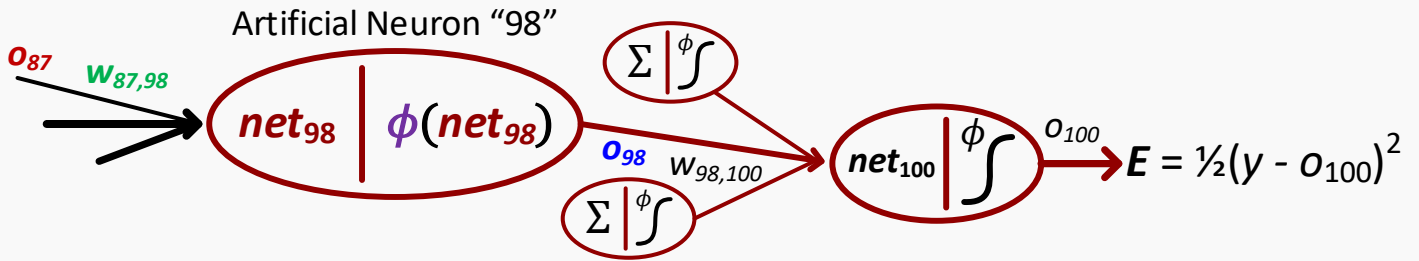


# backpropagation through neural network hidden layer

The following illustrates calculations passing back through an **Artificial Neural Network** layer:



Below represents the derivative of error  $E$  using the **chain rule** from calculus:

$$\frac{dE}{dw_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{dnet_{98}} \frac{dnet_{98}}{dw_{87,98}} \quad \phi(z) = \frac{1}{1 + e^z}$$

$$\phi'(z) = \frac{d\phi(z)}{dz} = \phi(z)(1 - \phi(z))$$

Computing the **last** derivative term:

$$\frac{dnet_{98}}{dw_{87,98}} = \frac{d(w_{87,98}o_{87} + w_{86,98}o_{86} + w_{85,98}o_{85} + \dots)}{dw_{87,98}} = o_{87}$$

Computing the **middle** derivative term:

$$\frac{do_{98}}{dnet_{98}} = \frac{d\phi(net_{98})}{dnet_{98}} = \phi'(net_{98})(1 - \phi(net_{98})) = o_{98}(1 - o_{98})$$

Computing the **first** derivative term:

$$\frac{dE}{do_{98}} = \frac{dE}{dnet_{100}} \frac{dnet_{100}}{do_{98}}$$

since  $\frac{dE}{dnet_{100}} = \delta_{100}$  and  $\frac{dnet_{100}}{do_{98}} = \frac{d(w_{99,100}o_{99} + w_{98,100}o_{98} + \dots)}{do_{98}} = w_{98,100}$

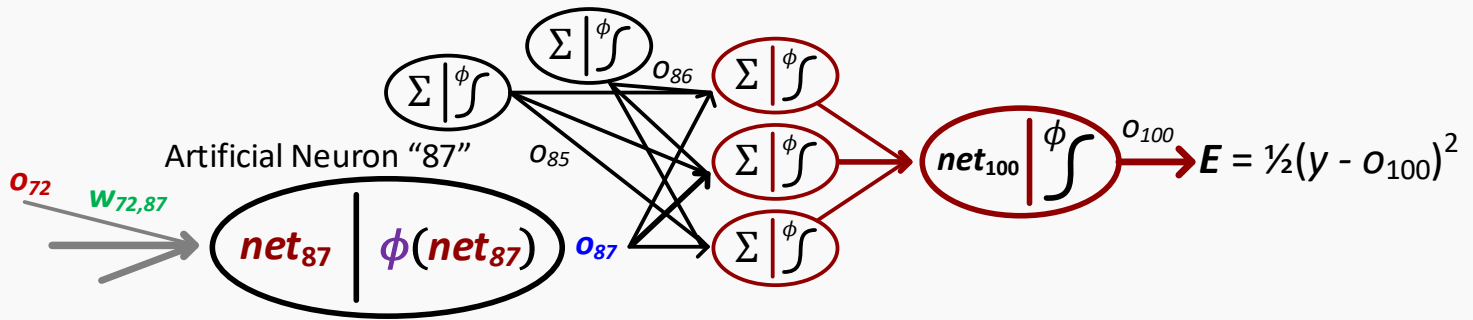
then  $\frac{dE}{do_{98}} = \delta_{100}w_{98,100}$

The **result** of the calculation above is as follows:

$$\frac{dE}{dw_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{dnet_{98}} \frac{dnet_{98}}{dw_{87,98}} = \delta_{100}w_{98,100}o_{98}(1 - o_{98})o_{87}$$

