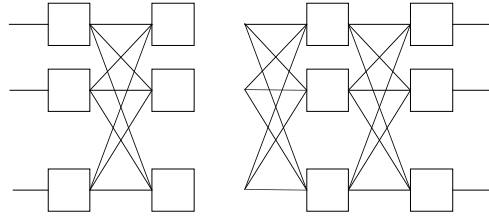


How to Train a MLP?

- With a single neuron, it is not too hard to see how to adjust the weights based upon the error values. We've already seen a couple of ways.
- With a multi-layer network, it is less obvious. For one thing, what is the "error" for the neurons in non-final layers? Without these, we don't know how to adjust.
- This is called the "credit assignment" problem (maybe should be "blame assignment").

Multi-Layer Network

Each box has a row-vector of weights and a bias.



Each layer has a matrix of weights and a column vector of biases.

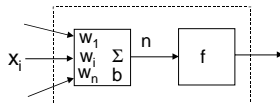
Multi-Layer Network

- Given an input vector, can compute the outputs.
- Given a sample, can compute the errors in output.
- Knowing gradient, can adjust the weights.
- Big Question:
How to compute the gradient?

Multi-Layer Network

- Recall that the gradient consists of components $\partial J / \partial w$ where J is the mean-squared error and w is some weight (or bias) in the network.
- For the Adaline, already derived: $\partial J / \partial w_i = -2 \epsilon x_i f'(n)$, where x_i is the input corresponding to weight w_i , and $n(\text{net})$ is the weighted sum. This works **as is** for the multi-layer case at the **output** layer.

Inside one neuron



$$\begin{aligned} \partial J / \partial w_i &= (\partial J / \partial n) (\partial n / \partial w_i) && \text{chain rule} \\ &= (\partial (d-f(n))^2 / \partial n) (\partial n / \partial w_i) \\ &= -2 \epsilon f'(n) x_i \\ &= s x_i \end{aligned}$$

where $s = (\partial J / \partial n)$ is called the *sensitivity*

Chain Rule Refresher

$$\frac{df(n(w))}{dw} = \frac{df(n)}{dn} \times \frac{dn(w)}{dw}$$

Example

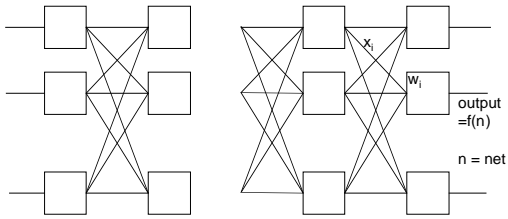
$$f(n) = \cos(n) \quad n = e^{2w} \quad f(n(w)) = \cos(e^{2w})$$

$$\frac{df(n(w))}{dw} = \frac{df(n)}{dn} \times \frac{dn(w)}{dw} = (-\sin(n))(2e^{2w}) = (-\sin(e^{2w}))(2e^{2w})$$

Application to Gradient Calculation

$$\frac{\partial J}{\partial w_{i,j}^m} = \frac{\partial J}{\partial n_i^m} \times \frac{\partial n_i^m}{\partial w_{i,j}^m}$$

Last Layer

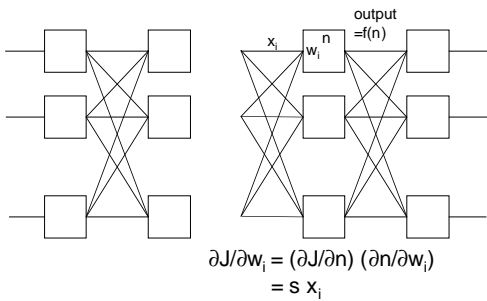


Utility of Sensitivity

$$\frac{\partial J}{\partial w_i} = \left(\frac{\partial J}{\partial n}\right) \left(\frac{\partial n}{\partial w_i}\right) = s x_i$$

anywhere in the network, not just at the final layer.

e.g. Next-Last Layer

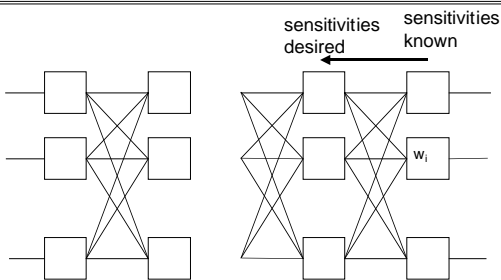


Utility of Sensitivity

Unlike at the output, it is less clear how to compute the sensitivity at an **arbitrary** neuron.

