

# **Intelligens Számítási Módszerek** **Mesterséges neurális hálózatok,** **felügyelt tanulás**

**2005/2006. tanév, II. félév**

**Dr. Kovács Szilveszter**

**E-mail: [szkovacs@iit.uni-miskolc.hu](mailto:szkovacs@iit.uni-miskolc.hu)**

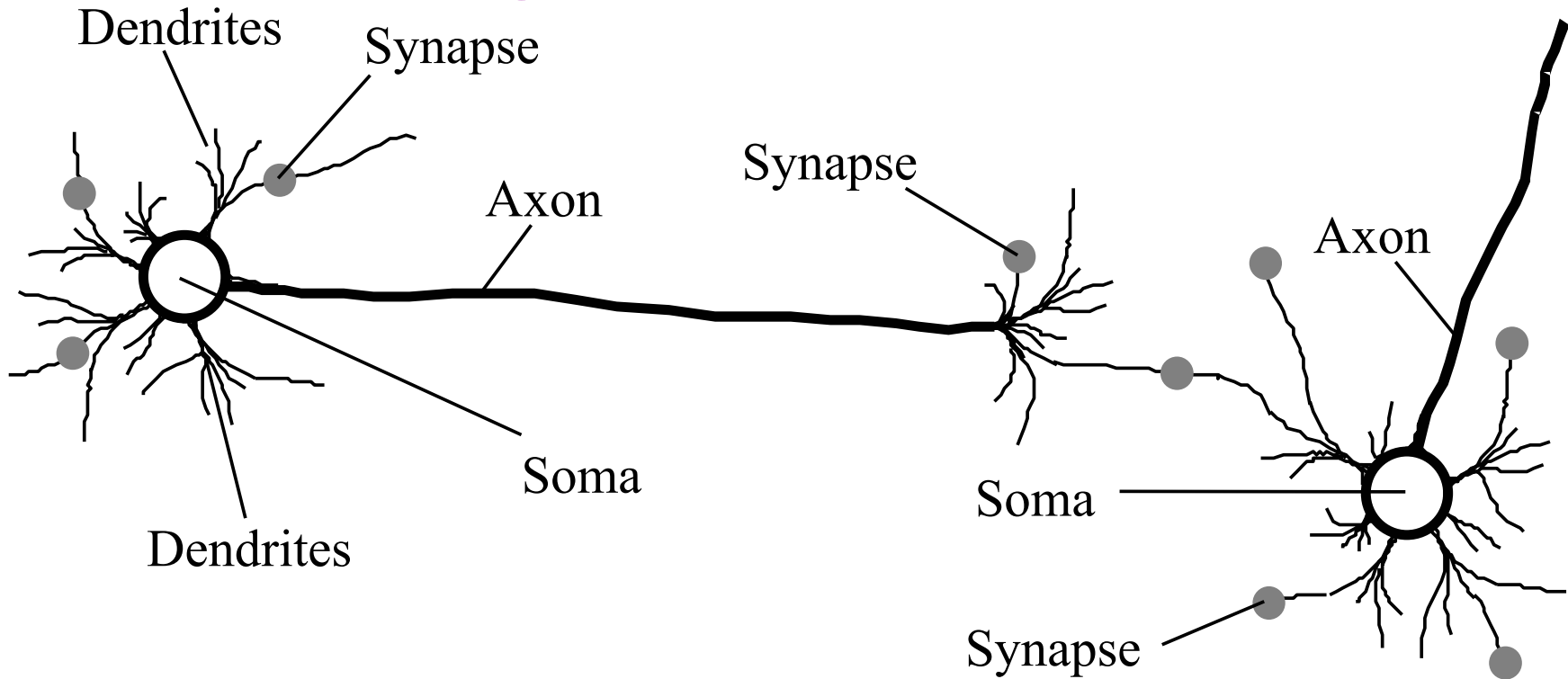
**Informatikai Intézet 106. sz. szoba**

**Tel: (46) 565-111 / 21-06**

# Introduction

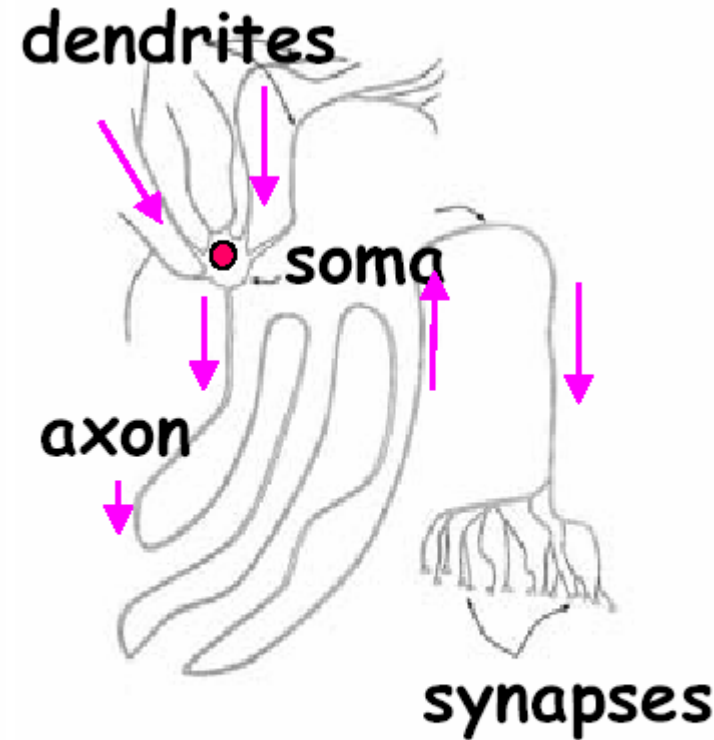
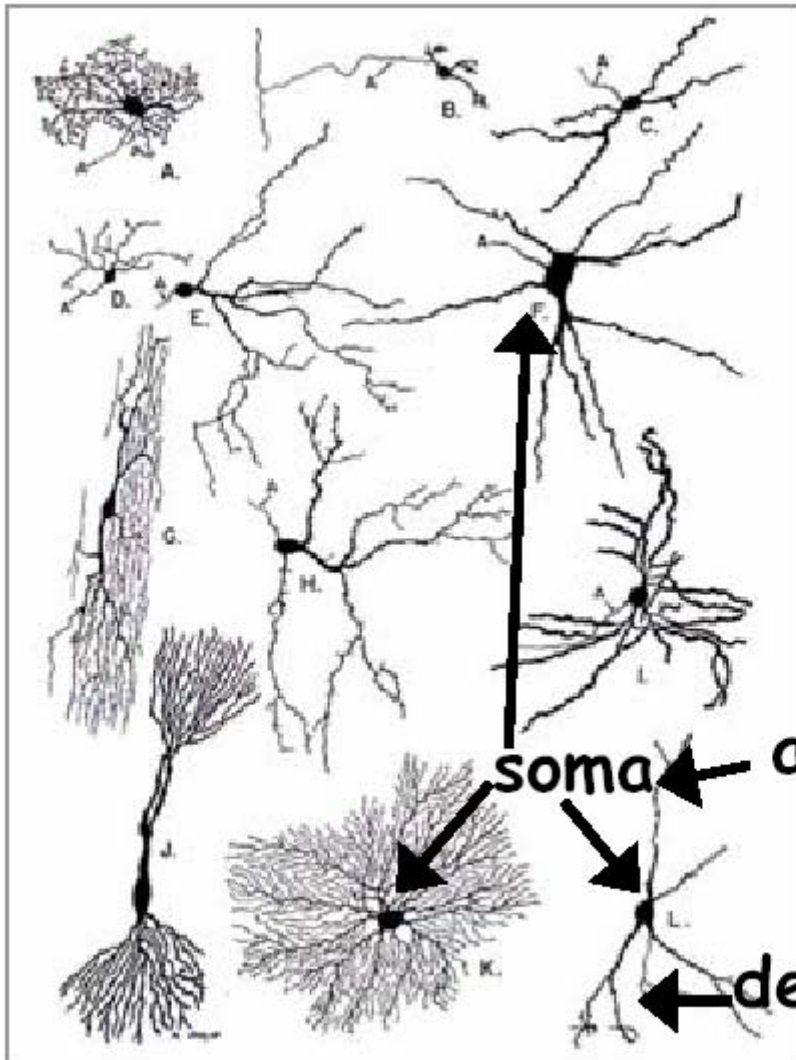
- A **neural network** can be defined as a model of reasoning based on the human brain. The brain consists of a densely interconnected set of nerve cells, or basic information-processing units, called **neurons**.
- The human brain incorporates nearly 10 billion neurons and 60 trillion connections, *synapses*, between them.
- By using multiple neurons simultaneously, the brain can perform its functions much faster than the fastest computers in existence today.

# Biological neural network



- A neuron consists of a cell body, **soma**, a number of fibers called **dendrites**, and a single long fiber called the **axon**.

# Biological neural network



(source: *An Introduction to Neural Networks - Anderson*)

(source: *An Introduction to Neural Networks - Anderson*)

