial w_{ij}^k} = \frac{\partial}{\partial w_{ij}^k} \frac{1}{\#L} \sum_{q=0}^{\#L} (O_q - t_q)^2$$

$$= \frac{1}{\#L} \sum_{q=0}^{\#L} \frac{\partial (O_q - t_q)^2}{\partial w_{ij}^k}$$

$$\nabla E_{ij}^k = \frac{2}{\#L} \cdot \sum_{q=0}^{\#L} (O_q - t_q) \frac{\partial a_q^L}{\partial w_{ij}^k} \quad (1)$$

Given that we usually choose  $g$  with known derivatives

our remaining unknown is  $\frac{\partial a_q^l}{\partial w_{ij}^k}$  with  $l=L$ .

As it can be anticipated we will need to analyse  $\frac{\partial a_q^l}{\partial w_{ij}^k}$  for  $l=0, \dots, L$

Let's proceed our analysis by cases:

If  $l=0 \Rightarrow \frac{\partial a_q^0}{\partial w_{ij}^k} = 0$  ( $a_q^0 = 0$ )

Else if  $q = \#l$  (bias neuron)  $\Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = 0$  ( $a_{\#l}^l = 1$ )

Else if  $k \geq l \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = 0$  ( $w_{ij}^k$  is posterior to  $a_q^l$  in the network)

Else if  $k = l-1 \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = \frac{\partial g(h_q^l)}{\partial h_q^l} \cdot \frac{\partial h_q^l}{\partial w_{ij}^k}$

And we have

$$\frac{\partial h_q^l}{\partial w_{ij}^k} = \frac{\partial \sum_{r=0}^{\#l-1} a_r^{l-1} w_{r,q}^{l-1}}{\partial w_{ij}^{l-1}} = \frac{\partial \sum_{r=0, r \neq i}^{\#l-1} a_r^{l-1} w_{r,q}^{l-1}}{\partial w_{ij}^{l-1}} + \frac{\partial a_i^{l-1} w_{i,q}^{l-1}}{\partial w_{ij}^{l-1}} \quad \text{therefore:}$$

If  $i = q \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = g'(h_q^l) \cdot a_i^{l-1}$

If  $i \neq q \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = g'(h_q^l) \cdot 0 = 0$

Else if  $k < l-1 \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = \frac{\partial \sum_{r=0}^{\#l-1} a_r^{l-1} w_{r,q}^{l-1}}{\partial w_{ij}^{l-1}} = \sum_{r=0}^{\#l-1} w_{r,q}^{l-1} \frac{\partial a_r^{l-1}}{\partial w_{ij}^k}$  \* (bias term is 0)

$$\Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = g'(h_q^l) \cdot \sum_{r=0}^{\#l-1} w_{r,q}^{l-1} \frac{\partial a_r^{l-1}}{\partial w_{ij}^k}$$

In conclusion we can now recursively express  $\frac{\partial a_q^l}{\partial w_{ij}^k}$ :

$$\frac{\partial a_q^l}{\partial w_{ij}^k} = g'(h_q^l) \cdot \begin{cases} 0 & \text{if } l=0 \text{ or } k \geq l \text{ or } q = \#l \\ & \text{or } (k=l-1 \text{ and } j \neq q) \\ a_i^k & \text{if } k=l-1 \text{ and } j=q \\ \sum_{r=0}^{\#l-1} w_{r,q}^{l-1} \cdot \frac{\partial a_r^{l-1}}{\partial w_{ij}^k} & \text{else} \end{cases} \quad (2)$$

## Implementation Details

Is important to note that (2) holds for valid values of  $l$  and  $k$ ,

i.e:  $l, k \in \{0, \dots, L\}$ . The implementation will need to perform checks to ensure  $l$  and  $k$  are not out of range.

Also, while the cases for  $q = \#l$  and  $k \geq l$  are theoretically correct, when doing backpropagation those cases should not be reached, therefore we will assert so as well.

It will be also important to add an extra "Else if  $k=l-2$ " case after the "Else if  $k=l-1$ " guard. While this case is not necessary for completeness sake, it will provide a significant performance gain for our initial recursive and slow implementation.