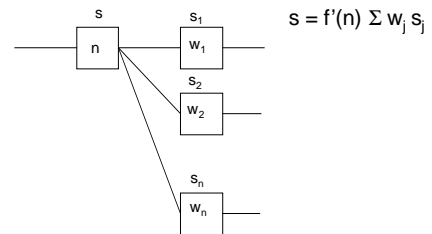
The key idea is to use an iterative approach that starts with the sensitivities at the last layer and works backward toward the first layer.

Backward Propagation of Sensitivity



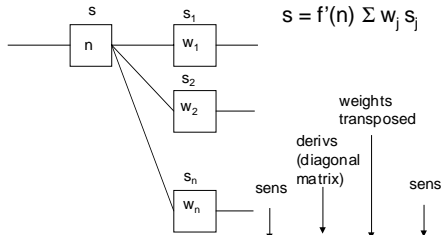
Backward Propagation of Sensitivity

Express desired as a weighted sum of known:



Vector Form

Express desired as a weighted sum of known:



Vector Form for entire layer : $s^m = \mathbf{F}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T s^{m+1}$

Correctness

Why should $s = f'(n) \sum w_j s_j$??

$$s = \frac{\partial J}{\partial n^m} = \sum (\frac{\partial n^{m+1}}{\partial n^m}) (\frac{\partial J}{\partial n^{m+1}})$$

from the vector form of the chain rule

$$s^m = \frac{\partial J}{\partial \mathbf{n}^m} = \left(\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} \right)^T \frac{\partial J}{\partial \mathbf{n}^{m+1}}$$

Vector Form Expansion

In the vector form $s^m = \frac{\partial J}{\partial \mathbf{n}^m} = \left(\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} \right)^T \frac{\partial J}{\partial \mathbf{n}^{m+1}}$

$$\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} \equiv \begin{bmatrix} \frac{\partial n_1^{m+1}}{\partial n_1^m} & \frac{\partial n_1^{m+1}}{\partial n_2^m} & \dots & \frac{\partial n_1^{m+1}}{\partial n_{s^m}^m} \\ \frac{\partial n_2^{m+1}}{\partial n_1^m} & \frac{\partial n_2^{m+1}}{\partial n_2^m} & \dots & \frac{\partial n_2^{m+1}}{\partial n_{s^m}^m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial n_{s^{m+1}}^{m+1}}{\partial n_1^m} & \frac{\partial n_{s^{m+1}}^{m+1}}{\partial n_2^m} & \dots & \frac{\partial n_{s^{m+1}}^{m+1}}{\partial n_{s^m}^m} \end{bmatrix}$$

the Jacobian

Correctness, continued

$$\begin{aligned} s^m &= \sum (\frac{\partial n^{m+1}}{\partial n^m}) (\frac{\partial J}{\partial n^{m+1}}) \\ &= (\frac{\partial}{\partial n^m}) \sum w_j f(n_j^m) s_j \quad \text{by def. of } n^{m+1} \\ &= \sum w_j f'(n^m) s_j \\ &= f'(n^m) \sum w_j s_j \end{aligned}$$