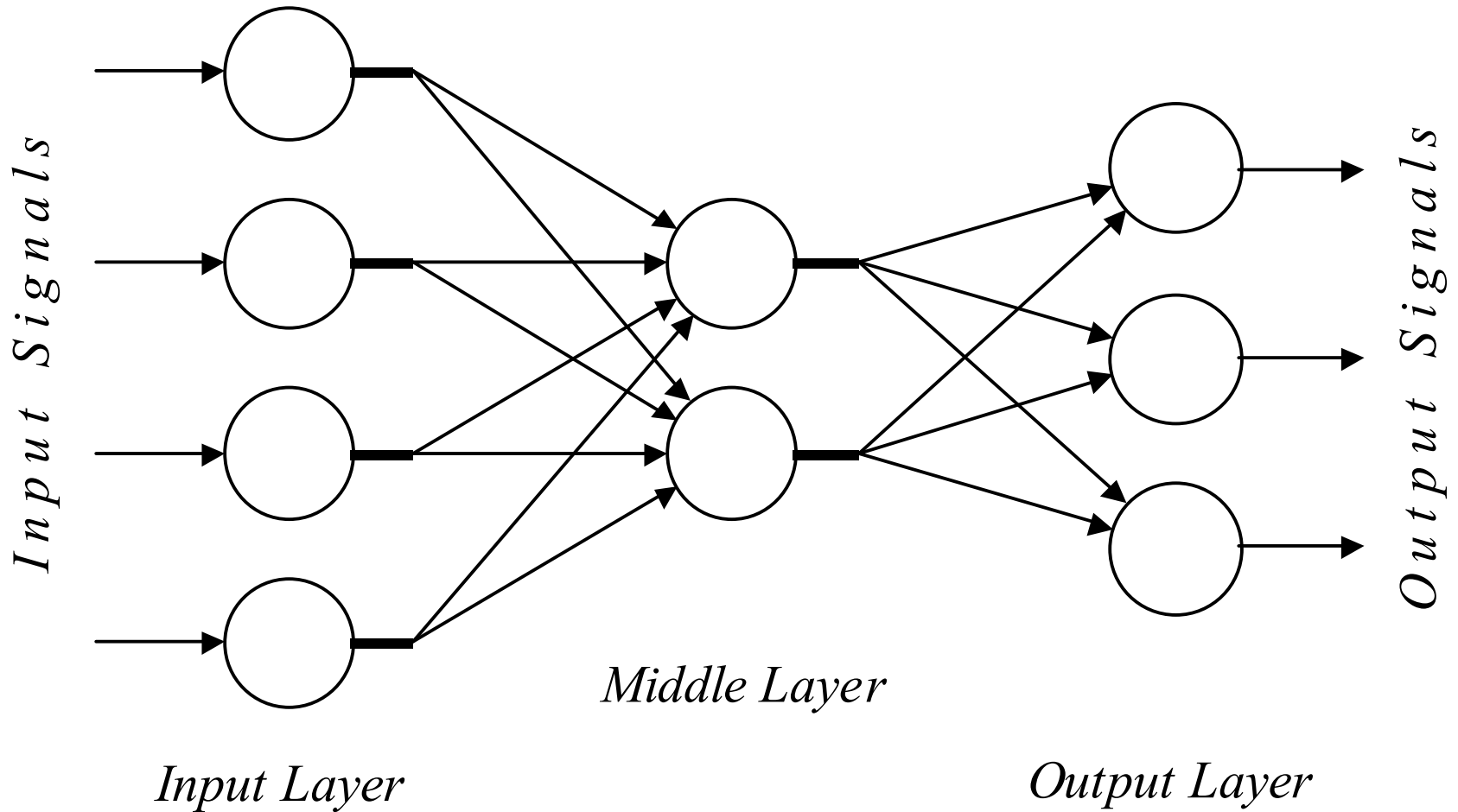
# Introduction

- Our brain can be considered as a highly complex, non-linear and parallel information-processing system.
- Information is stored and processed in a neural network simultaneously throughout the whole network, rather than at specific locations. In other words, in neural networks, both data and its processing are **global** rather than local.
- Learning is a fundamental and essential characteristic of biological neural networks. The ease with which they can learn led to attempts to emulate a biological neural network in a computer.

# Introduction

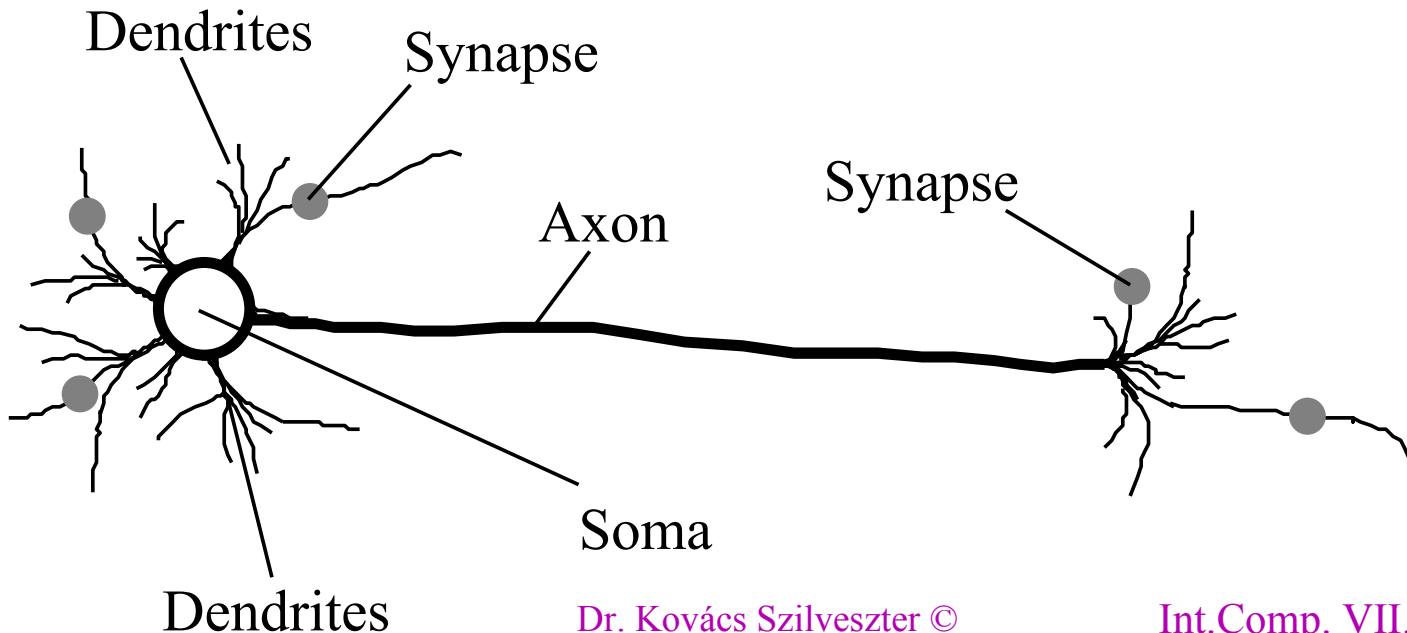
- An artificial neural network consists of a number of very simple processors, also called **neurons**, which are “analogous” to the biological neurons in the brain.
- The neurons are connected by weighted links passing signals from one neuron to another.
- The output signal is transmitted through the neuron’s outgoing connection. The outgoing connection splits into a number of branches that transmit the same signal. The outgoing branches terminate at the incoming connections of other neurons in the network.

# Architecture of a typical artificial neural network



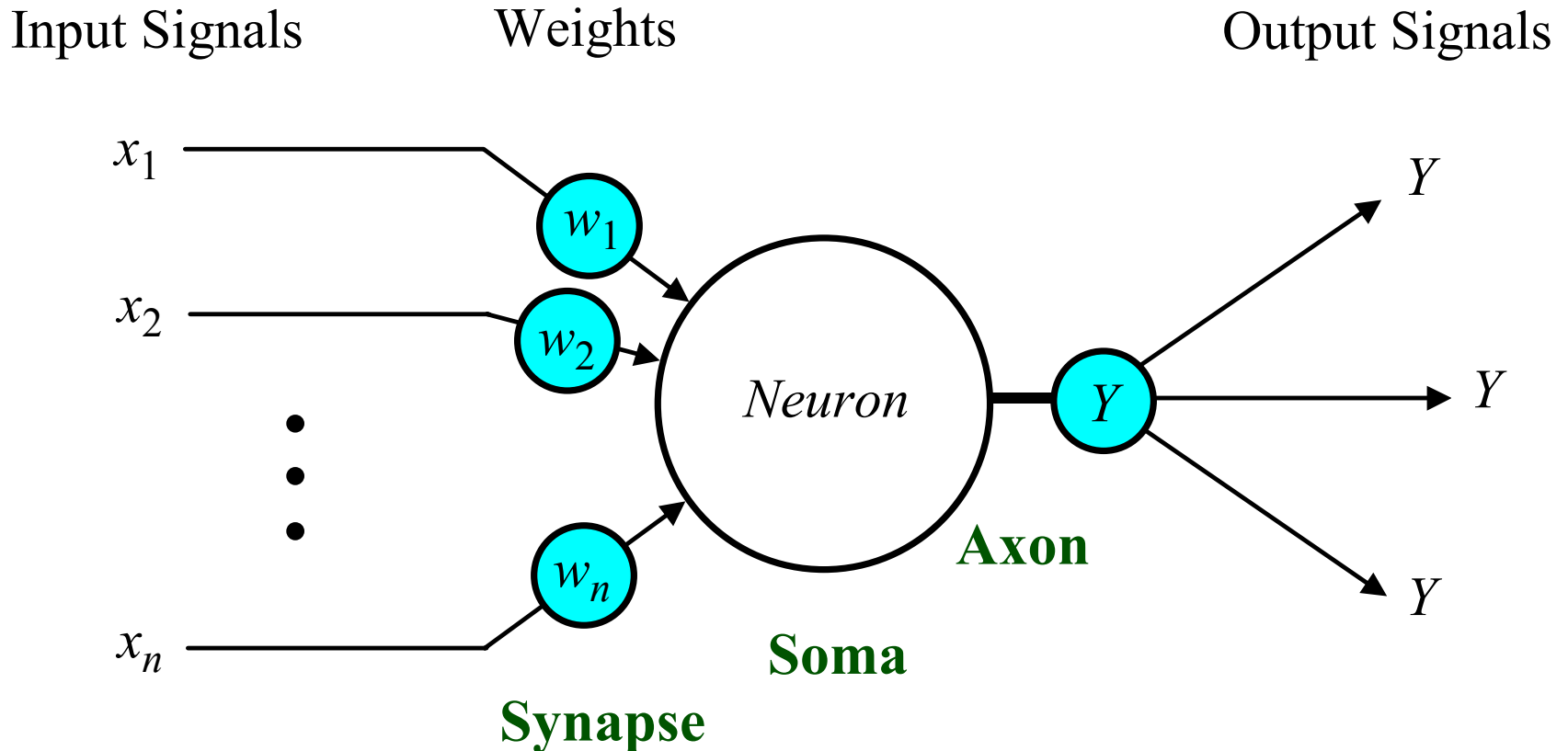
# Analogy between biological and artificial neural networks

<i>Biological Neural Network</i>	<i>Artificial Neural Network</i>
Soma	Neuron
Dendrite	Input
Axon	Output
Synapse	Weight



# The neuron as a simple computing element

## Diagram of a neuron



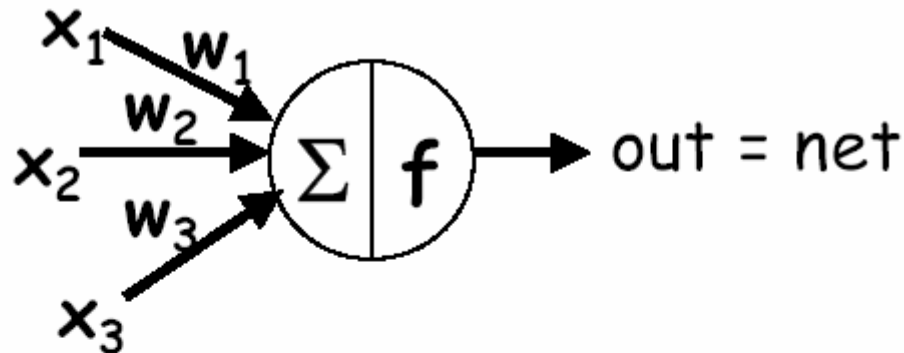
# Artificial neuron

- The neuron computes the weighted sum of the input signals and compares the result with a **threshold value**,  $\theta$ . If the net input is less than the threshold, the neuron output is  $-1$ . But if the net input is greater than or equal to the threshold, the neuron becomes activated and its output attains a value  $+1$ .
- The neuron uses the following transfer or **activation function**:

$$X = \sum_{i=1}^n x_i w_i \quad Y = \begin{cases} +1, & \text{if } X \geq \theta \\ -1, & \text{if } X < \theta \end{cases}$$

- This type of activation function is called a **sign function**.

