Gradient Descent — BEP —

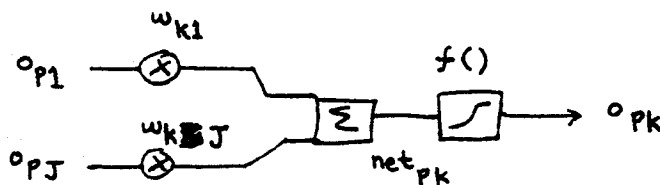5 May 1989
Neil E. Cotter

Chain rule and delta rule

We need to calculate ♥ $\Delta_p w_{ij} = -\eta \dfrac{\partial E_p}{\partial w_{ij}}$ .

$w_{ij}$ can be anywhere in a neural neural network of, say, three layers.



Each circle represents one neuron, complete with synaptic weights.

Consider first a neuron on the <u>output layer</u>:



$$\Delta_p w_{kj} = -\eta \frac{\partial E_p}{\partial w_{kj}}$$

1) $\quad \dfrac{\partial E_p}{\partial w_{kj}} = \dfrac{\partial}{\partial w_{kj}} \; \dfrac{1}{2} \left( t_{pk} - o_{pk} \right)^2$

     ↑ desired output    ↖ actual output

$$= \left( t_{pk} - o_{pk} \right) \frac{\partial}{\partial w_{kj}} \left( t_{pk} - o_{pk} \right)$$

• by chain rule
$\dfrac{\partial}{\partial x} f(g(x)) = f'(g(x)) \, g'(x)$

$$= -\left( t_{pk} - o_{pk} \right) \frac{\partial}{\partial w_{kj}} o_{pk} .$$

Now use the chain rule to compute $\dfrac{\partial}{\partial w_{kj}} o_{pk}$.

May 1989
Neil E. Cotter

2) $\quad \dfrac{\partial}{\partial w_{kj}} \, o_{Pk} = \dfrac{\partial}{\partial w_{kj}} \, f(net_{Pk})$

$$= f(net_{Pk})\left[1 - f(net_{Pk})\right] \dfrac{\partial}{\partial w_{kj}} net_{Pk}$$

- Recall that $f'(x) = f(x)\left[1 - f(x)\right]$
- Also use chain rule

Now compute $\dfrac{\partial}{\partial w_{kj}} net_{Pk}$:

3) $\quad \dfrac{\partial}{\partial w_{kj}} net_{Pk} = \dfrac{\partial}{\partial w_{kj}}\left(\displaystyle\sum_{j=1}^{J} w_{kj}\, o_{Pj} + \Theta_k\right)$

$$= \dfrac{\partial}{\partial w_{kj}}\left(\sum_{j=1}^{J} w_{kj}\, o_{Pj}\right)$$

$$= \; ?$$

Here we have to be careful about notation. "$\dfrac{\partial}{\partial w_{kj}}$" refers to a particular synaptic weight such as $w_{23}$, say. "$\sum w_{kj}\, o_{Pj}$" refers to the sum over various different $w_{kj}$ synapse values.

$\dfrac{\partial}{\partial w_{kj}} w_{kj} = 0$ unless both $j$'s have the same value.

$\therefore \quad \dfrac{\partial}{\partial w_{kj}} \displaystyle\sum_{j=1}^{J} w_{kj}\, o_{Pj} = o_{Pj}$ 
- One term, $w_{kj}\, o_{Pj}$, is extracted from the $\Sigma$

ex: $\quad \dfrac{\partial}{\partial w_{12}} \displaystyle\sum_{j=1}^{2} w_{1j}\, o_{Pj} = \dfrac{\partial}{\partial w_{12}}\left(w_{11}\, o_{P1} + w_{12}\, o_{P2}\right) = \dfrac{\partial}{\partial w_{12}} w_{12}\, o_P$

$$= o_{P2}$$

$\therefore \quad \dfrac{\partial}{\partial w_{kj}} net_{Pk} = o_{Pj}$

May 1989
Neil E Cotter

Put results together:

    <u>Output layer</u> learning rule:

$$\nabla_p w_{kj} = -\eta \frac{\partial E_p}{\partial w_{kj}} = -\eta(-1)(t_{pk} - o_{pk})\, f(net_{pk})\left[1 - f(net_{pk})\right.$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}$$

$$\cdot o_{pj}$$

$$\uparrow$$

input from previous layer outputs

$$= \eta \underbrace{(t_{pk} - o_{pk})\, f(net_{pk})\left[1 - f(net_{pk})\right]}\; o_{pj}$$

$$|||$$

$$\delta_{pk}$$

- We identify a quantity, $\delta_{pk}$, that is a function of things that happen after signals pass through a synapse. The $\delta_{pk}$ information is therefore propagated <u>back</u> through the network to the synapse. Unfortunately, no physiological mechanism for this has been identified. But it works. It learns by gradient descent.

- The $o_{pj}$ is the input to the synapse, $(w_{kj})$ and $o_{pj}$ comes from the previous layer of the network.

- $\eta$ is the step size for the gradient descent.

A subtle point: Our gradient descent changes the weights of the network so as to minimize $E_p$ for a particular pattern $p$. What we really want to do is minimize the <u>average</u> error over all patterns: $\text{Ave } E_p = \frac{1}{P}\sum_{p=1}^{P} E_p$.

The theory is that using a small value for $\eta$ and changing to a different learning pattern $p$ at each step of the gradient descent will give small values of averaged error.

Gradient Descent — BEP —
Chain rule and delta rule (cont.)
'5 May 1989    Neurons in hidden layers, learning rule
Neil E. Cotter

- A hidden layer is any layer other than the output layer.

$$\Delta_p w_{ji} = -\eta \frac{\partial E_p}{\partial w_{ji}}$$

Proceed as before for steps (1) and (2). We hit a snag at step (3):

3)  $\dfrac{\partial}{\partial w_{ji}} net_{pk} = \dfrac{\partial}{\partial w_{ji}} \left( \displaystyle\sum_{j=1}^{J} w_{kj} \, o_{pj} + \theta_k \right)$

$$= \frac{\partial}{\partial w_{ji}} \left( \sum_{j=1}^{J} w_{kj} \, o_{pj} \right)$$

$= ?$   Here, $o_{pj}$ is a function of $w_{ji}$. Furthermore, changing $w_{ji}$ will change every output in the output layer since the output $o_{pj}$ goes to every neuron in the ~~out~~ output layer.

$\therefore$  Total error is  $E_p = \displaystyle\sum_{k=1}^{K} \frac{1}{2} \left( t_{pk} - o_{pk} \right)^2$

Start over.

1)  $\dfrac{\partial E_p}{\partial w_{ji}} = \dfrac{\partial}{\partial w_{ji}} \displaystyle\sum_{k=1}^{K} \frac{1}{2} \left( t_{pk} - o_{pk} \right)^2$

$$= \sum_{k=1}^{K} \left( t_{pk} - o_{pk} \right) \frac{\partial (-o_{pk})}{\partial w_{ji}}$$

$$= - \sum_{k=1}^{K} \left( t_{pk} - o_{pk} \right) \frac{\partial \, o_{pk}}{\partial w_{ji}}$$

2)  $\dfrac{\partial}{\partial w_{ji}} o_{pk} = f(net_{pk}) \left[ 1 - f(net_{pk}) \right] \dfrac{\partial}{\partial w_{ji}} net_{pk}$    as before

Gradient Descent — BEP —
Chain rule and delta rule (cont.)

3) $\quad \dfrac{\partial}{\partial w_{ji}} \ net_{pk} = \dfrac{\partial}{\partial w_{ji}} \left( \displaystyle\sum_{j=1}^{J} w_{kj} \, o_{pj} \right)$

$\qquad\qquad = \dfrac{\partial}{\partial w_{ji}} \ w_{kj} \, o_{pj}$  •  Picks out one term

$\qquad\qquad = w_{kj} \dfrac{\partial}{\partial w_{ji}} o_{pj}$  •  $w_{kj}$ is _not_ a function of $w_{ji}$.