Vector Form

$$s^m = \frac{\partial J}{\partial \mathbf{n}^m} = \left(\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} \right)^T \frac{\partial J}{\partial \mathbf{n}^{m+1}} = \mathbf{F}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \frac{\partial J}{\partial \mathbf{n}^{m+1}}$$

$$s^m = \mathbf{F}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T s^{m+1}$$

$$\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} = \mathbf{W}^{m+1} \mathbf{F}^m(\mathbf{n}^m) \quad \mathbf{F}^m(\mathbf{n}^m) = \begin{bmatrix} f^m(n_1^m) & 0 & \dots & 0 \\ 0 & f^m(n_2^m) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f^m(n_{s^m}^m) \end{bmatrix}$$

Fully-Subscripted Alternatives to the Vector Forms

$$n_i^m = \sum_{j=1}^{s^{m-1}} w_{i,j}^{m-1} n_j^{m-1} + b_i^m$$

$$\frac{\partial n_i^m}{\partial w_{i,j}^m} = a_j^{m-1} \quad \frac{\partial n_i^m}{\partial b_i^m} = 1$$

Sensitivity

$$s_i^m \equiv \frac{\partial J}{\partial n_i^m}$$

Gradient

$$\frac{\partial J}{\partial w_{i,j}^m} = s_i^m a_j^{m-1} \quad \frac{\partial J}{\partial b_i^m} = s_i^m$$

Fully-Subscribed Alternatives to the Vector Forms

$$\frac{\partial n_i^{m+1}}{\partial n_j^m} = \frac{\partial \left(\sum_{l=1}^{s^m} w_{i,l}^{m+1} a_l^m + b_i^{m+1} \right)}{\partial n_j^m} = w_{i,j}^{m+1} \frac{\partial a_j^m}{\partial n_j^m}$$

$$\frac{\partial n_i^{m+1}}{\partial n_j^m} = w_{i,j}^{m+1} \frac{\partial f^m(n_j^m)}{\partial n_j^m} = w_{i,j}^{m+1} f^m(n_j^m)$$

$$f^m(n_j^m) = \frac{\partial f^m(n_j^m)}{\partial n_j^m}$$

Vector Form for Last Layer

$$s_i^M = \frac{\partial \hat{F}}{\partial n_i^M} = \frac{\partial (\mathbf{t} - \mathbf{a})^T (\mathbf{t} - \mathbf{a})}{\partial n_i^M} = \frac{\partial \sum_{j=1}^{s^M} (t_j - a_j)^2}{\partial n_i^M} = -2(t_i - a_i) \frac{\partial a_i}{\partial n_i^M}$$

$$\frac{\partial a_i}{\partial n_i^M} = \frac{\partial a_i^M}{\partial n_i^M} = \frac{\partial f^M(n_i^M)}{\partial n_i^M} = f^M(n_i^M)$$

$$s_i^M = -2(t_i - a_i) f^M(n_i^M)$$

$$\mathbf{s}^M = -2\mathbf{F}^M(\mathbf{n}^M)(\mathbf{t} - \mathbf{a})$$

Backpropagation Training Cycle

- **Forward propagation:** Derive the activation values (the inputs to the activation functions) at each neuron, and the final output.
- **Compute the error** in the output.
- **Backpropagate** the error through the network to get "sensitivities" at each neuron. (The gradient approximation is derivable from the sensitivities.)
- Use the sensitivities to **derive weight changes**.
- Apply the weight changes.

Backpropagation (Sensitivities)

The sensitivities are computed by starting at the last layer, and then propagating backwards through the network to the first layer.

$$\mathbf{s}^M \rightarrow \mathbf{s}^{M-1} \rightarrow \dots \rightarrow \mathbf{s}^2 \rightarrow \mathbf{s}^1$$

$$\mathbf{s}^M = -2\mathbf{F}^M(\mathbf{n}^M)(\mathbf{t} - \mathbf{a}) \quad \leftarrow \text{basis}$$

$$\mathbf{s}^m = \mathbf{F}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \mathbf{s}^{m+1} \quad \leftarrow \text{induction step}$$

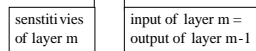
diagonal matrix of activation function derivative values

$$\mathbf{F}^m(\mathbf{n}^m) = \begin{bmatrix} f^m(n_1^m) & 0 & \dots & 0 \\ 0 & f^m(n_2^m) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f^m(n_{s^m}^m) \end{bmatrix}$$

Weight and Bias Update

Here we are using α instead of η for the learning rate.

