

Neil E. Cotter

When ~~in~~ learning ~~patterns~~ with neural networks, we want to minimize average error over all possible input patterns. We define the true average error to be

$$\bar{E} = \frac{1}{A} \int_{\text{input patterns } p} \frac{(t_p - y_p)^2}{2} dp$$

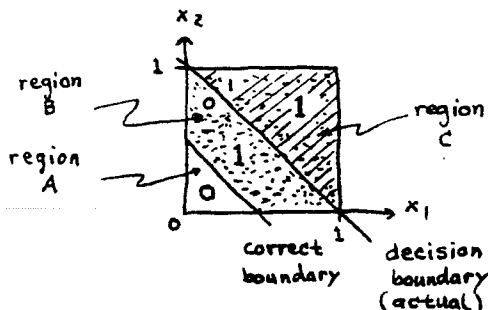
where  $A \equiv$  area of input pattern space

$p \equiv$  particular input pattern

$t_p \equiv$  desired output for pattern  $p$

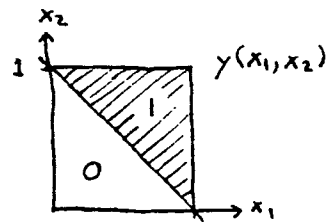
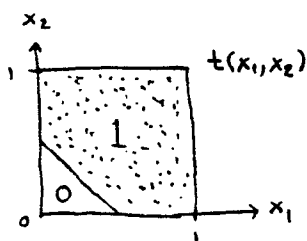
$y_p \equiv$  actual " " " "

ex:



Perceptron with incorrect decision boundary. Error in dotted, unhatched region, (B).

$$\bar{E} = \frac{1}{[0,1] \times [0,1]} \iint_{0,0}^{1,1} \left[ \frac{t(x_1, x_2) - y(x_1, x_2)}{2} \right]^2 dx_1 dx_2$$



$$\begin{aligned} \therefore \bar{E} &= \frac{1}{1} \left[ \iint_A \left( \frac{0-0}{2} \right)^2 dx_1 dx_2 + \iint_B \left( \frac{1-0}{2} \right)^2 dx_1 dx_2 + \iint_C \left( \frac{1-1}{2} \right)^2 dx_1 dx_2 \right] \\ &= \iint_B \frac{1}{2} dx_1 dx_2 = \frac{1}{2} \cdot \frac{3}{8} = \frac{3}{16} \end{aligned}$$

Apr 1990

Neil E Cotter

Gradient Descent - Error Surface Estimation (cont.)

In most neural network applications, we do not have access to  $\bar{E}$ . This is because  $\bar{E}$  is really  $E(\vec{w})$  and changes value each time we change  $\vec{w}$ . To find the new value of  $\bar{E}$  after each change in  $\vec{w}$ , we would have to present every input pattern to the network and measure the average output error. This approach is impractical.

Instead, we estimate  $E(\vec{w})$  in a crude way:

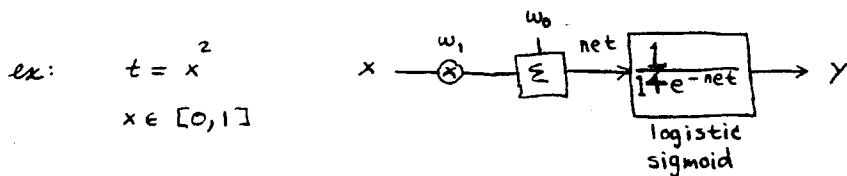
$$\bar{E}(\vec{w}) \approx E(\vec{x}) \quad \text{where} \quad E(\vec{x}) = \frac{1}{2}(t-y)^2$$

$t \equiv$  target output for present input,  $\vec{x}$ .  
 $y \equiv$  actual " " " " " "

In other words, we use the error,  $E(\vec{x})$ , for a single input pattern,  $\vec{x}$ , to estimate the entire error surface. This works fairly well, provided we pick values of  $\vec{x}$  at random as gradient descent proceeds.

Note: error acts like noise that helps avoid local min (or  $\nabla_{\vec{w}} E(\vec{x})$ )

We use  $E(\vec{x})$  because we can easily calculate  $\frac{\partial E(\vec{x})}{\partial \vec{w}}$ , the gradient of  $E$ .



$$\bar{E}(\vec{w}) = \frac{1}{1} \int_0^1 \left[ \frac{t(x)-y(x)}{2} \right]^2 dx = \int_0^1 \left[ x^2 - \frac{1}{1+e^{-w_0-w_1x}} \right]^2 dx$$

Very hard to compute

$$E(\vec{x}) = \left[ x^2 - \frac{1}{1+e^{-w_0-w_1x}} \right]^2$$

$$\frac{\partial E(\vec{x})}{\partial w_1} = \underbrace{f'(x)}_{-f(x)[1-f(x)]} \cdot \underbrace{(t-y)}_{[x^2 - f(x)]} w_1$$

Easy comp

where  $f(x) = \frac{1}{1+e^{-w_0-w_1x}}$