

Make Your Own NN p80-p100

2021.2.2

Backpropagating Errors with Matrix Multiplication

The starting point is the errors that emerge from the NN at the final output layer.

$$\text{error}_{\text{output}} = \begin{pmatrix} e_1 \\ e_2 \end{pmatrix}$$

Next we want to construct the matrix for the hidden layer errors.

$$\text{error}_{\text{hidden}} = \begin{pmatrix} & \frac{w_{11}}{w_{11} + w_{21}} & \frac{w_{12}}{w_{12} + w_{22}} \\ \frac{w_{21}}{w_{21} + w_{11}} & & \frac{w_{22}}{w_{22} + w_{12}} \end{pmatrix} \cdot \begin{pmatrix} e_1 \\ e_2 \end{pmatrix}$$

You can see that the most important thing is the multiplication of the output errors e_n with the linked weights w_{ij} . The larger the weight, the more of the output error is carried back to the hidden layer. If we ignored that factor, we'd only lose the scaling of the errors being fed back. That is, $\frac{w_{11}}{w_{11}+w_{21}} * e_1$ would become the much simpler $w_{11} * e_1$.

The weight matrix we use here is actually the **transpose** matrix which we talked before. Then we have

$$\text{error}_{\text{hidden}} = W_{\text{hidden_output}}^T \cdot \text{error}_{\text{output}}$$

Key Points

- ★ Back-propagating the error can be expressed as a matrix multiplication.
- ★ Both feeding signals forward and error back-propagation can be made efficient using matrix calculations.

How do we actually update weights: Gradient Descent
Easy for amath student, skip :)

E is error function, w_{jk} is linked weight.

$$\begin{aligned}\frac{\partial E}{\partial w_{jk}} &= \frac{\partial}{\partial w_{jk}} \sum_n (t_n - o_n)^2 \\ \frac{\partial E}{\partial w_{jk}} &= \frac{\partial E}{\partial o_k} \cdot \frac{\partial o_k}{\partial w_{jk}} \\ &= -2(t_k - o_k) \cdot \frac{\partial o_k}{\partial w_{jk}} \\ &= -2(t_k - o_k) \cdot \frac{\partial}{\partial w_{jk}} \text{sigmoid}(\sum_j w_{jk} \cdot o_j)\end{aligned}$$

where $\frac{\partial}{\partial x} \text{sigmoid}(x) = \text{sigmoid}(x) \cdot (1 - \text{sigmoid}(x))$

$$\frac{\partial E}{\partial w_{jk}} = -(t_k - o_k) \cdot \text{sigmoid}(\sum_j w_{jk} \cdot o_j)(1 - \text{sigmoid}(\sum_j w_{jk} \cdot o_j)) \cdot o_j$$

$$\frac{\partial E}{\partial w_{jk}} = -(t_k - o_k) \cdot \text{sigmoid}(\sum_j w_{jk} \cdot o_j)(1 - \text{sigmoid}(\sum_j w_{jk} \cdot o_j)) + o_j$$

$$\text{new } w_{jk} = \text{old } w_{jk} - \alpha \cdot \frac{\partial E}{\partial w_{jk}}$$

Basically, it is Newton Method, where α is learning rate.

$$\begin{pmatrix} \Delta w_{1,1} & \Delta w_{1,2} & \Delta w_{1,3} & \dots \\ \Delta w_{2,1} & \Delta w_{2,2} & \Delta w_{2,3} & \dots \\ \Delta w_{3,1} & \Delta w_{3,2} & \Delta w_{3,3} & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix} = \begin{pmatrix} E_1 * S_1(1-S_1) \\ E_2 * S_2(1-S_2) \\ E_3 * S_3(1-S_3) \\ \dots \end{pmatrix} \cdot \begin{pmatrix} o_1 & o_2 & o_3 & \dots \end{pmatrix}$$

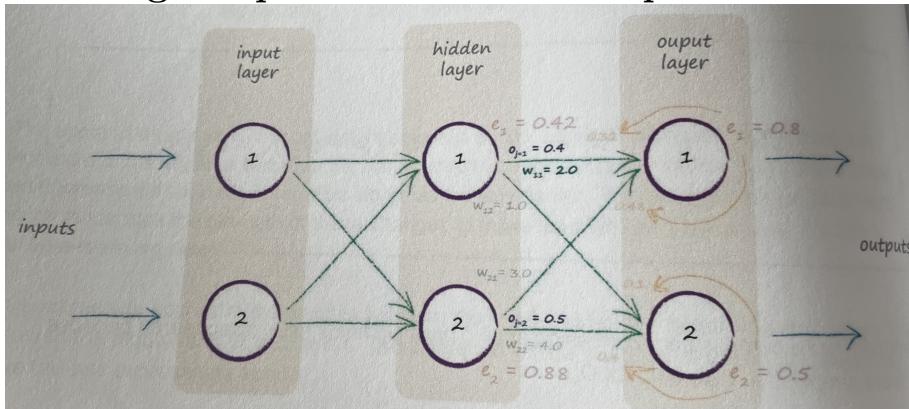
values from next layer values from previous layer

$$\Delta W_{jk} = \alpha \cdot E_k \cdot O_k(1 - O_k) \cdot O_j^T$$

Key Points

- * NN's error is a function of the internal link weights

Weight Update Worked Example



Suppose we want to update the weight w_{11} between the hidden and output layers, which currently has the value 2.0

$$\frac{\partial E}{\partial w_{jk}} = -(t_k - o_k) \cdot \text{sigmoid}(\sum_j w_{jk} \cdot o_j) (1 - \text{sigmoid}(\sum_j w_{jk} \cdot o_j)) \cdot o_j$$

- The first bit $(t_k - o_k)$ is the error $e_1 = 0.8$, just as we saw before.
- The sum inside the sigmoid functions $\sum_j w_{jk} \cdot o_j$ is $(2.0 * 0.4) + 3.0 * 0.5 = 2.3$
- The sigmoid $1/(1 + e^{-2.3})$ is then 0.909. That middle expression is then $0.909 * (1 - 0.909) = 0.083$
- $o_{j=1} = 0.4$

Multiplying all these three bits together and not forget minus sign at the start gives us -0.0265 . If we have learning rate $\alpha = 0.1$, then $-(0.1 * -0.0265) = +0.00265$. The new w_{11} is $2.0 + 0.00265 = 2.00265$