# Artificial neuron



$\Sigma$  = weighted sum =

$$f[(w_1 \cdot x_1) + (w_2 \cdot x_2) + (w_3 \cdot x_3)]$$

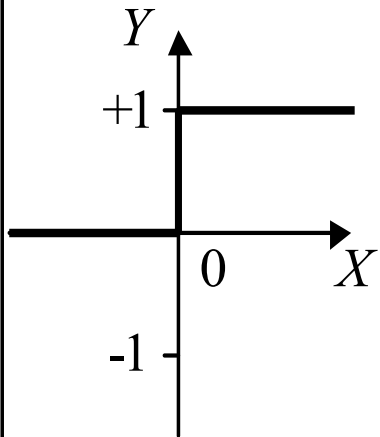
$f$  = activation (transfer) function

$w$  = weights

$x$  = inputs

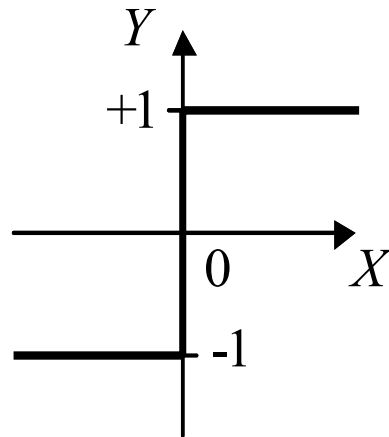
# Activation functions of a neuron

*Step function*



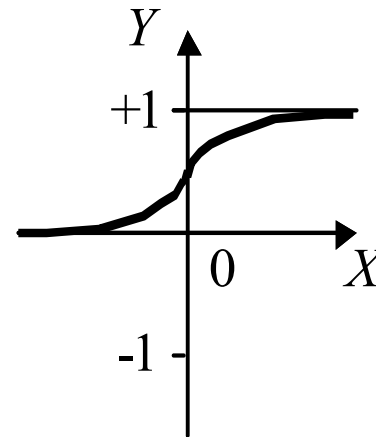
$$y^{step} = \begin{cases} 1, & \text{if } X \geq 0 \\ 0, & \text{if } X < 0 \end{cases}$$

*Sign function*



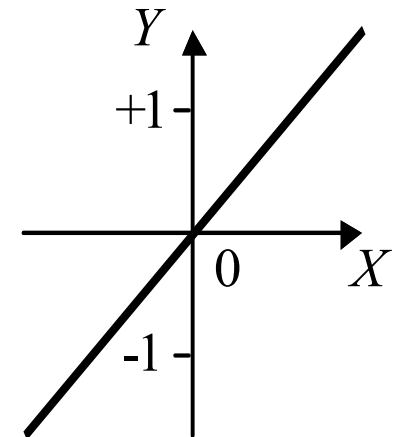
$$y^{sign} = \begin{cases} +1, & \text{if } X \geq 0 \\ -1, & \text{if } X < 0 \end{cases}$$

*Sigmoid function*



$$y^{sigmoid} = \frac{1}{1+e^{-X}}$$

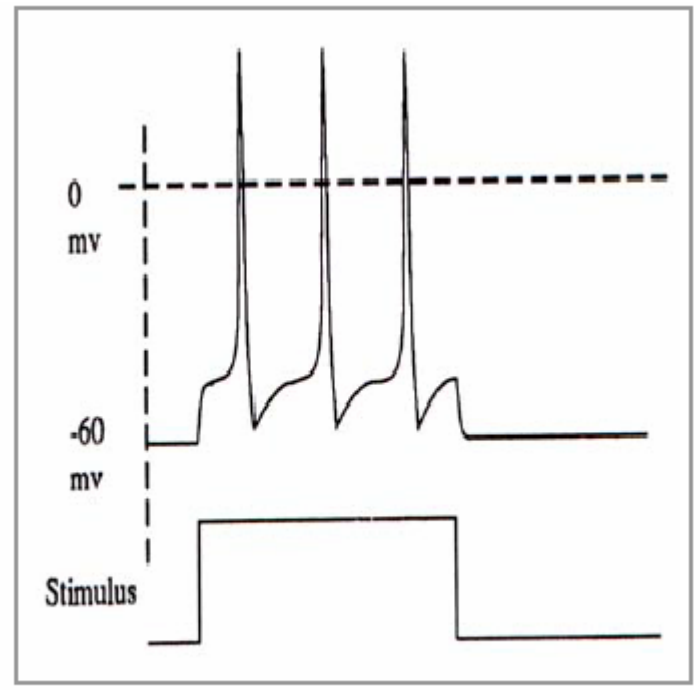
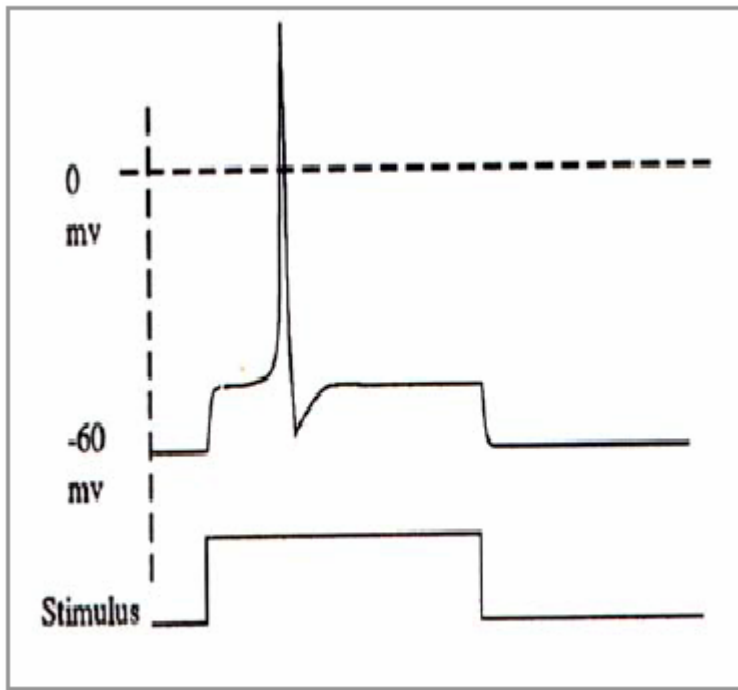
*Linear function*



$$y^{linear} = X$$

# Biological neural network

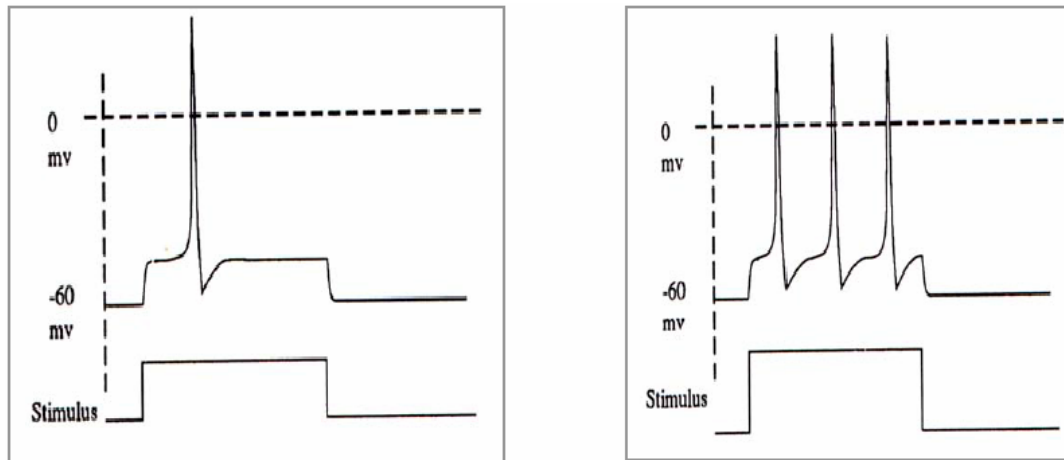
- **Action potential (stimulus) must exceed a threshold for the neuron to fire**
- **Stimulus can be excitatory or inhibitory**



(source: Information Processing & Neural Networks Course Publication - King's College, U of London)

# Biological neural network

- use not even approximately threshold devices
- it is assumed they use a non-linear summation method
- produce a sequence of pulses (not a single output level)
- do not have the same fixed delay ( $t \rightarrow t+1$ )
- are not updated synchronously
- amount of transmitter substance varies unpredictably



(source: Information Processing & Neural Networks Course Publication - King's College, U of London)