## Backpropagation through another Hidden Layer in a Neural Network

The following illustrates calculations passing back through another **Artificial Neural Network** layer:



Below represents the derivative of error  $E$  using the **chain rule** from calculus:

$$\frac{dE}{dw_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{dnet_{87}} \frac{dnet_{87}}{dw_{72,87}} \quad \phi(z) = \frac{1}{1 + e^z}$$

$$\phi'(z) = \frac{d\phi(z)}{dz} = \phi(z)(1 - \phi(z))$$

Computing the **last** derivative term:

$$\frac{dnet_{87}}{dw_{72,87}} = \frac{d(w_{72,87}o_{72} + w_{71,87}o_{71} + w_{70,87}o_{70} + \dots)}{dw_{72,87}} = o_{72}$$

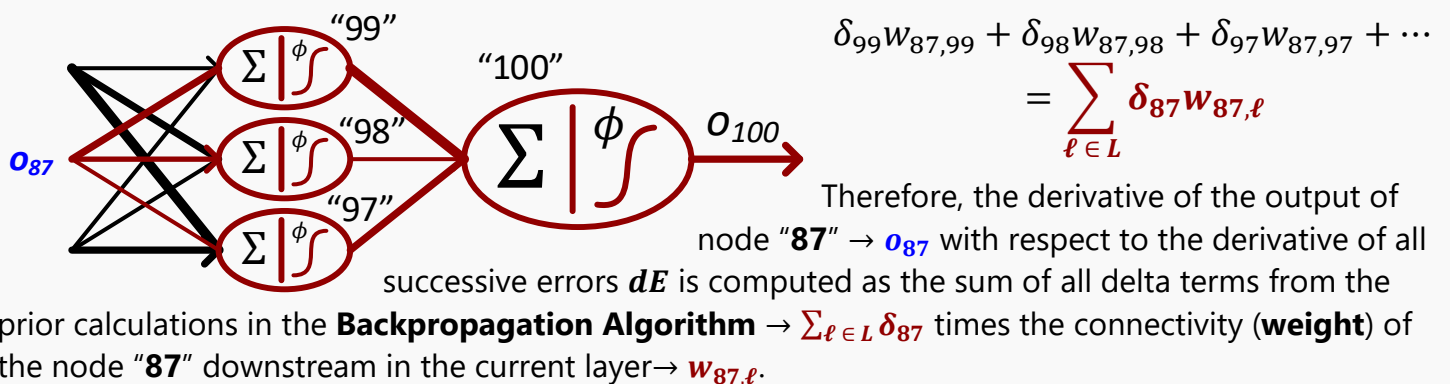
Computing the **middle** derivative term:

$$\frac{do_{87}}{dnet_{87}} = \frac{d\phi(net_{87})}{dnet_{87}} = \phi'(net_{87})(1 - \phi(net_{87})) = o_{87}(1 - o_{87})$$

Computing the **first** derivative term:

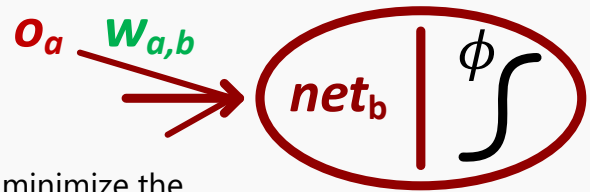
$$\frac{dE}{do_{87}} = \frac{dE}{dnet_{99}} \frac{dnet_{99}}{do_{87}} + \frac{dE}{dnet_{98}} \frac{dnet_{98}}{do_{87}} + \frac{dE}{dnet_{97}} \frac{dnet_{97}}{do_{87}} + \dots$$

As the layers become deeper in the **Artificial Neural Network**, the computation of the **first derivative term** in calculating the weight for a given node's output becomes increasingly complex. This is due to the output of nodes within the hidden layer being dependent of the derivative of the error  $dE$  of all successive nodes to the final output  $o_{100}$  as computed and illustrated below:



Writing the above computation in a more general way to understand **Backpropagation Intuition**:

$$\frac{dE}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{dnet_b} \frac{dnet_b}{dw_{a,b}}$$



Assuming Neuron "**a**" connected to Neuron "**b**". The

**optimization objective** is to adjust the **weight**  $w_{a,b}$  to minimize the derivative of the error with respect to the **weight**  $\frac{dE}{dw_{a,b}}$ ; using the **chain rule**:

$$\frac{dE}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{dnet_b} \frac{dnet_b}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{dnet_b} o_a$$

The **last** derivative term is always computed as the output of the prior node  $o_a$ .

