

Derivation of Backpropagation

1 Introduction

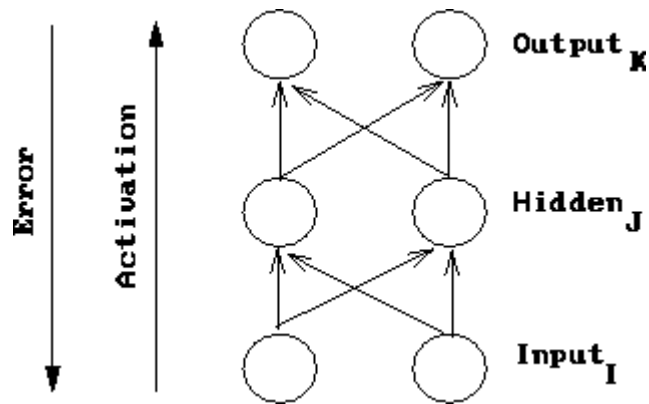


Figure 1: Neural network processing

Conceptually, a network forward propagates activation to produce an output and it backward propagates error to determine weight changes (as shown in Figure 1). The weights on the connections between neurons mediate the passed values in both directions.

The Backpropagation algorithm is used to learn the weights of a multilayer neural network with a fixed architecture. It performs gradient descent to try to minimize the sum squared error between the network's output values and the given target values.

Figure 2 depicts the network components which affect a particular weight change. Notice that all the necessary components are locally related to the weight being updated. This is one feature of backpropagation that seems biologically plausible. However, brain connections appear to be unidirectional and not bidirectional as would be required to implement backpropagation.

2 Notation

For the purpose of this derivation, we will use the following notation:

- The subscript k denotes the output layer.
- The subscript j denotes the hidden layer.
- The subscript i denotes the input layer.

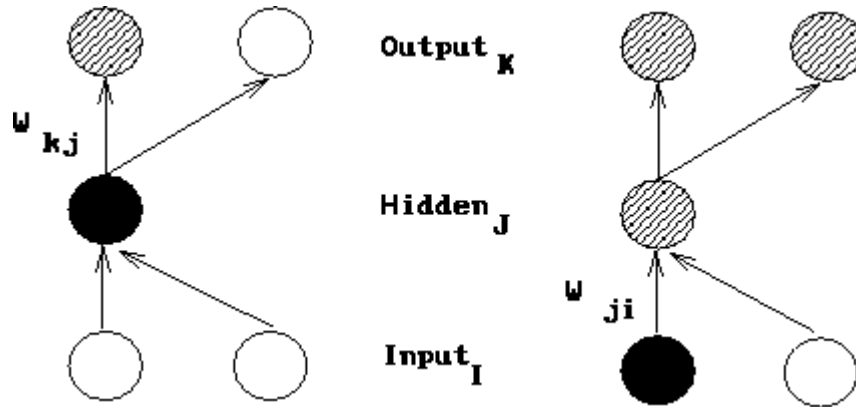


Figure 2: The change to a hidden to output weight depends on error (depicted as a lined pattern) at the output node and activation (depicted as a solid pattern) at the hidden node. While the change to a input to hidden weight depends on error at the hidden node (which in turn depends on error at all the output nodes) and activation at the input node.

- w_{kj} denotes a weight from the hidden to the output layer.
- w_{ji} denotes a weight from the input to the hidden layer.
- a denotes an activation value.
- t denotes a target value.
- net denotes the net input.

3 Review of Calculus Rules

$$\frac{d(e^u)}{dx} = e^u \frac{du}{dx} \quad \frac{d(g+h)}{dx} = \frac{dg}{dx} + \frac{dh}{dx} \quad \frac{d(g^n)}{dx} = ng^{n-1} \frac{dg}{dx}$$

4 Gradient Descent on Error

We can motivate the backpropagation learning algorithm as gradient descent on sum-squared error (we square the error because we are interested in its magnitude, not its sign). The total error in a network is given by the following equation (the $\frac{1}{2}$ will simplify things later).

$$E = \frac{1}{2} \sum_k (t_k - a_k)^2$$

We want to adjust the network's weights to reduce this overall error.

$$\Delta W \propto -\frac{\partial E}{\partial W}$$

We will begin at the output layer with a particular weight.

$$\Delta w_{kj} \propto -\frac{\partial E}{\partial w_{kj}}$$

However error is not directly a function of a weight. We expand this as follows.

$$\Delta w_{kj} = -\varepsilon \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial net_k} \frac{\partial net_k}{\partial w_{kj}}$$

Let's consider each of these partial derivatives in turn. Note that only one term of the E summation will have a non-zero derivative: the one associated with the particular weight we are considering.

4.1 Derivative of the error with respect to the activation

$$\frac{\partial E}{\partial a_k} = \frac{\partial(\frac{1}{2}(t_k - a_k)^2)}{\partial a_k} = -(t_k - a_k)$$

Now we see why the $\frac{1}{2}$ in the E term was useful.

4.2 Derivative of the activation with respect to the net input

$$\frac{\partial a_k}{\partial net_k} = \frac{\partial(1 + e^{-net_k})^{-1}}{\partial net_k} = \frac{e^{-net_k}}{(1 + e^{-net_k})^2}$$

We'd like to be able to rewrite this result in terms of the activation function. Notice that:

$$1 - \frac{1}{1 + e^{-net_k}} = \frac{e^{-net_k}}{1 + e^{-net_k}}$$

Using this fact, we can rewrite the result of the partial derivative as:

$$a_k(1 - a_k)$$

4.3 Derivative of the net input with respect to a weight

Note that only one term of the net summation will have a non-zero derivative: again the one associated with the particular weight we are considering.

$$\frac{\partial net_k}{\partial w_{kj}} = \frac{\partial(w_{kj}a_j)}{\partial w_{kj}} = a_j$$

4.4 Weight change rule for a hidden to output weight

Now substituting these results back into our original equation we have:

$$\Delta w_{kj} = \varepsilon \overbrace{(t_k - a_k)a_k(1 - a_k)}^{\delta_k} a_j$$

Notice that this looks very similar to the Perceptron Training Rule. The only difference is the inclusion of the derivative of the activation function. This equation is typically simplified as shown below where the δ term represents the product of the error with the derivative of the activation function.

$$\Delta w_{kj} = \varepsilon \delta_k a_j$$

4.5 Weight change rule for an input to hidden weight

Now we have to determine the appropriate weight change for an input to hidden weight. This is more complicated because it depends on the error at all of the nodes this weighted connection can lead to.

$$\Delta w_{ji} \propto - \left[\sum_k \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial net_k} \frac{\partial net_k}{\partial a_j} \right] \frac{\partial a_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}}$$

$$= \varepsilon \left[\sum_k \overbrace{(t_k - a_k) a_k (1 - a_k) w_{kj}}^{\delta_k} \right] a_j (1 - a_j) a_i$$

$$= \varepsilon \left[\sum_k \overbrace{\delta_k w_{kj}}^{\delta_j} \right] a_j (1 - a_j) a_i$$

$$\Delta w_{ji} = \varepsilon \delta_j a_i$$