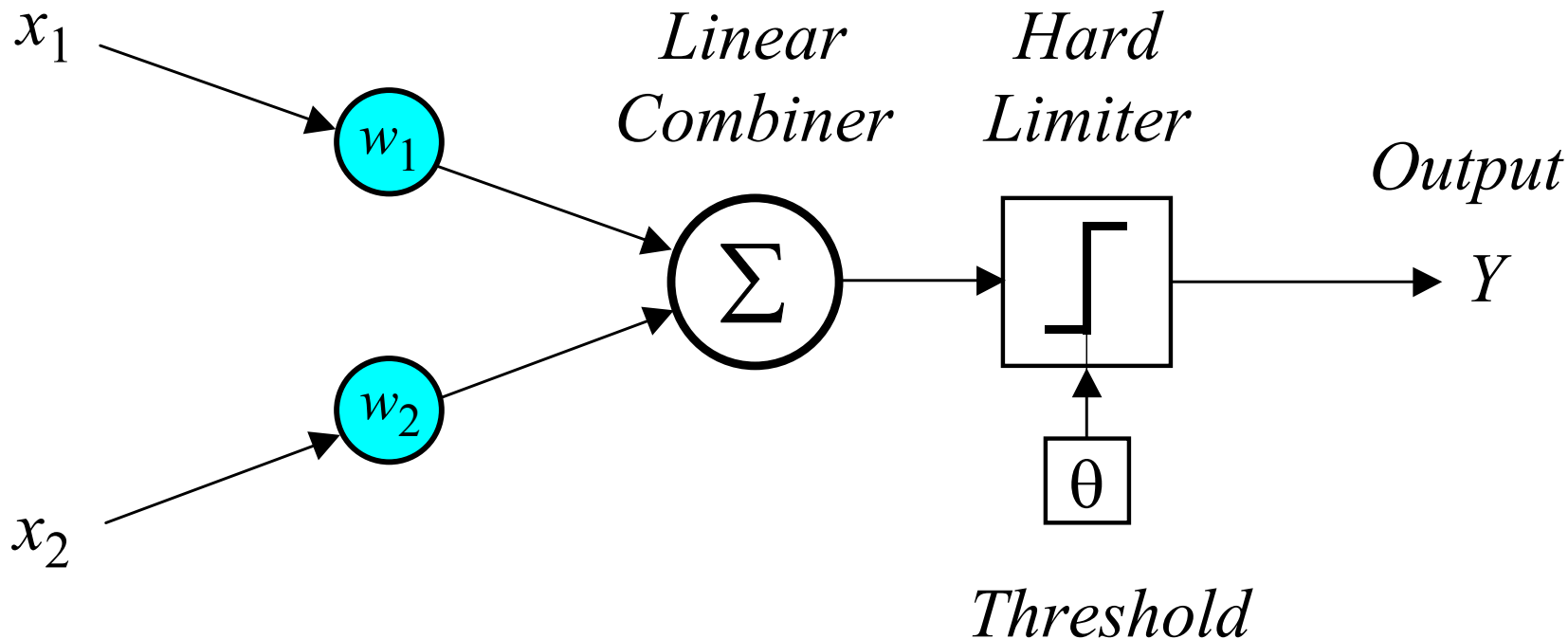
# Can a single neuron learn a task?

- In 1958, **Frank Rosenblatt** introduced a training algorithm that provided the first procedure for training a simple ANN: a **perceptron**.
- The perceptron is the simplest form of a neural network. It consists of a single neuron with *adjustable* synaptic weights and a *hard limiter*.

# Single-layer two-input perceptron

McCulloch and Pitts neuron model (1943)

*Inputs*



# The Perceptron

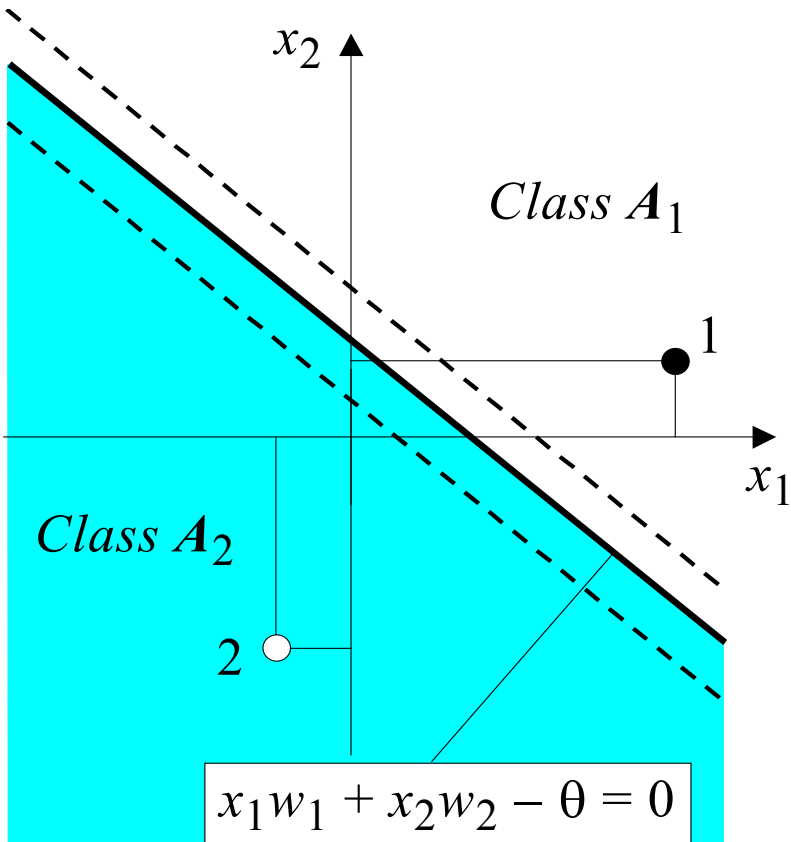
- The operation of Rosenblatt's perceptron is based on the **McCulloch and Pitts neuron model**.  
(McCulloch and Pitts, *A Logical Calculus of the Ideas Immanent in Nervous Activity*, **1943** )
- The model consists of a linear combiner followed by a hard limiter.
- The weighted sum of the inputs is applied to the **hard limiter**, which produces an output equal to **+1** if its input is positive and **-1** if it is negative.

# The Perceptron

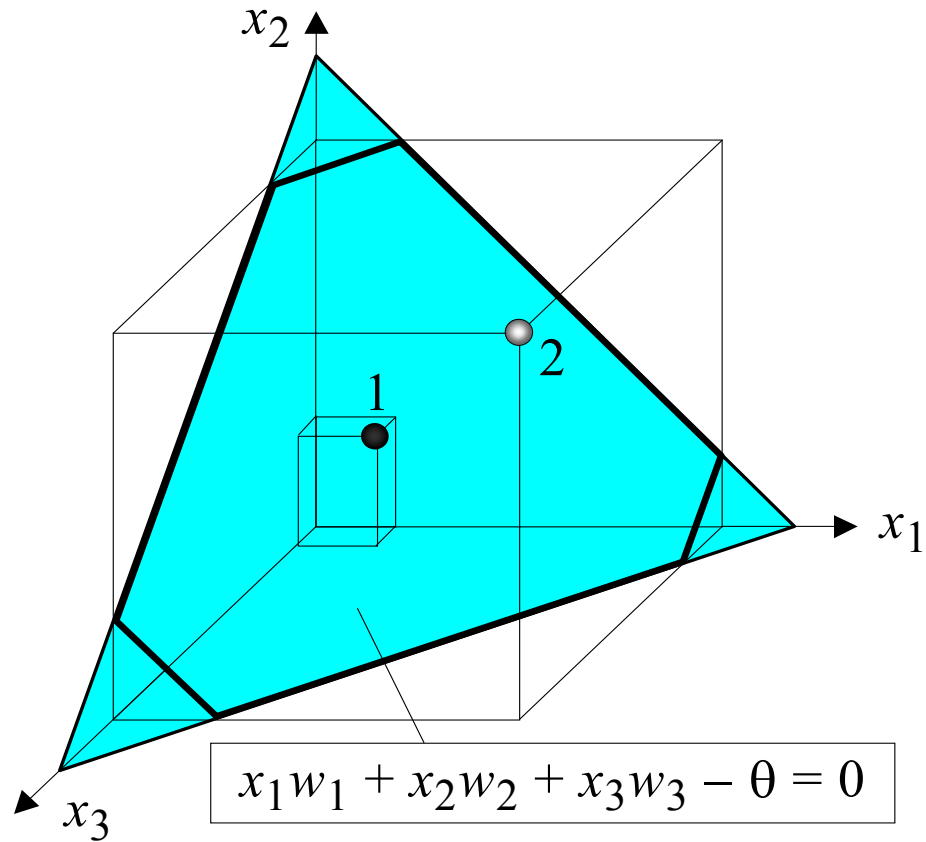
- The aim of the perceptron is to classify inputs,  $x_1, x_2, \dots, x_n$ , into one of two classes, say  $A_1$  and  $A_2$ .
- In the case of an elementary perceptron, the  $n$ -dimensional space is divided by a *hyperplane* into two decision regions. The hyperplane is defined by the *linearly separable function*:

$$\sum_{i=1}^n x_i w_i - \theta = 0$$

# Linear separability in the perceptrons



(a) Two-input perceptron.



(b) Three-input perceptron.

# How does the perceptron learn its classification tasks?

- This is done by making small adjustments in the weights to reduce the difference between the actual and desired outputs of the perceptron.
- The initial weights are randomly assigned, usually in the range  $[-0.5, 0.5]$ , and then updated to obtain the output consistent with the training examples.

# The Perceptron

- If at iteration  $p$ , the actual output is  $Y(p)$  and the desired output is  $Y_d(p)$ , then the error is given by:

$$e(p) = Y_d(p) - Y(p) \quad \text{where } p = 1, 2, 3, \dots$$

- Iteration  $p$  here refers to the  $p$ th training example presented to the perceptron.
- If the error,  $e(p)$ , is positive, we need to increase perceptron output  $Y(p)$ , but if it is negative, we need to decrease  $Y(p)$ .

# The perceptron learning rule

$$w_i(p+1) = w_i(p) + \alpha \cdot x_i(p) \cdot e(p)$$

where  $p = 1, 2, 3, \dots$

$\alpha$  is the **learning rate**, a positive constant less than unity.

The perceptron learning rule was first proposed by **Rosenblatt** in **1960**.

- Using this rule we can derive the perceptron training algorithm for classification tasks.

# Perceptron's training algorithm

## Step 1: Initialisation

Set initial weights  $w_1, w_2, \dots, w_n$  and threshold  $\theta$  to random numbers in the range  $[-0.5, 0.5]$ .

If the error,  $e(p)$ , is positive, we need to increase perceptron output  $Y(p)$ , but if it is negative, we need to decrease  $Y(p)$ .

# Perceptron's training algorithm

## Step 2: Activation

Activate the perceptron by applying inputs  $x_1(p)$ ,  $x_2(p), \dots, x_n(p)$  and desired output  $Y_d(p)$ .

Calculate the actual output at iteration  $p = 1$

$$Y(p) = \text{step} \left[ \sum_{i=1}^n x_i(p) w_i(p) - \theta \right]$$

where  $n$  is the number of the perceptron inputs, and  $\text{step}$  is a step activation function.

$$Y^{\text{step}} = \begin{cases} 1, & \text{if } X \geq 0 \\ 0, & \text{if } X < 0 \end{cases}$$

# Perceptron's training algorithm

## Step 3: Weight training

Update the weights of the perceptron

$$w_i(p+1) = w_i(p) + \Delta w_i(p)$$

where  $\Delta w_i(p)$  is the weight correction at iteration  $p$ .

