

Lecture Series on
Intelligent Control

Lecture 7
Artificial Neural Networks
Backpropagation

Kwang Y. Lee
 Professor of Electrical & Computer Engineering
 Baylor University
 Waco, TX 76706, USA
 Kwang_Y_Lee@baylor.edu

1

Neural Network Terminology

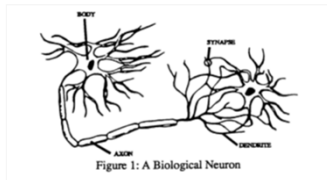


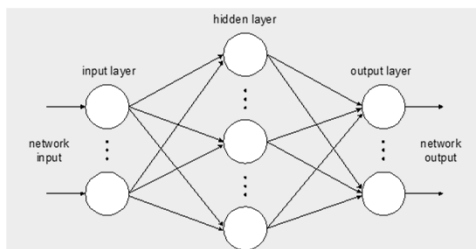
Figure 1: A Biological Neuron

A Biological Neuron Model:

- Each neuron in the brain is composed of *body*, *axon*, and a multitude of *dendrites*.
- Axon is a long tube which splits into branches terminating in endbulbs, almost touching the dendrites of other cells.
- The small gap between an endbulb and a dendrite is called *synapse*.

2

Neural Network Terminology

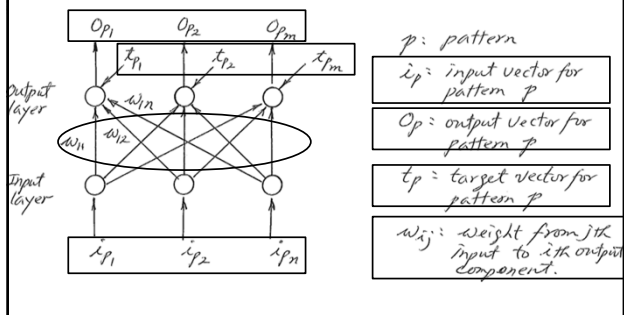


A Feedforward Multi-layer Network: each circle corresponds to a node and each arrow represents a weighted link.

3

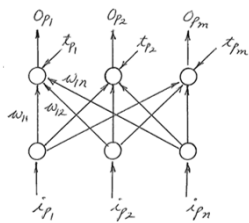
Backpropagation

Delta Rule



4

Delta Rule:



For a pattern p , the error is
 $E_p = \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2$, $E = \sum_p E_p$

The gradient is, by chain rule,

$$\frac{\partial E_p}{\partial w_{ji}} = \frac{\partial E_p}{\partial o_{pj}} \frac{\partial o_{pj}}{\partial w_{ji}}$$

$$\frac{\partial E_p}{\partial o_{pj}} = -(t_{pj} - o_{pj}) \triangleq -\delta_{pj}$$

Assuming linear units,

$$o_{pj} = \sum_i w_{ji} i_{pi} \rightarrow \frac{\partial o_{pj}}{\partial w_{ji}} = i_{pi}$$

Therefore, the negative gradient is
 $-\frac{\partial E_p}{\partial w_{ji}} = \delta_{pj} i_{pi}$

\therefore The weight w_{ji} can be updated by

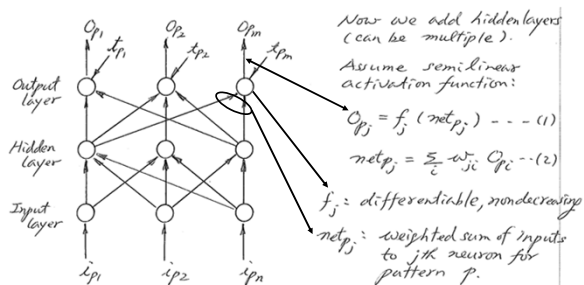
$$\Delta_p w_{ji} = \eta \delta_{pj} i_{pi}, \quad \eta: \text{learning rate}$$

$$\text{Since } \frac{\partial E}{\partial w_{ji}} = \sum_p \frac{\partial E_p}{\partial w_{ji}},$$

the Delta rule implements a gradient descent on E .

