# Backpropagation: (FOR MLP)

Consider a network with $L$ layers & $N_\ell$ neurons in the $\ell$-th layer

$W^\ell = \left(W_{ij}^\ell\right)$ $\rightarrow$ weight for the connection between the $i$-th neuron in layer $\ell$ & $j$-th neuron in layer $\ell-1$

$b^\ell \rightarrow$ column vector with $b_i^\ell$ being the bias for the $i$-th neuron in layer $\ell$

$\sigma^\ell \rightarrow$ activation function of $\ell$-th layer, acting componentwise on any vector

$a^\ell \rightarrow$ output vector from layer $\ell$ after activation.

$$a^\ell = \sigma^\ell\left(W^\ell a^{\ell-1} + b^\ell\right)$$

$$W^\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}, \quad b^\ell \in \mathbb{R}^{N_\ell \times 1}$$

$$a^\ell \in \mathbb{R}^{N_\ell \times 1}$$

Introduce the vector:
$$Z^\ell = W^\ell a^{\ell-1} + b^\ell$$

## Note:
1) The first layer is the input layer. No computations occur here. It only provides an input signal.
2) The final output of the network is $a^L$
3) Traditionally, $\sigma^L$ is not called an activation function but an output function. (this is only for the final output layer)

Consider an input $X \in \mathbb{R}^{N_I} \equiv \mathbb{R}^{N_1}$, which has the true output $Y \in \mathbb{R}^{N_O} \equiv \mathbb{R}^{N_L}$.

Clearly $a^1 = X$

Define an error/cost function to be minimized

$$E := E(\hat{Y}, Y)$$

where $\hat{Y} = a^L$ is the network prediction. For instance, this could be the mean squared error.

Having evaluated $\hat{Y}$ (this corresponds to a "forward pass"), we now wish to evaluate

$$\frac{\partial E}{\partial W_{ij}^\ell}, \quad \frac{\partial E}{\partial b_i^\ell}$$

needed to update the weights & biases of the network.

Now

$$Z_k^\ell = \sum_{j=1}^{N_{\ell-1}} W_{kj}^\ell \, a_j^{\ell-1} + b_k^\ell$$

$$\qquad\qquad\qquad\qquad\qquad k = 1, \cdots, N_\ell$$

$$a_k^\ell = \sigma^\ell(Z_k^\ell)$$

Thus, we have in the $(\ell+1)$-th layer

$$Z_m^{\ell+1} = \sum_{k=1}^{N_\ell} W_{mk}^{\ell+1} \, a_k^\ell + b_m^\ell \qquad\qquad m = 1, \ldots, N_{\ell+1}$$

Using chain rule, we have

$$\frac{\partial E}{\partial W_{kj}^\ell} = \frac{\partial E}{\partial Z_k^\ell} \frac{\partial Z_k^\ell}{\partial W_{kj}^\ell}$$

$$= \frac{\partial E}{\partial a_k^\ell} \frac{\partial a_k^\ell}{\partial Z_k^\ell} \, a_j^{\ell-1}$$

$$= \left[ \sum_{m=1}^{N_{\ell+1}} \frac{\partial E}{\partial z_m^{\ell+1}} \frac{\partial z_m^{\ell+1}}{\partial a_k^{\ell}} \right] \cdot \sigma^{\ell'}(z_k^{\ell}) \, a_j^{\ell-1}$$

Define the "error" corresponding to a neuron $k$ in the layer $\ell$ as

$$\delta_k^{\ell} = \frac{\partial E}{\partial z_k^{\ell}} = \left[ \sum_{m=1}^{N_{\ell+1}} \frac{\partial E}{\partial z_m^{\ell+1}} \frac{\partial z_m^{\ell+1}}{\partial a_k^{\ell}} \right] \sigma^{\ell'}(z_k^{\ell})$$

$$\Rightarrow \boxed{\delta_k^{\ell} = \left[ \sum_{m=1}^{N_{\ell+1}} \delta_m^{\ell+1} \, W_{mk}^{\ell+1} \right] \sigma^{\ell'}(z_k^{\ell})} \quad - \textcircled{1}$$

Thus,

$$\boxed{\frac{\partial E}{\partial W_{kj}^{\ell}} = \delta_k^{\ell} \cdot a_j^{\ell-1}} \quad - \textcircled{2}$$

Similarly,

$$\boxed{\frac{\partial E}{\partial b_k^{\ell}} = \frac{\partial E}{\partial z_k^{\ell}} \frac{\partial z_k^{\ell}}{\partial b_k^{\ell}} = \delta_k^{\ell}} \quad - \textcircled{3}$$

The errors are evaluated by propagating backwards. Thus, we need

$$\delta_j^{L} = \frac{\partial E}{\partial z_j^{L}} = \frac{\partial E}{\partial a_j^{L}} \frac{\partial a_j^{L}}{\partial z_j^{L}} = \frac{\partial E}{\partial a_j^{L}} \sigma^{L'}(z_j^{L})$$

Since $a_j^{L} = \hat{Y}_j$

$$\boxed{\delta_j^{L} = \frac{\partial E}{\partial \hat{Y}_j} \sigma^{L'}(z_j^{L})} \quad - \textcircled{4}$$

Summarizing the steps to evaluate the
derivatives of the cost/error junction wrt the
weight and biases:

1. Forward pass:
For the input $X$, do an evaluation of the network
while storing all $z^l$ & $a^l$ for each layer $l$.

2. Evaluate $\delta_j^L$ for the final layer $L$, which
will depend on the choice of error junction

3. Backward pass:
Traverse the network backwards to evaluate
$\delta_j^l$ from ① & ④

4. Find the derivatives using ② & ③