The weight correction is computed by the **delta rule**:

$$\Delta w_i(p) = \alpha \cdot x_i(p) \cdot e(p)$$

## Step 4: Iteration

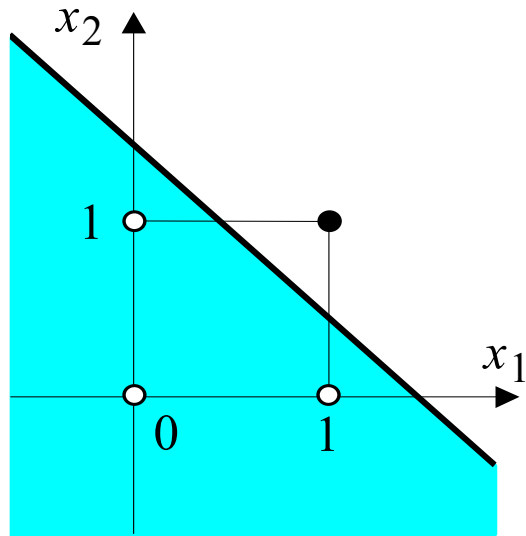
Increase iteration  $p$  by one, go back to *Step 2* and repeat the process until convergence.

# Example of perceptron learning: the logical operation *AND*

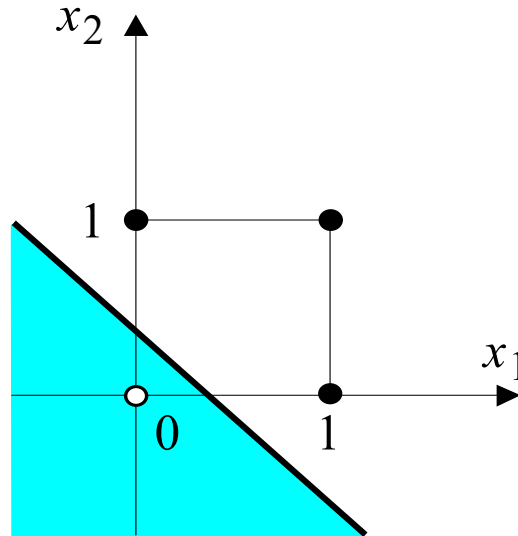
Epoch	Inputs		Desired output $Y_d$	Initial weights		Actual output $Y$	Error $e$	Final weights	
	$x_1$	$x_2$		$w_1$	$w_2$			$w_1$	$w_2$
1	0	0	0	0.3	-0.1	0	0	0.3	-0.1
	0	1	0	0.3	-0.1	0	0	0.3	-0.1
	1	0	0	0.3	-0.1	1	-1	0.2	-0.1
	1	1	1	0.2	-0.1	0	1	0.3	0.0
2	0	0	0	0.3	0.0	0	0	0.3	0.0
	0	1	0	0.3	0.0	0	0	0.3	0.0
	1	0	0	0.3	0.0	1	-1	0.2	0.0
	1	1	1	0.2	0.0	1	0	0.2	0.0
3	0	0	0	0.2	0.0	0	0	0.2	0.0
	0	1	0	0.2	0.0	0	0	0.2	0.0
	1	0	0	0.2	0.0	1	-1	0.1	0.0
	1	1	1	0.1	0.0	0	1	0.2	0.1
4	0	0	0	0.2	0.1	0	0	0.2	0.1
	0	1	0	0.2	0.1	0	0	0.2	0.1
	1	0	0	0.2	0.1	1	-1	0.1	0.1
	1	1	1	0.1	0.1	1	0	0.1	0.1
5	0	0	0	0.1	0.1	0	0	0.1	0.1
	0	1	0	0.1	0.1	0	0	0.1	0.1
	1	0	0	0.1	0.1	0	0	0.1	0.1
	1	1	1	0.1	0.1	1	0	0.1	0.1

Threshold:  $\theta = 0.2$ ; learning rate:  $\alpha = 0.1$

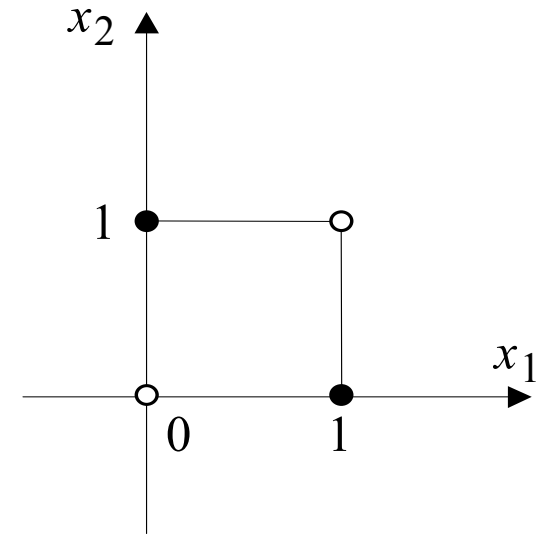
# Two-dimensional plots of basic logical operations



(a) *AND* ( $x_1 \cap x_2$ )



(b) *OR* ( $x_1 \cup x_2$ )



(c) *Exclusive-OR*  
( $x_1 \oplus x_2$ )

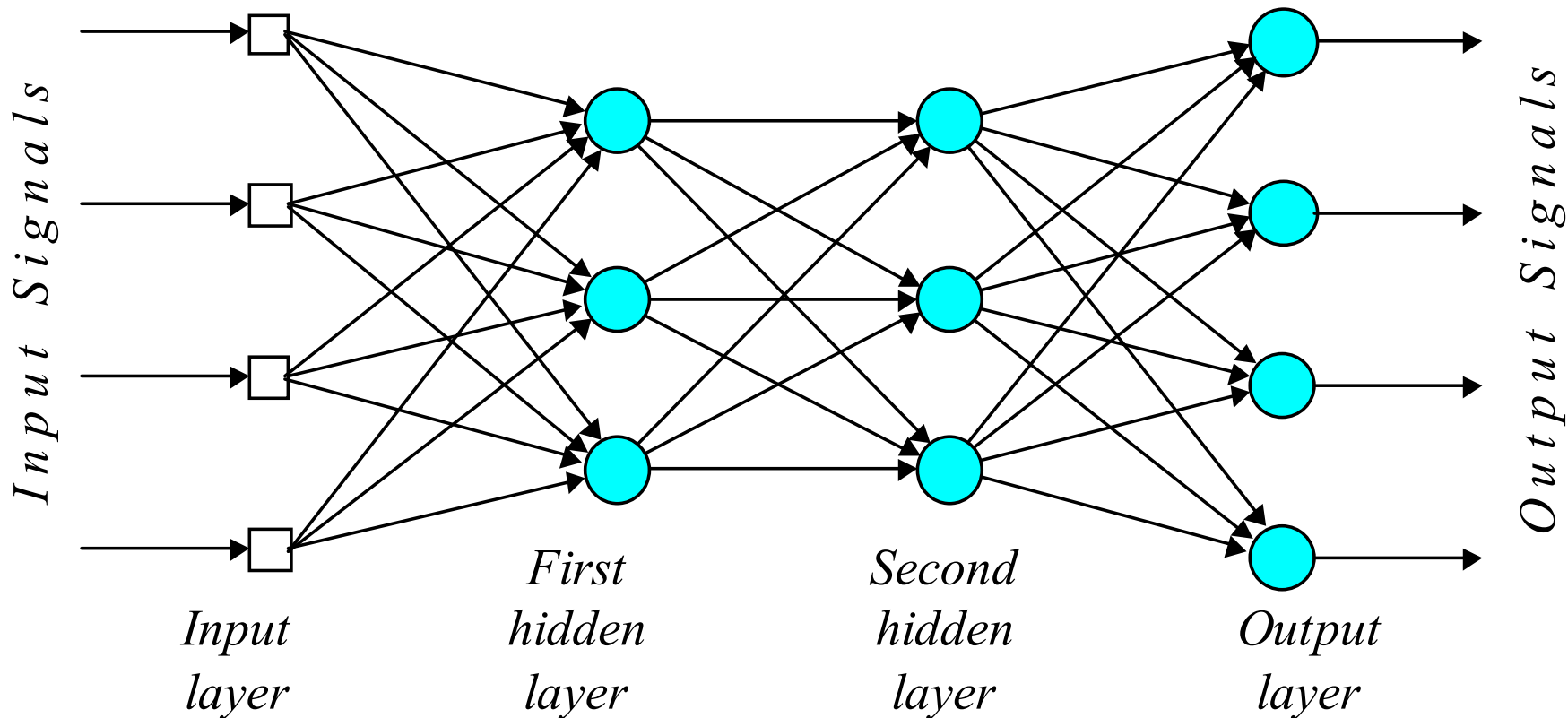
A **perceptron** can learn the operations ***AND*** and ***OR***,  
**but not *Exclusive-OR***.

- ...which (together with a lack of proper training techniques for multi-layer perceptrons) all but killed interest in neural nets in the 70s and early 80s.
- the **Dark Ages of NNs**