5

Generalized Delta Rule:



Note: inputs to j th neuron are outputs of i th neuron in the previous layer.

6

Generalized Delta Rule:

$$O_{pj} = f_j(\text{net}_{pj}) \dots (1)$$

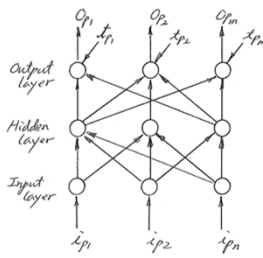
$$\text{net}_{pj} = \sum_i w_{ji} O_{pi} \dots (2)$$

Define the error function:

$$\epsilon_p = \frac{1}{2} (t_{pj} - O_{pj})^2, \quad \epsilon = \sum_p \epsilon_p$$

Then the gradient is *the negative gradient is*

$$\frac{\partial \epsilon_p}{\partial w_{ji}} = \underbrace{\frac{\partial \epsilon_p}{\partial \text{net}_{pj}}}_{O_{pi} \text{ from (2)}} \frac{\partial \text{net}_{pj}}{\partial w_{ji}} \Rightarrow - \frac{\partial \epsilon_p}{\partial w_{ji}} = \delta_{pj} O_{pi} \dots (3)$$



7

Generalized Delta Rule:

$$\epsilon_p = \frac{1}{2} (t_{pj} - O_{pj})^2, \quad \epsilon = \sum_p \epsilon_p$$

Then the gradient is *the negative gradient is*

$$\frac{\partial \epsilon_p}{\partial w_{ji}} = \underbrace{\frac{\partial \epsilon_p}{\partial \text{net}_{pj}}}_{O_{pi} \text{ from (2)}} \frac{\partial \text{net}_{pj}}{\partial w_{ji}} \Rightarrow - \frac{\partial \epsilon_p}{\partial w_{ji}} = \delta_{pj} O_{pi} \dots (3)$$

Define $\delta_{pj} \triangleq - \frac{\partial \epsilon_p}{\partial \text{net}_{pj}}$ $O_{pj} = f_j(\text{net}_{pj}) \dots (1)$
 $\text{net}_{pj} = \sum_i w_{ji} O_{pi} \dots (2)$

$$= - \frac{\partial \epsilon_p}{\partial O_{pj}} \underbrace{\frac{\partial O_{pj}}{\partial \text{net}_{pj}}}_{f'_j(\text{net}_{pj}) \text{ from (1)}} \dots (4)$$

Now we have two different cases to compute the first term: $\frac{\partial \epsilon_p}{\partial O_{pj}}$

8

Generalized Delta Rule:

Case I. Output layer

$$\frac{\partial \epsilon_p}{\partial O_{pj}} = -(t_{pj} - O_{pj}) \quad \epsilon_p = \frac{1}{2} (t_{pj} - O_{pj})^2,$$

Thus, from (4), $\delta_{pj} = (t_{pj} - O_{pj}) f'_j(\text{net}_{pj}) \dots (5)$

$$\delta_{pj} \triangleq - \frac{\partial \epsilon_p}{\partial \text{net}_{pj}} = - \frac{\partial \epsilon_p}{\partial O_{pj}} \underbrace{\frac{\partial O_{pj}}{\partial \text{net}_{pj}}}_{f'_j(\text{net}_{pj}) \text{ from (1)}} \dots (4)$$

9

Generalized Delta Rule:

Class II. Hidden layers

Note that when the output O_j of hidden layer changes the net inputs net_{jk} to the output layer changes.

