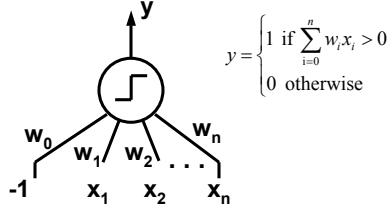


Neural Nets for Dummies

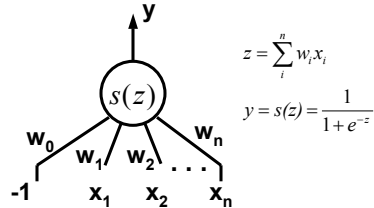
- Training: • Choose connection weights that minimize error
- Prediction: • Propagate input feature values through the network of "artificial neurons"
- Advantages: • Fast prediction
• Does feature weighting
• Very generally applicable
- Disadvantages: • Very slow training
• Overfitting is easy

Perceptron Unit



Creates a decision plane (line) in feature space

Sigmoid Unit



Creates a "soft" decision plane (line) in feature space

Adopted from Tomas Lozano Perez's 6.034 Recitation Notes

The simplest two-layer sigmoid Neural Net

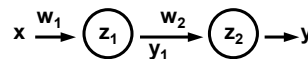
$$E = \frac{1}{2} (y^* - F(\bar{x}, \bar{w}))^2 \quad y = F(\bar{x}, \bar{w})$$

$\frac{\partial E}{\partial w_j} = -(y^* - y) \frac{\partial y}{\partial w_j}$

$y^* = \text{Desired output}$

$y = F(\bar{x}, \bar{w}) = s(z_2) = s(w_2 y_1) = s(w_2 s(z_1)) = s(w_2 s(w_1 x))$

Goal: find the weight vector that minimizes the error



Approach: Gradient Descent

(How does the error change as we twiddle with the weights?)

$$\frac{\partial y}{\partial w_2} = \frac{\partial s(z_2)}{\partial z_2} \frac{\partial z_2}{\partial w_2} = \frac{\partial s(z_2)}{\partial z_2} y_1 \quad \text{recall } z_2 = w_2 y_1 \text{ so, } \frac{\partial z_2}{\partial w_2} = y_1$$

$$\frac{\partial E}{\partial w_2} = \left[\frac{\partial s(z_2)}{\partial z_2} (y - y^*) \right] y_1 \delta_2$$

$$\frac{\partial y}{\partial w_1} = \frac{\partial s(z_2)}{\partial z_2} \frac{\partial z_2}{\partial s(z_1)} \frac{\partial s(z_1)}{\partial z_1} \frac{\partial z_1}{\partial w_1} = \frac{\partial s(z_2)}{\partial z_2} w_2 \frac{\partial s(z_1)}{\partial z_1} x \quad \text{recall } z_2 = w_2 s(z_1) \text{ so, } \frac{\partial z_2}{\partial s(z_1)} = w_2$$

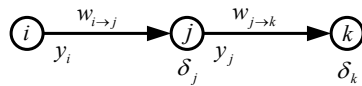
$$\frac{\partial E}{\partial w_1} = \left[\frac{\partial s(z_1)}{\partial z_1} \delta_2 w_2 \right] x \delta_1$$

recall $z_1 = w_1 x$ so, $\frac{\partial z_1}{\partial w_1} = x$

Descent rule:

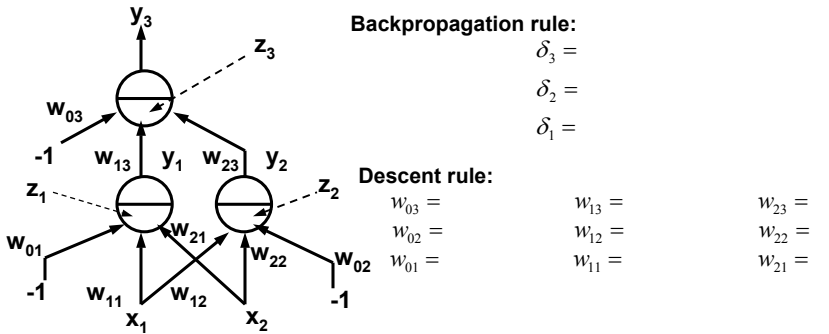
$$w_{i \rightarrow j} = w_{i \rightarrow j} - r \delta_j y_i$$

Backpropagation rule: $\delta_j = \frac{ds(z_j)}{dz_j} \sum_k \delta_k w_{j \rightarrow k}$



y_i is x_i for input layer

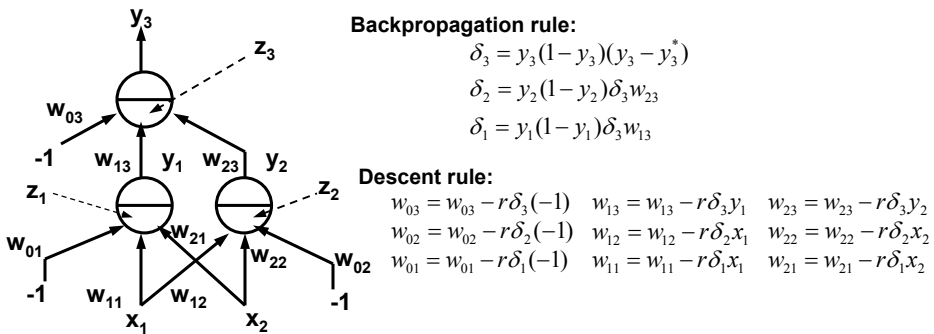
Example of Backpropagation



Initial Conditions: all weights are zero, learning rate is 8. Input: $(x_1, x_2) = (0, 1)$

$y^* =$	$\delta_3 =$	$w_{03} =$	$w_{13} =$
$z_1 =$	$\delta_2 =$	$w_{02} =$	$w_{12} =$
$z_2 =$	$\delta_1 =$	$w_{01} =$	$w_{11} =$
$y_1 =$			$w_{23} =$
$y_2 =$			$w_{22} =$
$z_3 =$			$w_{21} =$
$y_3 =$			

Example of Backpropagation



Initial Conditions: all weights are zero, learning rate is 8. Input: $(x_1, x_2) = (0, 1)$

$y^* = 1$	$\delta_3 = -1/8$	$w_{03} = -1$	$w_{13} = 1/2$
$z_1 = 0$	$\delta_2 = 0$	$w_{02} = 0$	$w_{12} = 0$
$z_2 = 0$	$\delta_1 = 0$	$w_{01} = 0$	$w_{11} = 0$
$y_1 = 1/2$			$w_{23} = 1/2$
$y_2 = 1/2$			$w_{22} = 0$
$z_3 = 0$			$w_{21} = 0$
$y_3 = 1/2$			