# Multilayer neural networks

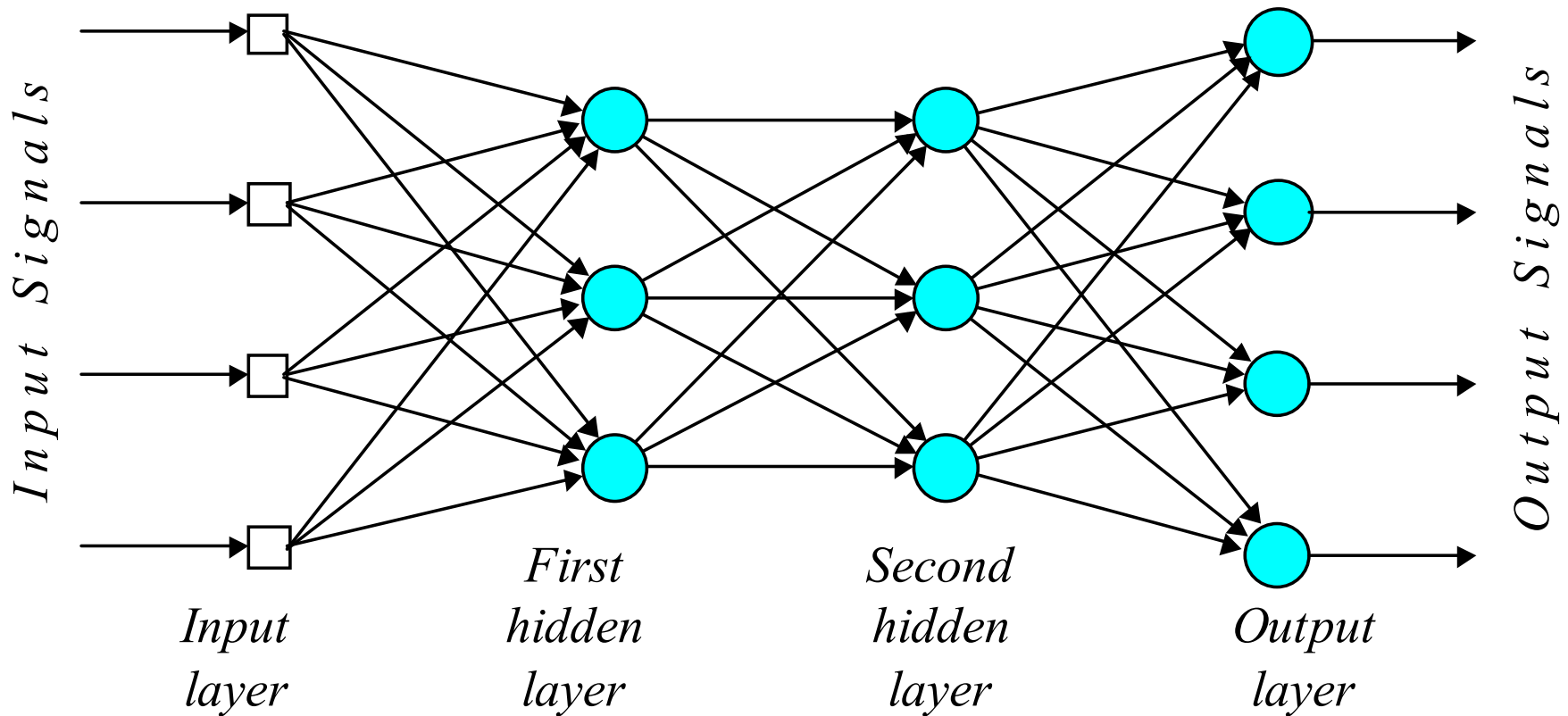
- A multilayer perceptron is a feedforward neural network with one or more hidden layers.
- The network consists of an **input layer** of source neurons, at least one middle or **hidden layer** of computational neurons, and an **output layer** of computational neurons.
- The input signals are propagated in a forward direction on a layer-by-layer basis.

# Multilayer perceptron with two hidden layers



# Supervised Neural Networks

- **Requirement:**  
known input-output relations



# What does the middle layer hide?

- A hidden layer “hides” its desired output. Neurons in the hidden layer cannot be observed through the input/output behaviour of the network. There is no obvious way to know what the desired output of the hidden layer should be.
- Commercial ANNs incorporate three and sometimes four layers, including one or two hidden layers. Each layer can contain from 10 to 1000 neurons. Experimental neural networks may have five or even six layers, including three or four hidden layers, and utilise millions of neurons.

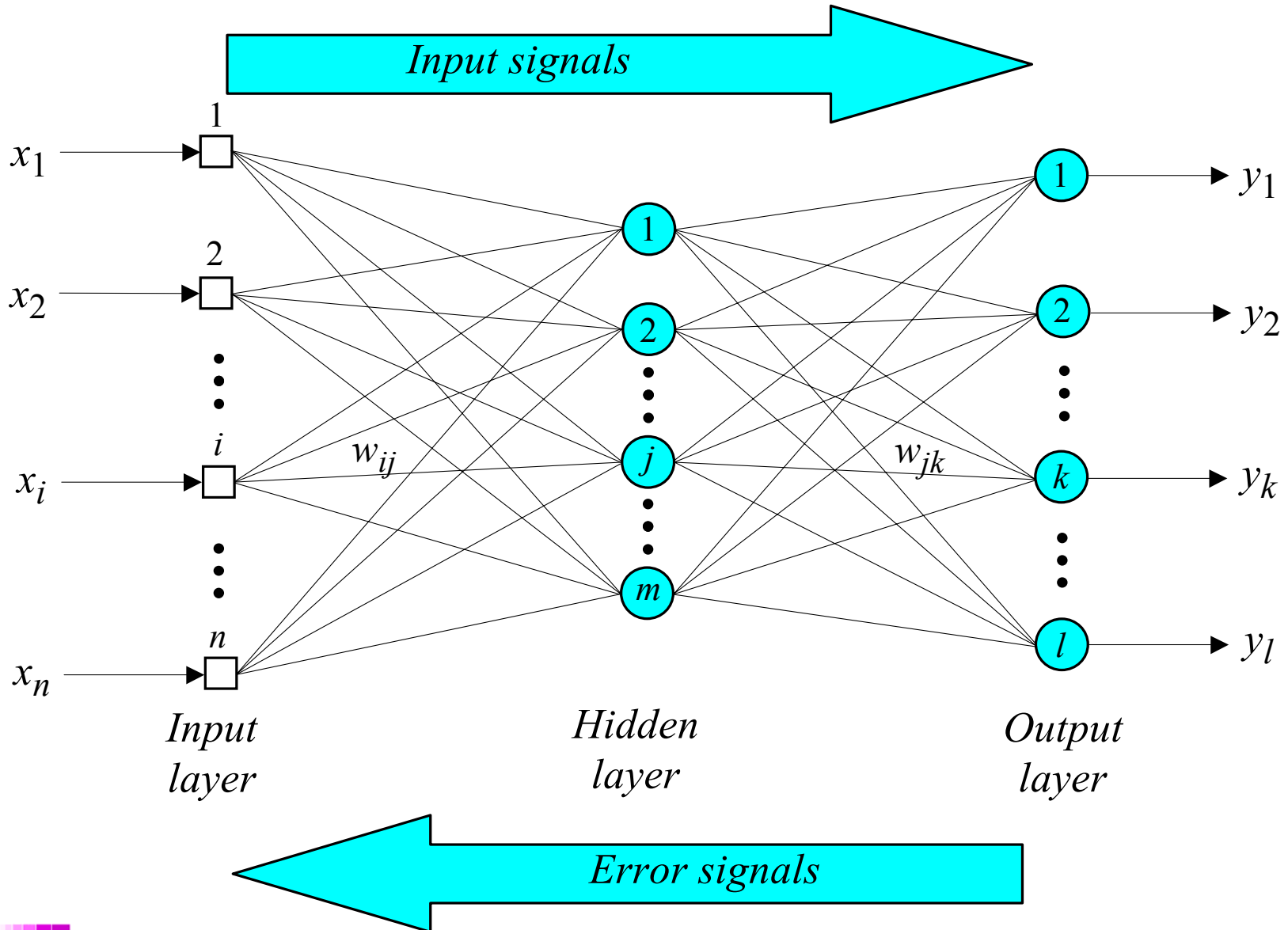
# Back-propagation neural network

- Learning in a multilayer network proceeds the same way as for a perceptron.
- A training set of input patterns is presented to the network.
- The network computes its output pattern, and if there is an error – or in other words a difference between actual and desired output patterns – the weights are adjusted to reduce this error.

# Back-propagation neural network

- In a back-propagation neural network, the learning algorithm has two phases.
- First, a training input pattern is presented to the network input layer. The network propagates the input pattern from layer to layer until the output pattern is generated by the output layer.
- If this pattern is different from the desired output, an error is calculated and then propagated backwards through the network from the output layer to the input layer.
- The weights are modified as the error is propagated.

# Three-layer back-propagation neural network



# The back-propagation training algorithm

## Step 1: Initialisation

Set all the weights and threshold levels of the network to random numbers uniformly distributed inside a small range:

$$\left( -\frac{2.4}{F_i}, +\frac{2.4}{F_i} \right)$$

where  $F_i$  is the total number of inputs of neuron  $i$  in the network. The weight initialisation is done on a neuron-by-neuron basis.

# The back-propagation training algorithm

## Step 2: Activation

Activate the back-propagation neural network by applying inputs  $x_1(p), x_2(p), \dots, x_n(p)$  and desired outputs  $y_{d,1}(p), y_{d,2}(p), \dots, y_{d,n}(p)$ .

(a) Calculate the actual outputs of the neurons in the hidden layer:

$$y_j(p) = \textit{sigmoid} \left[ \sum_{i=1}^n x_i(p) \cdot w_{ij}(p) - \theta_j \right]$$

where  $n$  is the number of inputs of neuron  $j$  in the hidden layer, and *sigmoid* is the *sigmoid* activation function.

$$Y^{\textit{sigmoid}} = \frac{1}{1 + e^{-x}}$$

# The back-propagation training algorithm

## Step 2: Activation

(b) Calculate the actual outputs of the neurons in the output layer:

$$y_k(p) = \textit{sigmoid} \left[ \sum_{j=1}^m x_{jk}(p) \cdot w_{jk}(p) - \theta_k \right]$$

where  $m$  is the number of inputs of neuron  $k$  in the output layer.

$$Y^{\textit{sigmoid}} = \frac{1}{1 + e^{-x}}$$

## Step 3: Weight training

Update the weights in the back-propagation network propagating backward the errors associated with output neurons.

(a) Calculate the **error gradient** for the neurons in the output layer:

$$\delta_k(p) = y_k(p) \cdot [1 - y_k(p)] \cdot e_k(p)$$

where  $e_k(p) = y_{d,k}(p) - y_k(p)$

Calculate the **weight corrections**:

$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p)$$

Update the weights at the output neurons:

$$w_{jk}(p+1) = w_{jk}(p) + \Delta w_{jk}(p)$$

# Error gradient

$$Y^{\text{sigmoid}} = \frac{1}{1 + e^{-X}}$$

$$\delta_k(p) = \frac{\partial y_k(p)}{\partial X_k(p)} \cdot e_k(p)$$

$$y_k(p) = \frac{1}{1 + e^{-X_k(p)}}$$

$$\delta_k(p) = \frac{\partial \left\{ \frac{1}{1 + \exp[-X_k(p)]} \right\}}{\partial X_k(p)} \cdot e_k(p) = \frac{\exp[-X_k(p)]}{\{1 + \exp[-X_k(p)]\}^2} \cdot e_k(p)$$

$$\begin{aligned} \frac{\exp[-X_k(p)]}{\{1 + \exp[-X_k(p)]\}^2} &= \frac{1 + \exp[-X_k(p)]}{\{1 + \exp[-X_k(p)]\}^2} - \frac{1}{\{1 + \exp[-X_k(p)]\}^2} \\ &= \frac{1}{1 + \exp[-X_k(p)]} - \frac{1}{\{1 + \exp[-X_k(p)]\}^2} = y_k(p) - y_k(p)^2 \end{aligned}$$

$$\delta_k(p) = y_k(p) \cdot [1 - y_k(p)] \cdot e_k(p)$$

## Step 3: Weight training

(b) Calculate the error gradient for the neurons in the hidden layer:

$$\delta_j(p) = y_j(p) \cdot [1 - y_j(p)] \cdot \sum_{k=1}^l \delta_k(p) w_{jk}(p)$$