$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T$$



$$\mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha \mathbf{s}^m$$

Fully-Subscribed Version of Weight Update

$$w_{i,j}^m(k+1) = w_{i,j}^m(k) - \alpha s_i^m a_j^{m-1} \quad b_i^m(k+1) = b_i^m(k) - \alpha s_i^m$$

$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T \quad \mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha \mathbf{s}^m$$

$$\mathbf{s}^m = \frac{\partial J}{\partial \mathbf{n}^m} = \begin{bmatrix} \frac{\partial J}{\partial n_1^m} \\ \frac{\partial J}{\partial n_2^m} \\ \vdots \\ \frac{\partial J}{\partial n_{s^m}^m} \end{bmatrix}$$

Summary

Forward Propagation

$$\mathbf{a}^0 = \mathbf{p}$$

$$\mathbf{a}^{m+1} = \mathbf{f}^{m+1}(\mathbf{W}^{m+1} \mathbf{a}^m) \quad m = 0, 2, \dots, M-1$$

$$\mathbf{a} = \mathbf{a}^M$$

Backpropagation

$$\mathbf{s}^M = -2\mathbf{F}^M(\mathbf{n}^M)(\mathbf{t} - \mathbf{a})$$

$$\mathbf{s}^m = \mathbf{F}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \mathbf{s}^{m+1} \quad m = M-1, \dots, 2, 1$$

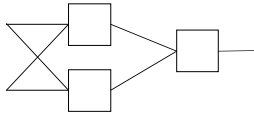
Weight Update

$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T \quad \mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha \mathbf{s}^m$$

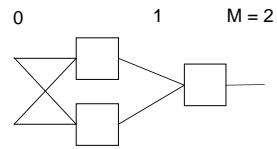
Exercise

- Derive the backprop equations symbolically for a simple network.
- Then use the equations to train the network.

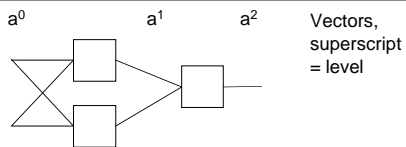
A simple network



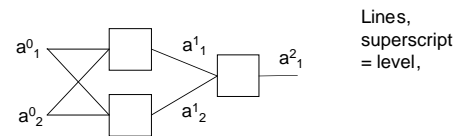
Label the Levels



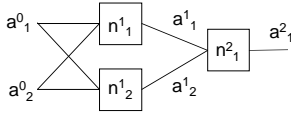
Label the Signal Vectors or Lines



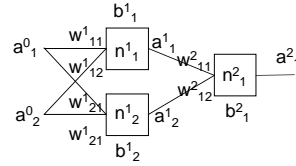
Label the Signal Vectors or Lines



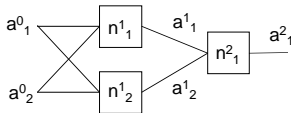
Label the Net (Activation) Values



Label the Weights and Biases



Write the forward equations for activations



$$n^1_1 = w^1_{11} a^0_1 + w^1_{12} a^0_2 + b^1_1$$

$$a^1_1 = f(n^1_1)$$

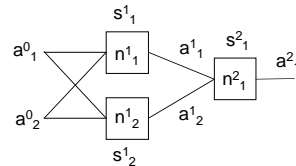
$$n^1_2 = w^1_{21} a^0_1 + w^1_{22} a^0_2 + b^1_2$$

$$a^1_2 = f(n^1_2)$$

$$n^2_1 = w^2_{11} a^1_1 + w^2_{12} a^1_2 + b^2_1$$

$$a^2_1 = f(n^2_1)$$

Write the backward equations for sensitivities



$$s^1_1 = w^2_{11} s^2_1 f'(n^1_1)$$

$$s^1_2 = w^2_{12} s^2_1 f'(n^1_2)$$

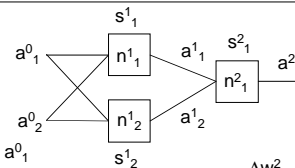
$$s^2_1 = -2(d^2_1 - a^2_1) f'(n^2_1)$$

↑ desired ↑ actual

Note

- The summations for the backpropagated sensitivities have only one term in this example, since the following layer has only one neuron.