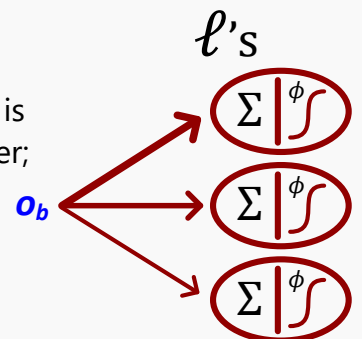
$$\frac{dE}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{dnet_b} \frac{dnet_b}{dw_{a,b}} = \frac{dE}{do_b} o_b(1 - o_b) o_a$$

The **middle** derivative term is always computed as the derivative of phi  $\phi \rightarrow o_b(1 - o_b)$ .

The **first** derivative term is the sum of the deltas  $\delta$  from all of the downstream nodes and the downstream weights ( $\sum_{\ell \in L} \delta_{\ell} w_{b,\ell}$ )

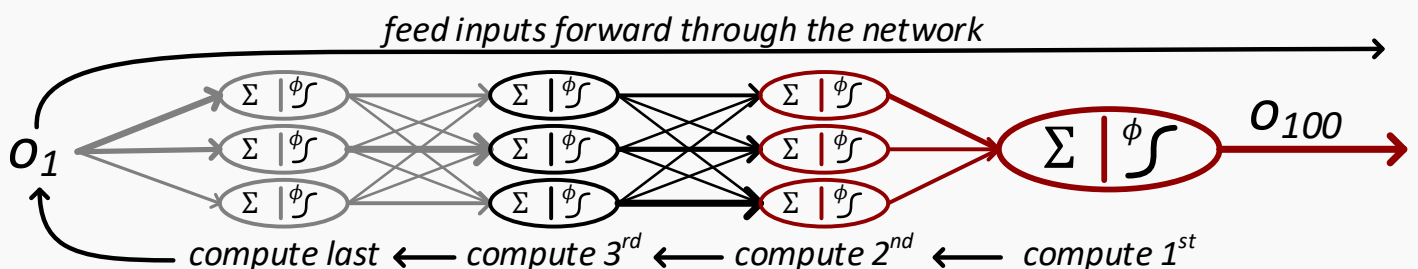
$$\frac{dE}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{dnet_b} \frac{dnet_b}{dw_{a,b}} = \left( \sum_{\ell \in L} \delta_{\ell} w_{b,\ell} \right) o_b(1 - o_b) o_a$$

It is important that the **L**'s represent the indices of **downstream neurons**. It is necessary to compute all of the  $\delta_{\ell}$ 's subsequent to the current node in a layer; working backward along the neural network in order to compute. The computation as a whole represents the **gradient**:



$$\left( \sum_{\ell \in L} \delta_{\ell} w_{b,\ell} \right) o_b(1 - o_b) o_a \rightarrow \text{Gradient Computation}$$

In other words, the implication of the above intuition is that the **gradients** are computed with respect to the **final output layer**, then the **gradients** are computed with respect to the **second last hidden layer**, then the **gradients** are computed with respect to the **third last hidden layer**, and so on... using the deltas  $\delta$ 's each time that were computed downstream of the current layer.



Given an understanding to compute the derivative  $dE$  with respect to the weight  $w_{a,b}$  for all weights  $w_{a,b}$ 's, the following explains the mathematics behind **gradient descent**:

$$w_{a,b} \leftarrow w_{a,b} - \alpha \frac{dE}{dw_{a,b}}$$

Each weight  $w_{a,b}$  is taken, with a **step-sized**  $\alpha$  **weight**  $w_{a,b}$  taken down the **gradient**.

In standard **Backpropagation**,  $\alpha$  is between **0** and **1** and referred to as the **learning rate**.

If the **learning rate**  $\alpha$  is too **low**, the neural network will learn very slowly. Conversely, if the **learning rate**  $\alpha$  is too **high**, the weights can cause an objective to **diverge** from the **optimal minimum**.

The general logic is that the weight changes depending on how **steep** the **gradient** is at that point.

Once the **errors** are **propagated** back through the neural network and the **weights** have been updated, the **error** must be updated; achieved by feeding the **input forward** through the network.

## Neural Network Summarized

**Backpropagation**: Repeat the process going **backwards** (**gradient** calculation), adjusting the **weights**, and feeding the **input forward** (**error** calculation) by many interactions for **learning**

### Advantages of Neural Networks

- Highly expressive nonlinear models
- Advanced in computer vision and speech that is unmatched through competing methods.
- Capable of capturing latent structure with the hidden network layers

### Disadvantages of Neural Networks

- Prone to becoming stagnant in local optima and produce poor solutions
- Another black-box algorithm that requires extensive parameter tuning (network structure).