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$  •  $o_{pj}$ _is_ a function of $w_j$

4) $\quad \dfrac{\partial}{\partial w_{ji}} o_{pj} = \dfrac{\partial}{\partial w_{ji}} \ f(net_{pj})$

This has the same form as a previous
derivation, and we get:

$$\dfrac{\partial}{\partial w_{ji}} o_{pj} = f(net_{pj}) \left[ 1 - f(net_{pj}) \right] o_{pi}$$

Putting the above results together, we have

Learning rule for neurons in next-to-last (penultimate) layer:

$$\Delta_p w_{ji} = + \eta \sum_{k=1}^{K} \left\{ (t_{pk} - o_{pk}) \ f(net_{pk}) \left[ 1 - f(net_{pk}) \right] w_{kj} \right.$$
$$\left. f(net_{pj}) \left[ 1 - f(net_{pj}) \right] o_{pi} \right\}$$

$$= \eta \ o_{pi} \ \underbrace{f(net_{pj}) \left[ 1 - f(net_{pj}) \right] \sum_{k=1}^{K} (t_{pk} - o_{pk}) \ f(net_{pk}) \left[ 1 - f(net_{pk}) \right]}_{\substack{||| \\ \delta_{pj}}}$$

•  Same general form as output layer learning rule
•  Notice that we have used the chain rule to propagate
the effects of output errors back through the network.
Hence the name "Backward Error Propagation."

Gradient Descent — BEP—
Chain rule and delta rule (cont.)

Now for the amazing result:

$$\delta_{Pj} = f(net_{Pj})\left[1 - f(net_{Pj})\right] \sum_{k=1}^{K} \delta_{Pk}\, w_{kj}$$

Amazing? Yes. We can compute the $\delta$'s for this layer if we know the $\delta$'s for the next layer. So start at the output layer and work backwards, calculating $\delta$'s as you go.

Even more amazing is the similarity between the calculation of $\delta$'s and the calculations performed by the neural net. We have weighted sums of $\delta$'s just as we have weighted sums of inputs in the forward direction. The backward error propagation computation is nearly equivalent in complexity to the forward computation performed by the network.

Best of all, learning requires only information that a synapse might actually have access to: $o_{Pi}$, the incoming signal from the previous layer $\delta_{Pj}$, the $\delta_{Pj}$ is the same for all synapses on this neuron (neuron j) and $\delta_{Pj}$ is computed from $\delta_{Pk}$ on the next layer and from $w_{kj}$ synapses connecting this neuron to the next layer.

By recursion we get the general delta rule
for Backward Error Propagation:

$$\nabla_P w_{ji} = \eta\, \delta_{Pj}\, o_{Pi} \qquad \text{where} \qquad \delta_{Pj} = \sum_{k=1}^{K} \delta_{Pk}\, w_{kj}$$

$\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}$ Sum over next layer.

$o_{Pi}$ = previous layer's output
 is this layer's input

6 May 1988     Learning   Rule        · Backward Error Propagation

Neil E Cotter                               · Is a ~~type~~ of delta rule

$$w_{ji} \text{ new} = w_{ji} \text{ old} + \underbrace{\Delta_p w_{ji}}$$

synapse connects neuron i in previous layer to neuron j in this layer

change in weight when input pattern p is training vector

Simplistic δ rule     $\overbrace{\text{change in weight}}$

Use     $\Delta_p w_{ji} = \eta \, (t_{pj} - o_{pj}) \, o_{pi} \equiv \eta \, \delta_{pj} \, o_{pi}$

$\underbrace{\phantom{(t_{pj} - o_{pj})}}$
desired out
— actual out
$= \text{error} = \delta_{pj}$

to neuron j
input from previous layer neuron i

· $\eta$ is a constant that controls step size in this gradient descent algorithm.