Write the Equations for Weight and Bias Update



$$\Delta w^1_{11} = -\alpha s^1_1 a^0_1$$

$$\Delta w^1_{12} = -\alpha s^1_1 a^0_2$$

$$\Delta b^1_1 = -\alpha s^1_1$$

$$\Delta w^1_{21} = -\alpha s^1_2 a^0_1$$

$$\Delta w^1_{22} = -\alpha s^1_2 a^0_2$$

$$\Delta b^1_2 = -\alpha s^1_2$$

$$\Delta w^2_{11} = -\alpha s^2_1 a^1_1$$

$$\Delta w^2_{12} = -\alpha s^2_1 a^1_2$$

$$\Delta b^2_1 = -\alpha s^2_1$$

Exercises

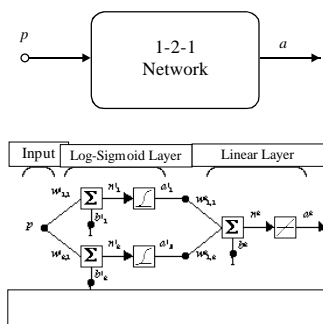
- Use the derivation to train the network to realize xor, where the first layer activations are logsig and the second layer is a hardlim.
- Carry out the preceding type of derivation for a 2-3-2 network (2 inputs, 3 middle neurons, 2 outputs).

Note on Training vs. Use

- Discontinuous functions such as hardlim, hardlims, etc. don't have derivatives.
- Therefore we train the network with continuous approximations to these functions, then replace them with the discontinuous versions during usage:

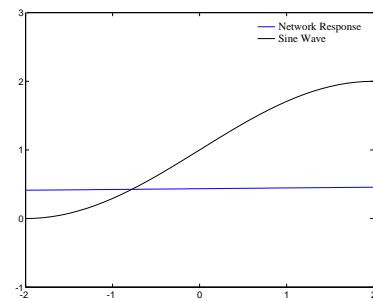
<u>usage</u>	<u>train with</u>
hardlim	logsig
hardlims	tansig

Numeric Example



Initial Conditions

$$\mathbf{W}^1(0) = \begin{bmatrix} -0.27 \\ -0.41 \end{bmatrix} \quad \mathbf{b}^1(0) = \begin{bmatrix} -0.48 \\ -0.13 \end{bmatrix} \quad \mathbf{W}^2(0) = [0.09 \ -0.17] \quad \mathbf{b}^2(0) = [0.48]$$



Forward Propagation

$$a^0 = p = 1$$

$$\mathbf{a}^1 = \mathbf{f}^1(\mathbf{W}^1 \mathbf{a}^0 + \mathbf{b}^1) = \text{logsig} \left(\begin{bmatrix} -0.27 \\ -0.41 \end{bmatrix} [1] + \begin{bmatrix} -0.48 \\ -0.13 \end{bmatrix} \right) = \text{logsig} \left(\begin{bmatrix} -0.75 \\ -0.54 \end{bmatrix} \right)$$

$$\mathbf{a}^1 = \begin{bmatrix} \frac{1}{1 + e^{0.75}} \\ \frac{1}{1 + e^{0.54}} \end{bmatrix} = \begin{bmatrix} 0.321 \\ 0.368 \end{bmatrix}$$

$$a^2 = \mathbf{f}^2(\mathbf{W}^2 \mathbf{a}^1 + \mathbf{b}^2) = \text{purelin} \left([0.09 \ -0.17] \begin{bmatrix} 0.321 \\ 0.368 \end{bmatrix} + [0.48] \right) = [0.446]$$

$$e = t - a = \left\{ 1 + \sin\left(\frac{\pi}{4}p\right) \right\} - a^2 = \left\{ 1 + \sin\left(\frac{\pi}{4}\right) \right\} - 0.446 = 1.261$$

