

# Backpropagation

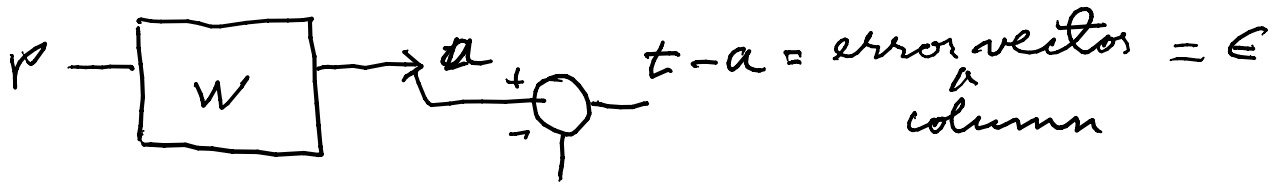
summary on pp. 11-24-25

EE 434  
02/08/05

scheme for training

$t$  = desired outputs for an input  $p$

$a$  = output of net for input  $p$



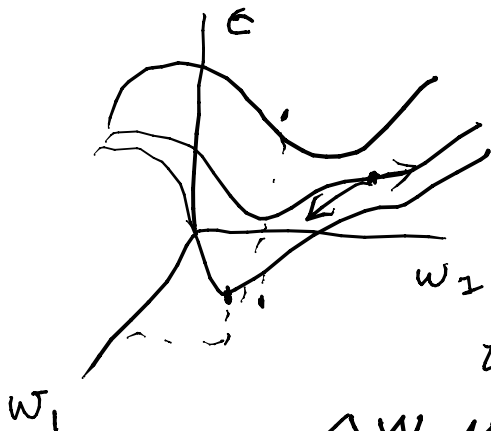
form  $(t-a)^T(t-a)$  ;  $V^T$  means transpose  
= sum of squares

adjust the weights (including biases) to  
try to minimize  $E = (t-a)^T(t-a)$

use Newton iteration, index by integers  $k$   
for  $k$ th step

$$W(k+1) = W(k) + \Delta W(k+1)$$

"change in weights"

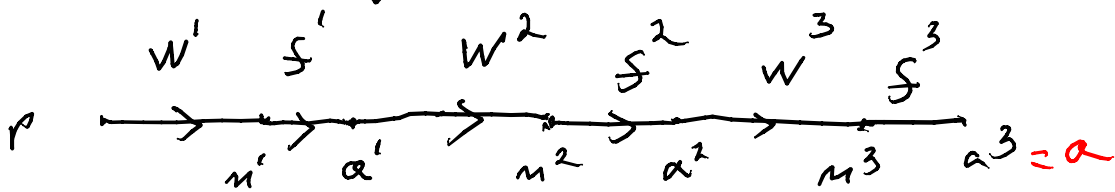


desire smallest  $E$  by  
choice of  $W \rightarrow \Delta W$   
makes  $E$  smaller

to go downhill desire  
 $\Delta W$  proportional to  $-\nabla_W E$

as desire to not overshoot the minimum of  $E$  put in a scale factor  $\alpha$ ,  $0 < \alpha < 1$

This works ok if all the weights directly in  $E$ , but for many sectioned feed forward networks  $E$  is not a direct function of the layer weights of internal layers. So have backpropagation.



$$a^3 = f^3(n^3) = f^3(W^3 a^2) = f^3(W^3 f^2(n^2))$$

$$= f^3(W^3 f^2(W^2 a^1)) = f^3(W^3 f^2(W^2 f^1(W^1 p)))$$

desire to set the  $W^1$ ,  $\frac{\partial E}{\partial W^1_{i,j}}$

$$\frac{\partial E}{\partial W^3_{i,j}} \Rightarrow (t-a)^T(t-a) = E = t^T t - \underbrace{a^T t - t^T a}_{\text{equal as } (a^T t)^T = t^T a^T} + a^T a$$

$$E = \underbrace{t^T t}_{\text{fixed by choice}} - 2t^T a + a^T a$$

$$a = f^3(n^3)$$

$$\frac{\partial E}{\partial W^3_{i,j}} = 0 - 2t^T \frac{\partial a}{\partial W^3_{i,j}} + \frac{\partial a^T}{\partial W^3_{i,j}} \cdot a + a^T \frac{\partial a}{\partial W^3_{i,j}}$$

$$= (-2t^T + 1a^T) \frac{\partial a}{\partial w_{i,j}^3}; \quad \frac{\partial a}{\partial w_{i,j}^3} = \sum_{k=1}^3 \frac{\partial f^3(m_k^3)}{\partial m_k^3} \cdot \frac{\partial m_k^3}{\partial w_{i,j}^3}$$

$$m^3 = W^3 \cdot a^2 \Rightarrow m_k^3 = \sum_{j=1}^2 W_{k,j}^3 a_j^2 \text{ only get if } k=i$$

$$\text{here } m_i^3 = \sum_{l=1}^2 W_{i,l}^3 a_l^2, \quad \frac{\partial m_i^3}{\partial w_{i,j}^3} = 1 \cdot a_j^2 = a_j^2$$

$$\frac{\partial E}{\partial w_{i,j}^3} = -2(t-a)^T \cdot \frac{\partial a_{i,j}^3}{\partial w_{i,j}^3} = -2e^T \cdot \frac{\partial f^3(m_i^3)}{\partial m_i^3} \cdot a_j^2$$

known from running on p.

need the activation function  $f_i(\cdot)$

$$\text{Ex: } f_i(x) = \tanh x$$

$$\text{need } \frac{\partial f_i(x)}{\partial x} = \frac{df_i(x)}{dx} = \frac{d}{dx} \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right)$$

$$= \frac{e^x + e^{-x}}{e^x + e^{-x}} - \frac{(e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} = 1 - (\tanh x)^2$$

$$\text{in our case } a_i^3 = \text{known}, \quad \frac{\partial a_i^3}{\partial m_i^3} = 1 - (a_i^3)^2$$

understood as only one activation function on each neuron

works well for output layer, but  $\frac{\partial E}{\partial w_{i,j}^2}, \frac{\partial E}{\partial w_{i,j}^1}$

are not so easy to get as don't know desired outputs for 1st & second layer

$M \leftarrow$  last one  
 $\Delta = -2 \frac{\partial \dot{F}^M(n^M)}{\partial n^M}(t-a) =$  sensitivity vectors for last  $(M)$  layers  
*derivative*

$$\Delta^m = \dot{F}^m(n^m) (W^{m+1})^T \Delta^{m+1} \quad \text{for } m = M-1, \dots, 1$$

p.11-25

$$F^m = \begin{bmatrix} f(n_{i_1}^m) \\ \vdots \\ f(n_{s_m}^m) \end{bmatrix} = \text{diag}[f(n_{i_1}^m)]$$

$$\dot{F}^m = \begin{bmatrix} \frac{\partial f(n_{i_1}^m)}{\partial n_{i_1}^m} \\ \vdots \\ \frac{\partial f(n_{s_m}^m)}{\partial n_{s_m}^m} \end{bmatrix}$$

$S^m =$  size of the  $m$ th  $n^m$  vector  
 $=$  # of neurons in  $m$ th layer

$$W^m(k+1) = W^m(k) - \alpha \Delta^m (\Delta^{m-1})^T \quad ; \quad 0 < \alpha < 1$$

up to us to choose