(more positive)

Now   $o_{pj} \propto w_{ji} \, o_{pi}$     · ∴ if $o_{pi} > 0$ then larger $w_{ji} \Rightarrow$ more pos
                                      · i.e. if $o_{pi} > 0$ and $w_{ji} \uparrow$ then $o_{pj}$

| In | | out | error that gave $\Delta w_{ji}$ | |
|---|---|---|---|---|
| $o_{pi}$ | $\Delta w_{ji}$ | $o_{pj}$ | $t_{pj} - o_{pj}$ | result: $\lvert t_{pj} - o_{pj} \rvert$ |
| $>0$ | $\uparrow$ | $\uparrow$ | $>0$ | $\downarrow$ |
| $>0$ | $\downarrow$ | $\downarrow$ | $<0$ | $\downarrow$ |
| $<0$ | $\downarrow$ | $\uparrow$ | $>0$ | $\downarrow$ |
| $<0$ | $\uparrow$ | $\downarrow$ | $<0$ | $\downarrow$ |

· move this col here; makes more sense

$\lvert \text{error} \rvert$ reduced in every case

Generalized δ rule     (·One actually used in Backward Error Prop.

def:     $E_p \equiv \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2 \equiv$ total squared error for pattern p as input
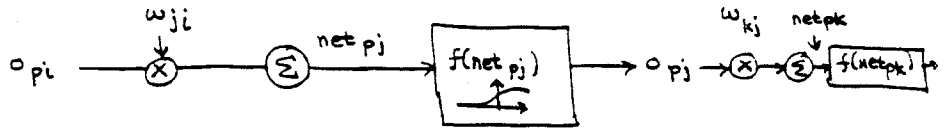
$\sum$ over all output layer neurons

· Remember that   $t_{pj} \equiv$ desired (or target) output for neur
              $o_{pj} \equiv$ actual output for neuron j

Then  use   $\delta_{pj} \equiv - \partial E_p / \partial net_{pj}$

Gradient Descent — BEP —
Chain rule and delta rule (cont.)

May 1988

Neil E. Cotter

Calculation of $\delta_{pj}$

prelims:



· Next layer (if $o_{pj}$ not output lay

eqn's:

$$net_{pj} = \sum_i w_{ji} o_{pi} + \theta_j$$

$$o_{pj} = f(net_{pj})$$

$$E_p = \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2$$

$\delta_{pj}$ for output-layer neuron $j$

$$\delta_{pj} \equiv -\frac{\partial E_p}{\partial net_{pj}} = -\frac{\partial E_p}{\partial o_{pj}} \frac{\partial o_{pj}}{\partial net_{pj}}$$

· Use the chain rule for differentiation

1)  $$\frac{\partial E_p}{\partial o_{pj}} = \frac{\partial \frac{1}{2}\sum_j (t_{pj}-o_{pj})^2}{\partial o_{pj}} \overset{?}{=} \sum_j \frac{\partial \frac{1}{2}(t_{pj}-o_{pj})^2}{\partial o_{pj}}$$  · Be careful

$$= \frac{1}{2} 2 (t_{pj}-o_{pj})(-1)$$

$$= -(t_{pj}-o_{pj})$$

· The only term in the $\sum$ th isn't constant relative to $o_{pj}$ is the $j^{th}$ term.

· Note the difference between $o_{pj}$ inside and outside the $\sum$. Inside the $\sum$ the $j$ is a dummy variable (similar to dummy variables of integra

2)  $$\frac{\partial o_{pj}}{\partial net_{pj}} = \frac{\partial f(net_{pj})}{\partial net_{pj}} = f' \Big|_{net_{pj}} = f(net_{pj})(1- f(net_{pj}))$$

$$= o_{pj}(1-o_{pj})$$

· Recall that for $f(u) = 1/(1+e^{-u})$, $\partial f/\partial u = f(1-f)$

∴ Output layer learning rule:

$$\overbrace{\qquad\qquad}^{\delta_{pj}}$$

$$\Delta_p w_{ji} = \eta\, \delta_{pj}\, o_{pi} = \eta\, (t_{pj}-o_{pj})\, o_{pj}(1-o_{pj})\, o_{pi}$$