$$net_{jk} = \sum_i w_{ji} O_i$$

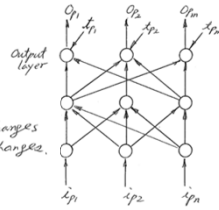
$$\text{Thus, } \frac{\partial \epsilon_p}{\partial O_j} = \sum_k \frac{\partial \epsilon_p}{\partial net_{jk}} \frac{\partial net_{jk}}{\partial O_j}$$

$$= \sum_k \frac{\partial \epsilon_p}{\partial net_{jk}} w_{kj}$$

- δ_{pk} : already calculated at the output layer.

$$= - \sum_k \delta_{pk} w_{kj}$$

Thus, from (4), $\delta_{pj} = f'_j(net_{pj}) \sum_k \delta_{pk} w_{kj}$: error is propagated backward (6)



10

Generalized Delta Rule:

The negative gradient is

$$-\frac{\partial \epsilon_p}{\partial w_{ji}} = \delta_{pj} O_i \dots (3)$$

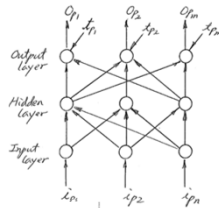
Thus, from (4), $\delta_{pj} = f'_j(net_{pj}) \sum_k \delta_{pk} w_{kj}$: error is propagated backward (6)

The weight adjustment is, therefore, in the direction of the negative gradient:

$$\Delta w_{ji} = \eta \delta_{pj} O_i \dots (7)$$

NOTE:

δ_{pj} is propagated backward (to neuron j) from the next layer.
 O_i is propagated forward (from neuron i in the previous layer)



11

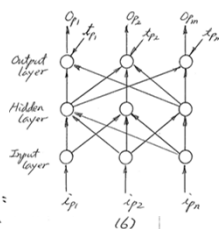
Generalized Delta Rule:

$$\text{Thus, from (4), } \delta_{pj} = f'_j(net_{pj}) \sum_k \delta_{pk} w_{kj} \dots (6)$$

The weight adjustment is, therefore, in the direction of the negative gradient:

$$\Delta w_{ji} = \eta \delta_{pj} O_i \dots (7)$$

Note that the errors in the output layer (5) are propagated backward through output weights (weights between output layer and hidden layer) (6) and multiplied by inputs propagated forward (7) to compute the weight adjustment.



12

Generalized Delta Rule:

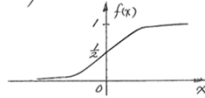
Thus, from (4), $\delta_{pj} = f'_j(\text{net}_{pj}) \sum_k \delta_{pk} w_{kj}$: error is propagated backward (6)

Derivative of the activation function, $f'_j(\text{net}_{pj})$, can be computed directly without differentiating numerically. For example, for a sigmoidal function

$$f(x) = \frac{1}{1 + e^{-(x+\tau)}}$$

it can be shown that

$$f'(x) = f(x)[1 - f(x)]$$



Activation function ($\tau=0$)

Therefore, from (1),

$$O_{pj} = f_j(\text{net}_{pj}) \quad \dots (1)$$

$$f'_j(\text{net}_{pj}) = O_{pj}[1 - O_{pj}]$$

$$\text{net}_{pj} = \sum_i w_{ji} O_{pi} \quad \dots (2)$$

13

Generalized Delta Rule:

Thus, from (4), $\delta_{pj} = f'_j(\text{net}_{pj}) \sum_k \delta_{pk} w_{kj}$: error is propagated backward (6)

The weight adjustment is, therefore, in the direction of the negative gradient:

$$\Delta_p w_{ji} = \eta \delta_{pj} O_{pi} \quad \dots (7)$$

In general, the momentum is added to avoid local minima:

$$\Delta_p w_{ji}^{\text{new}} = \eta \delta_{pj} O_{pi} + \alpha \Delta_p w_{ji}^{\text{prev}}$$

where the second term on the right hand side is the momentum term. This term adds a portion of the most recent weight change when computing the new weight change. The momentum term is supposed to give the neuron momentum in weight space, enabling it to pass through local minima.

14