Gradient Descent - BEP - Algorithm

Engistic Networks (Backward Error Propagation

 $\Delta w_{ji} = -\gamma \frac{\partial E}{\partial w_{ji}}$ (Gradient Descent

Equivalent to:

Mar 1990 Neil E Cotten

 $\Delta \omega_{ji} = \gamma \delta_{j} o_{i}$ (Delta Rule)

 $\delta_j = f'(\text{net}_j) \sum_{k} \delta_k w_{kj} = o_j(1-o_j) \sum_{k} \delta_k w_{kj}$

oj = output of j $\delta_{k} = \delta$ for next layer downstream, neuron k $w_{kj} = synaptic$ weights from neuron j

to neuron k in next layer

(propagating)
8's defined recursively working backward
from error at output layer:

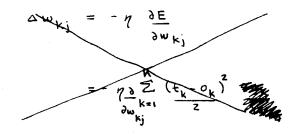
 $\delta_{k} = (t_{k} - o_{k}) f'(net_{k})$ for output layer $\delta_{j} = f'(net_{j}) \geq \delta_{k} w_{kj}$ for penultimate (note to last) layer j $\delta_{i} = f'(net_{i}) \geq \delta_{j} w_{ji}$ for second to (a layer i

we usually have only three layers.

Backward Error Propagation = Delta Rule = Gradient Descent

Gradient Descent - Backward Error Propagation (BEP)-Algorithm (zont.) $E = \frac{1}{2} \sum_{k=1}^{K} (t_k - o_k)$

Last Layer Update Rule:



$$\Delta w_{kj} = - \eta \frac{\partial E}{\partial w_{kj}}$$

$$= - \eta \frac{\partial E}{\partial w_{kj}} \frac{1}{2} \sum_{k=1}^{k} (+_k - o_k)^2$$

Bad notation because k in wkj easily confused with k in Z. My solution is to use & for the particular synaptic weight whose update rule we are finding. PDP book uses same k for both and is confusion. Same story for just j.

Now
$$\frac{\partial o_k}{\partial w_{kj}}$$
 is the change in output o_k when weight w_{kj} changes.

 $= + \eta \quad \frac{1}{2} \quad \frac{\sum_{k=1}^{K} \ z(t_k - o_k)}{\frac{\partial o_k}{\partial w_k}}$

But o_k will not change unless w_k is in neuron k.

$$\frac{\partial OK}{\partial w_{kj}} = 0 \quad \text{unless} \quad k = R$$

So we only get one term in the sum:

Now find dok :

Gradient Descent - BEP-algorithm (cont.)

$$\frac{\partial OR}{\partial w_{k_j}} = \frac{\partial}{\partial w_{k_j}} f(net_k)$$

$$= \frac{\partial f(net_R)}{\partial net_R} \frac{\partial net_R}{\partial w_{R_i}}$$

logistic sigmoid f' = f(1-f)

$$\frac{\partial \text{ net } \underline{\lambda}}{\partial w_{k_{j}}} = \frac{\partial}{\partial w_{k_{j}}} \sum_{j=1}^{J} w_{k_{j}} \circ \underline{\lambda}; \qquad \qquad o_{j} \text{ is input to}$$

$$w_{k_{j}} = \sum_{j=1}^{J} w_{k_{j}} \circ \underline{\lambda}; \qquad \qquad w_{k_{j}} = \sum_{j=1}^{J} w_{k_{j}} \circ \underline{\lambda}; \qquad \qquad w_$$

Again we only get the term where j = j:

$$\Delta w_{kj} = \eta(t_k - o_k) f(net_k)[1 - f(net_k)] o_j$$

| III

 δ_k

We define Sp as shown, and we call our weight update rule a "delta rule." We shall find that weight updates for any layer can be written as a delta rule:

Better yet, 8's for previous layers are defined recursively in terms of 8's for later layers, (consequence of chain rule for der

20 Apr 1990 Peil E Cotter Gradient Descent - BEP-algorithm (cont.)

Penultimate Layer Update Rule: (Next-to-last)

$$\Delta w_{ji} = - \eta \frac{\partial E}{\partial w_{ji}}$$

$$= - \eta \frac{\partial}{\partial w_{ji}} \frac{1}{2} \sum_{k=1}^{K} (t_k - o_k)^2$$

Since the output of neuron of goes to every neuron on the last layer, with affects every o_k and we get all the terms in the \sum_{k} :

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

Now we repeatedly apply the chain rule to $\frac{\partial O_K}{\partial W_{ji}}$. The first few steps are the same as for the last layer update, and we can immediately write down:

$$\frac{\partial o_{k}}{\partial w_{ji}} = f(net_{k}) \left[1 - f(net_{k}) \right] \frac{\partial net_{k}}{\partial w_{ji}}$$

$$\frac{\partial net_{k}}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \sum_{j=1}^{N} w_{kj} \circ_{j}$$

Since with affects only of we get only the of term:

This is similar to dop / dwg, and we get:

Apr 1990 Gradient Descent - BEP-algorithm (ant.)
New E Cotton

$$\frac{\partial o_{j}}{\partial w_{j}i} = f(net_{j})[1-f(net_{j})] o_{i}$$

Putting it all together:

$$\Delta w_{jik} = \eta \sum_{k=1}^{K} (t_k - o_k) f(net_k) [1 - f(net_k)] w_{kj} f(net_j) [1 - f(net_j)]$$

$$\delta_k$$
III

Notice that δ_k is embedded in the update rule. We can write:

$$\Delta w_{ji} = \eta \sum_{k=1}^{K} \delta_{k} w_{kj} f(net_{j})[1-f(net_{j})] o_{i}$$

$$= \eta \delta_{j} o_{i}$$

So we have a delta rule once again. The δ ; is defined in terms of δ_k and synaptic weights from neuron j to neuron k: $\delta_j = f(\text{net}_j) \left[1 - f(\text{net}_j) \right] \sum_{k=1}^{\infty} \delta_k w_{kj}$

(We can take f(1-f) out of sum since it does not depend on k.)

Apr 1990 Gradient Descent - BEP-algorithm (cont.)
Neil E Cotten

Previous Layer Update Rule:

We discover by calculation that the update rule for any layer can be written in terms of 8's for the layer after it.

The form of result for Dwji carries over to all previous layers. In particular, for Dwil we obtain:

DWil = n Siih

We use if rather than of because we are on the first layer of our network, and the inputs to our neuron do not actually come from a previous layer. (Just notation)

where $S_i = f(\text{net}_i) \left[1 - f(\text{net}_i)\right] \sum_{j=1}^{\infty} S_j w_{ji}$

Now we see where the name backward-error-propagation comes from: we start with the output error, t_k -o_k, and propagate it back through the network via the recursive definition of the 8.5.

The computation of the update rule is very similar in form to the neural network computation. Hence, the complexity of weight updates is roughly the same as the complexity of the network computation. Learning slows us down by about half.