Transfer Function Derivatives

$$f^1(n) = \frac{d}{dn} \left(\frac{1}{1 + e^{-n}} \right) = \frac{-e^{-n}}{(1 + e^{-n})^2} = \left(1 - \frac{1}{1 + e^{-n}} \right) \left(\frac{1}{1 + e^{-n}} \right) = (1 - a^1)(a^1)$$

$$f^2(n) = \frac{d}{dn}(n) = 1$$

Backpropagation

$$\mathbf{s}^2 = -2\mathbf{F}'(\mathbf{a}^2)(\mathbf{t} - \mathbf{a}) = -2\left[f'(a^2)\right](1,261) = -2\left[1\right](1,261) = -2,522$$

$$\mathbf{s}^1 = \mathbf{F}'(\mathbf{a}^1)(\mathbf{W}^2)^T \mathbf{s}^2 = \begin{bmatrix} (1-a_1^1)(a_1^1) & 0 \\ 0 & (1-a_2^1)(a_2^1) \end{bmatrix} \begin{bmatrix} 0,09 \\ -0,17 \end{bmatrix} \begin{bmatrix} 2,522 \end{bmatrix}$$

$$\mathbf{s}^1 = \begin{bmatrix} (1-0,321)(0,321) & 0 \\ 0 & (1-0,368)(0,368) \end{bmatrix} \begin{bmatrix} 0,09 \\ -0,17 \end{bmatrix} \begin{bmatrix} 2,522 \end{bmatrix}$$

$$\mathbf{s}^1 = \begin{bmatrix} 0,218 & 0 \\ 0 & 0,233 \end{bmatrix} \begin{bmatrix} -0,227 \\ 0,429 \end{bmatrix} = \begin{bmatrix} -0,0495 \\ 0,0997 \end{bmatrix}$$

Weight Update

$$\alpha = 0.1$$

$$\mathbf{W}^2(1) = \mathbf{W}^2(0) - \alpha \mathbf{s}^2 (\mathbf{a}^1)^T = \begin{bmatrix} 0,09 & -0,17 \end{bmatrix} - 0,1 \begin{bmatrix} -2,522 \end{bmatrix} \begin{bmatrix} 0,321 & 0,368 \end{bmatrix}$$

$$\mathbf{W}^2(1) = \begin{bmatrix} 0,171 & -0,0772 \end{bmatrix}$$

$$\mathbf{b}^2(1) = \mathbf{b}^2(0) - \alpha \mathbf{s}^2 = \begin{bmatrix} 0,48 \end{bmatrix} - 0,1 \begin{bmatrix} -2,522 \end{bmatrix} = \begin{bmatrix} 0,732 \end{bmatrix}$$

$$\mathbf{W}^1(1) = \mathbf{W}^1(0) - \alpha \mathbf{s}^1 (\mathbf{a}^0)^T = \begin{bmatrix} -0,27 \\ -0,41 \end{bmatrix} - 0,1 \begin{bmatrix} -0,0495 \\ 0,0997 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \begin{bmatrix} -0,265 \\ -0,420 \end{bmatrix}$$

$$\mathbf{b}^1(1) = \mathbf{b}^1(0) - \alpha \mathbf{s}^1 = \begin{bmatrix} -0,48 \\ -0,13 \end{bmatrix} - 0,1 \begin{bmatrix} -0,0495 \\ 0,0997 \end{bmatrix} = \begin{bmatrix} -0,475 \\ -0,140 \end{bmatrix}$$