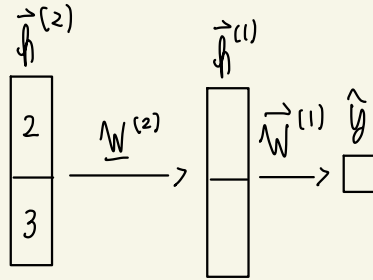Which can be rewritten into Martha's version, which I'll be using here:



Our network is initialized with:

$y$    True value of the labeled input        $y = 1$

$\hat{y}$    Model output        $\hat{y} = 0.2533$

$\vec{h}^{(1)} = \begin{bmatrix} h_1^{(1)} & h_2^{(1)} \end{bmatrix}^T$    Input values for our labeled input        $\vec{h}^{(1)} = \begin{bmatrix} 0.93 & 0.56 \end{bmatrix}^T$

$\vec{h}^{(2)} = \begin{bmatrix} h_1^{(2)} & h_2^{(2)} \end{bmatrix}^T$    Input values for our labeled input        $\vec{h}^{(2)} = \begin{bmatrix} 2 & 3 \end{bmatrix}^T$

$\vec{w}^{(1)} = \begin{bmatrix} w_1^{(1)} & w_2^{(1)} \end{bmatrix}$    Weights of the linear model at the end        $\vec{w}^{(1)} = \begin{bmatrix} 0.17 & 0.17 \end{bmatrix}$

$W^{(2)} = \begin{bmatrix} w_1^{(2)} & w_3^{(2)} \\ w_2^{(2)} & w_4^{(2)} \end{bmatrix}$    Weights in the input layer        $W^{(2)} = \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix}$

$\alpha$    The learning rate we use for backprop        $\alpha = 0.05$

$\ell : \mathbb{R}^2 \to \mathbb{R}$    Loss function for training        $\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

Since we'll need it later, we should write out our prediction y_hat in terms of the inputs

$$\hat{y} = \vec{W}^{(1)} \vec{h}^{(1)}$$

$$= \vec{W}^{(1)} W^{(2)} \vec{h}^{(2)}$$

$$= \begin{bmatrix} W_1^{(1)} & W_2^{(1)} \end{bmatrix} \begin{bmatrix} W_1^{(2)} & W_3^{(2)} \\ W_2^{(2)} & W_4^{(2)} \end{bmatrix} \begin{bmatrix} h_1^{(2)} \\ h_2^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} W_1^{(1)} & W_2^{(1)} \end{bmatrix} \begin{bmatrix} W_1^{(2)} h_1^{(2)} + W_3^{(2)} h_2^{(2)} \\ W_2^{(2)} h_1^{(2)} + W_4^{(2)} h_2^{(2)} \end{bmatrix}$$

$$= W_1^{(1)} \left( W_1^{(2)} h_1^{(2)} + W_3^{(2)} h_2^{(2)} \right) + W_2^{(1)} \left( W_2^{(2)} h_1^{(2)} + W_4^{(2)} h_2^{(2)} \right)$$

$$= W_1^{(1)} W_1^{(2)} h_1^{(2)} + W_1^{(1)} W_3^{(2)} h_2^{(2)} + W_2^{(1)} W_2^{(2)} h_1^{(2)} + W_2^{(1)} W_4^{(2)} h_2^{(2)}$$

$$\vec{W}_*^{(1)} = \vec{W}^{(1)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial \vec{w}^{(1)}}$$

$$= \vec{W}^{(1)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial \hat{y}} \, \frac{\partial \hat{y}}{\partial \vec{w}^{(1)}}$$

$$= \vec{W}^{(1)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial \hat{y}} \left[ \frac{\partial \hat{y}}{\partial w_1^{(1)}} \quad \frac{\partial \hat{y}}{\partial w_2^{(1)}} \right]^T$$

$$= \vec{W}^{(1)} - \alpha \, (\hat{y} - y) \left[ w_1^{(2)} h_1^{(2)} + w_3^{(2)} h_2^{(2)} \quad w_2^{(2)} h_1^{(2)} + w_4^{(2)} h_2^{(2)} \right]^T$$

$$= \vec{W}^{(1)} - 0.05 \, (-0.7467) \left[ 0.1581 \quad 0.0952 \right]^T$$

$$= \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix} + 0.0373 \begin{bmatrix} 0.1581 \\ 0.0952 \end{bmatrix}$$

$$= \begin{bmatrix} 0.1759 \\ 0.1736 \end{bmatrix}$$

$$\underline{W}_*^{(2)} = \underline{W}^{(2)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial \underline{w}^{(2)}}$$

$$= \underline{W}^{(2)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial (\vec{w}^{(1)} \vec{h}^{(1)})} \, \frac{\partial (\vec{w}^{(1)} \vec{h}^{(1)})}{\partial \underline{w}^{(2)}}$$

$$= \underline{W}^{(2)} - \alpha \, \frac{\partial l(\hat{y}, y)}{\partial (\vec{w}^{(1)} \vec{h}^{(1)})} \, \frac{\partial (\vec{w}^{(1)} \vec{h}^{(1)})}{\partial \underline{w}^{(2)}}$$

$$= \underline{W}^{(2)} - \alpha \left( \hat{y} - y \right) \begin{bmatrix} \partial \hat{y} / \partial w_1^{(2)} & \partial \hat{y} / \partial w_3^{(2)} \\ \partial \hat{y} / \partial w_2^{(2)} & \partial \hat{y} / \partial w_4^{(2)} \end{bmatrix}$$

$$= \underline{W}^{(2)} - \alpha \left( \hat{y} - y \right) \begin{bmatrix} w_1^{(1)} h_1^{(2)} & w_1^{(1)} h_2^{(2)} \\ w_2^{(1)} h_1^{(2)} & w_2^{(1)} h_2^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix} + 0.0373 \begin{bmatrix} 0.34 & 0.51 \\ 0.34 & 0.51 \end{bmatrix}$$

$$= \begin{bmatrix} 0.1327 & 0.1490 \\ 0.2427 & 0.1190 \end{bmatrix}$$

Now we do a forward pass with these new weights

$$\vec{h}^{(2)} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \qquad y = 1$$

$$\hat{y} = \vec{w}_*^{(1)} \underline{w}_*^{(2)} \vec{h}^{(2)}$$

$$= \begin{bmatrix} 0.1759 \\ 0.1736 \end{bmatrix} \begin{bmatrix} 0.1327 & 0.1490 \\ 0.2427 & 0.1190 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$= 0.2715$$

This is about 7% better than our previous output of 0.2533, so the backpropagation update clearly is working