Else if  $k=l-2$ :

$$\frac{\partial h_q^l}{\partial w_{ij}^k} = \frac{\partial \sum_{r=0}^{\#l-1} a_r^{l-1} w_{r,q}^{l-1}}{\partial w_{ij}^{l-2}} = \sum_{r=0}^{\#l-1} w_{r,q}^{l-1} \cdot \frac{\partial a_r^{l-1}}{\partial w_{ij}^{l-2}}$$

and because of our previous case considering  $r$  and  $j$  this is

$$= w_{j,q}^{l-1} \cdot g'(h_j^{l-1}) \cdot a_i^q \quad \text{therefore}$$

$$\boxed{k=l-2 \Rightarrow \frac{\partial a_q^l}{\partial w_{ij}^k} = g'(h_j^l) \cdot w_{j,q}^{l-1} \cdot g'(h_j^{l-1}) \cdot a_i^q} \quad (2')$$

# MATRICIZATION

Now let's translate our equations to matrices so that we can use numpy matrix operations and get a significant speed boost.

Let's define:

$$\alpha^l := [a_1 \dots a_{\#l}] = [a_i^l]_{ij} \quad \text{no bias } i=0, \dots, \#k$$

$$k = 0, \dots, L-1$$

$$\omega^k := \begin{bmatrix} \omega_{1,1}^k & \dots & \omega_{1,\#k}^k \\ \vdots & \ddots & \vdots \\ \omega_{\#k,1}^k & \dots & \omega_{\#k,\#k}^k \end{bmatrix} = [\omega_{ij}^k]_{ij} \quad \begin{array}{l} i = 0, \dots, (\#k)+1 \rightarrow \text{bias} \\ j = 0, \dots, \#(k+1) \end{array}$$

$$\nabla E^k := \begin{bmatrix} \frac{\partial E}{\partial \omega_{1,1}^k} & \dots & \frac{\partial E}{\partial \omega_{1,\#k}^k} \\ \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial \omega_{\#k,1}^k} & \dots & \frac{\partial E}{\partial \omega_{\#k,\#k}^k} \end{bmatrix} \stackrel{\textcircled{1}}{=} \left[ \frac{2}{\#L} \cdot \sum_{q=0}^{\#L} (o_q - t_q) \frac{\partial \alpha_q^L}{\partial \omega_{ij}^k} \right]_{ij}$$

$$A_k^{l,q} := \left[ \frac{\partial \alpha_q^L}{\partial \omega_{ij}^k} \right]_{ij} \Rightarrow \nabla E^k = \frac{2}{\#L} \cdot \sum_{q=0}^{\#L} (o_q - t_q) A_k^{L,q} \quad \textcircled{3}$$

Now let's use  $\textcircled{2}$  to obtain an expression by cases for  $A_k^{l,q}$

$$\text{If } l = k+1 \Rightarrow A_k^{l,q} = \left[ \frac{\partial \alpha_q^{k+1}}{\partial \omega_{ij}^k} \right]_{ij} \stackrel{\textcircled{2}}{=} g'(h_q^k) \begin{bmatrix} 0 & \text{if } j \neq q \\ \alpha_i^k & \text{if } j = q \end{bmatrix}_{ij} \quad \textcircled{4}$$

$$\text{Else if } l = k+2 \Rightarrow A_k^{l,q} = \left[ \frac{\partial \alpha_q^{k+2}}{\partial \omega_{ij}^k} \right]_{ij} \stackrel{\textcircled{2}}{=} \left[ g'(h_q^k) \omega_{j,q}^{k+1} g'(h_j^{k+1}) \alpha_i^k \right]_{ij} \quad \textcircled{5}$$

$$\text{Else if } l > k+1 \Rightarrow A_k^{l,q} = \left[ \frac{\partial \alpha_q^l}{\partial \omega_{ij}^k} \right]_{ij} \stackrel{\textcircled{2}}{=} \left[ g'(h_q^k) \sum_{r=0}^{l-1} \omega_{r,q}^{l-1} \cdot \frac{\partial \alpha_r^{l-1}}{\partial \omega_{ij}^k} \right]_{ij} \quad \textcircled{6}$$

$$= g'(h_q^k) \sum_{r=0}^{l-1} \omega_{r,q}^{l-1} A_k^{l-1,r}$$

In summary:

$$\textcircled{3} \quad \nabla E^k = \frac{2}{\#L} \cdot \sum_{q=0}^{\#L} (o_q - t_q) A_k^{L,q}$$

$$\textcircled{4} \quad A_k^{k+1,q} = g'(h_q^k) \begin{bmatrix} 0 & \text{if } j \neq q \\ \alpha_i^k & \text{if } j = q \end{bmatrix}_{ij}$$

$$\textcircled{5} \quad A_k^{k+2,q} = g'(h_q^k) \left[ \omega_{j,q}^{k+1} g'(h_j^{k+1}) \alpha_i^k \right]_{ij}$$

$$\textcircled{6} \quad A_k^{l,q} = g'(h_q^k) \sum_{r=0}^{l-1} \omega_{r,q}^{l-1} A_k^{l-1,r} \quad l > k+1$$