*Gradient* ←      ← *Error*

Calculate the weight corrections:

$$\Delta w_{ij}(p) = \alpha \cdot x_i(p) \cdot \delta_j(p)$$

Update the weights at the hidden neurons:

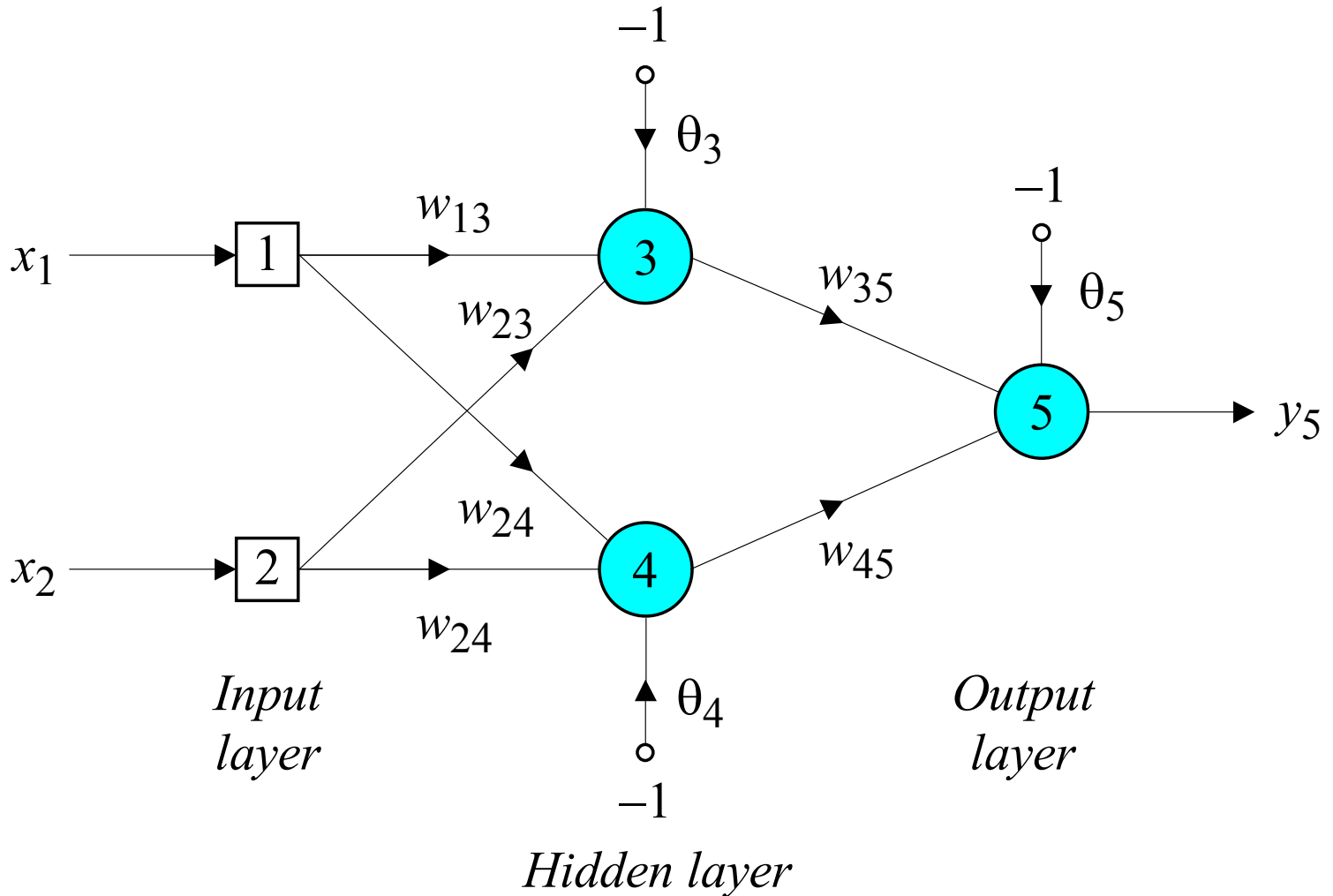
$$w_{ij}(p+1) = w_{ij}(p) + \Delta w_{ij}(p)$$

## Step 4: Iteration

Increase iteration  $p$  by one, go back to *Step 2* and repeat the process until the selected error criterion is satisfied.

**As an example**, we may consider the three-layer back-propagation network. Suppose that the network is required to perform logical operation ***Exclusive-OR***. Recall that a single-layer perceptron could not do this operation. Now we will apply the three-layer net.

# Three-layer network for solving the Exclusive-OR operation



# Exclusive-OR operation

- The effect of the threshold applied to a neuron in the hidden or output layer is represented by its weight,  $\theta$ , connected to a fixed input equal to  $-1$ .
- The initial weights and threshold levels are set randomly as follows:

$$w_{13} = 0.5, w_{14} = 0.9, w_{23} = 0.4, w_{24} = 1.0, w_{35} = -1.2, \\ w_{45} = 1.1, \theta_3 = 0.8, \theta_4 = -0.1 \text{ and } \theta_5 = 0.3.$$

- We consider a training set where inputs  $x_1$  and  $x_2$  are equal to 1 and desired output  $y_{d,5}$  is 0. The actual outputs of neurons 3 and 4 in the hidden layer are calculated as

$$y_3 = \textit{sigmoid}(x_1w_{13} + x_2w_{23} - \theta_3) = 1 / \left[ 1 + e^{-(1 \cdot 0.5 + 1 \cdot 0.4 - 1 \cdot 0.8)} \right] = 0.5250$$

$$y_4 = \textit{sigmoid}(x_1w_{14} + x_2w_{24} - \theta_4) = 1 / \left[ 1 + e^{-(1 \cdot 0.9 + 1 \cdot 1.0 + 1 \cdot 0.1)} \right] = 0.8808$$

- Now the actual output of neuron 5 in the output layer is determined as:

$$y_5 = \textit{sigmoid}(y_3w_{35} + y_4w_{45} - \theta_5) = 1 / \left[ 1 + e^{-(0.5250 \cdot 1.2 + 0.8808 \cdot 1.1 - 1 \cdot 0.3)} \right] = 0.5097$$

- Thus, the following error is obtained:

$$e = y_{d,5} - y_5 = 0 - 0.5097 = -0.5097$$

- The next step is weight training. To update the weights and threshold levels in our network, we propagate the error,  $e$ , from the output layer backward to the input layer.
- First, we calculate the error gradient for neuron 5 in the output layer:

$$\delta_5 = y_5 (1 - y_5) e = 0.5097 \cdot (1 - 0.5097) \cdot (-0.5097) = -0.1274$$

- Then we determine the weight corrections assuming that the learning rate parameter,  $\alpha$ , is equal to 0.1:

$$\Delta w_{35} = \alpha \cdot y_3 \cdot \delta_5 = 0.1 \cdot 0.5250 \cdot (-0.1274) = -0.0067$$

$$\Delta w_{45} = \alpha \cdot y_4 \cdot \delta_5 = 0.1 \cdot 0.8808 \cdot (-0.1274) = -0.0112$$

$$\Delta \theta_5 = \alpha \cdot (-1) \cdot \delta_5 = 0.1 \cdot (-1) \cdot (-0.1274) = -0.0127$$

- Next we calculate the error gradients for neurons 3 and 4 in the hidden layer:

$$\delta_3 = y_3(1 - y_3) \cdot \delta_5 \cdot w_{35} = 0.5250 \cdot (1 - 0.5250) \cdot (-0.1274) \cdot (-1.2) = 0.0381$$

$$\delta_4 = y_4(1 - y_4) \cdot \delta_5 \cdot w_{45} = 0.8808 \cdot (1 - 0.8808) \cdot (-0.1274) \cdot 1.1 = -0.0147$$

- We then determine the weight corrections:

$$\Delta w_{13} = \alpha \cdot x_1 \cdot \delta_3 = 0.1 \cdot 1 \cdot 0.0381 = 0.0038$$

$$\Delta w_{23} = \alpha \cdot x_2 \cdot \delta_3 = 0.1 \cdot 1 \cdot 0.0381 = 0.0038$$

$$\Delta \theta_3 = \alpha \cdot (-1) \cdot \delta_3 = 0.1 \cdot (-1) \cdot 0.0381 = -0.0038$$

$$\Delta w_{14} = \alpha \cdot x_1 \cdot \delta_4 = 0.1 \cdot 1 \cdot (-0.0147) = -0.0015$$

$$\Delta w_{24} = \alpha \cdot x_2 \cdot \delta_4 = 0.1 \cdot 1 \cdot (-0.0147) = -0.0015$$

$$\Delta \theta_4 = \alpha \cdot (-1) \cdot \delta_4 = 0.1 \cdot (-1) \cdot (-0.0147) = 0.0015$$

- At last, we update all weights and threshold:

$$w_{13} = w_{13} + \Delta w_{13} = 0.5 + 0.0038 = 0.5038$$

$$w_{14} = w_{14} + \Delta w_{14} = 0.9 - 0.0015 = 0.8985$$

$$w_{23} = w_{23} + \Delta w_{23} = 0.4 + 0.0038 = 0.4038$$

$$w_{24} = w_{24} + \Delta w_{24} = 1.0 - 0.0015 = 0.9985$$

$$w_{35} = w_{35} + \Delta w_{35} = -1.2 - 0.0067 = -1.2067$$

$$w_{45} = w_{45} + \Delta w_{45} = 1.1 - 0.0112 = 1.0888$$

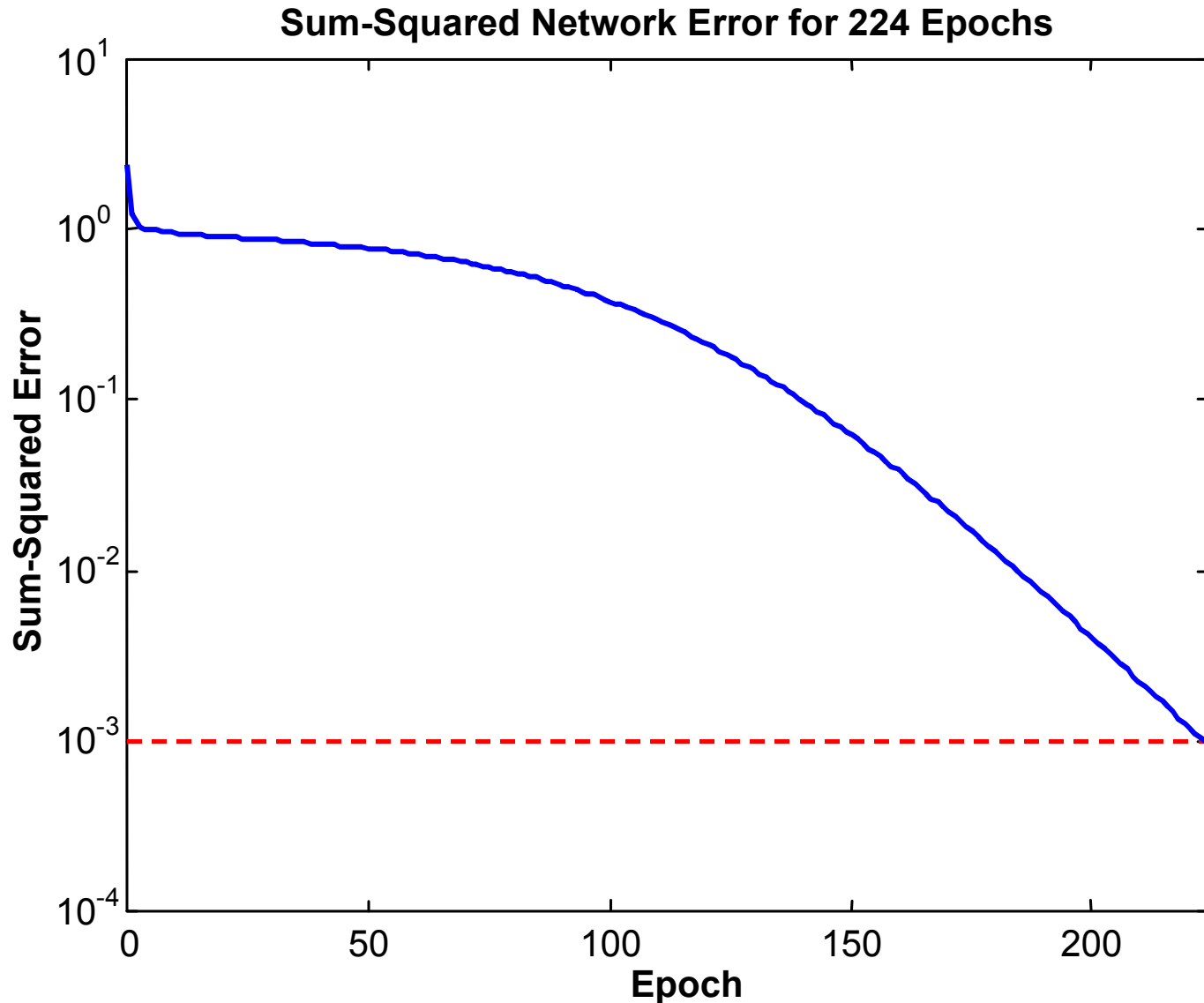
$$\theta_3 = \theta_3 + \Delta\theta_3 = 0.8 - 0.0038 = 0.7962$$

$$\theta_4 = \theta_4 + \Delta\theta_4 = -0.1 + 0.0015 = -0.0985$$

$$\theta_5 = \theta_5 + \Delta\theta_5 = 0.3 + 0.0127 = 0.3127$$

- The training process is repeated until the sum of squared errors is less than 0.001.

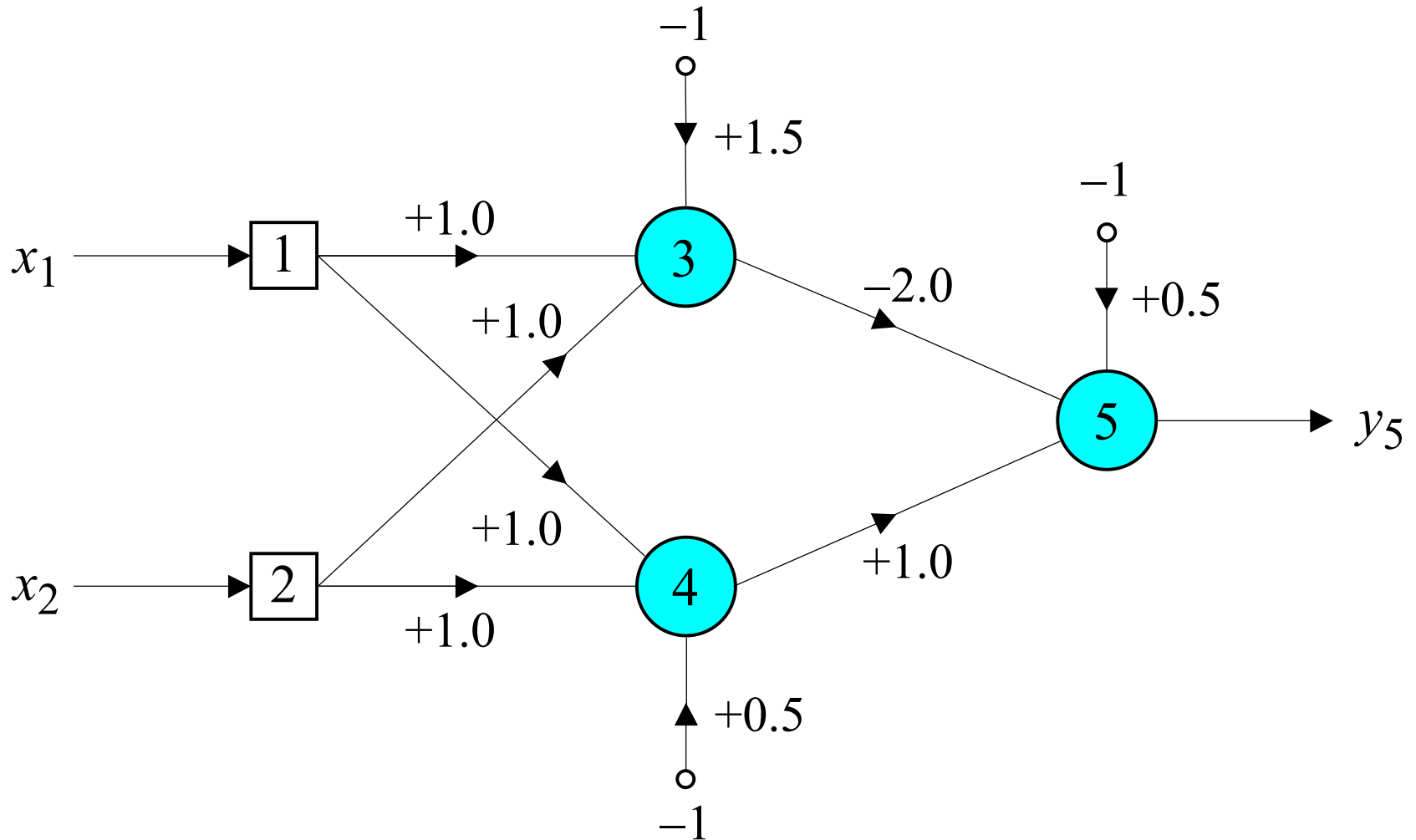
# Learning curve for operation Exclusive-OR



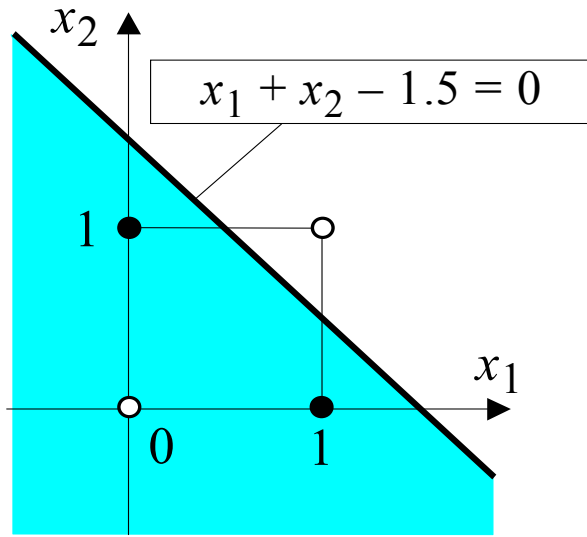
# Final results of three-layer network learning

Inputs		Desired output $y_d$	Actual output $y_5$	Error $e$	Sum of squared errors
$x_1$	$x_2$				
1	1	0	0.0155	-0.0155	0.0010
0	1	1	0.9849	0.0151	
1	0	1	0.9849	0.0151	
0	0	0	0.0175	-0.0175	

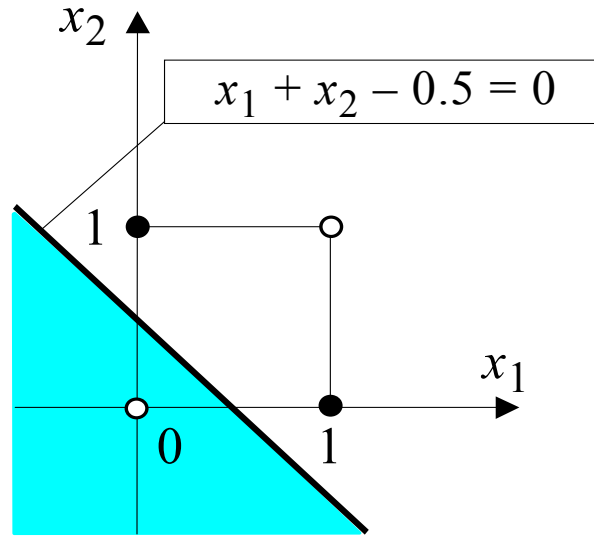
# Network represented by McCulloch-Pitts model for solving the *Exclusive-OR* operation



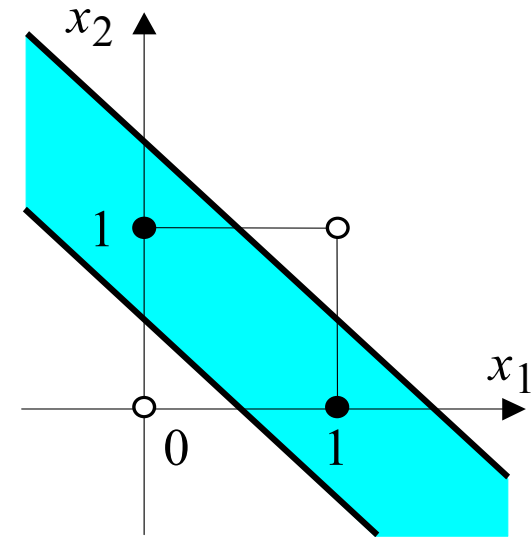
# Decision boundaries



(a)



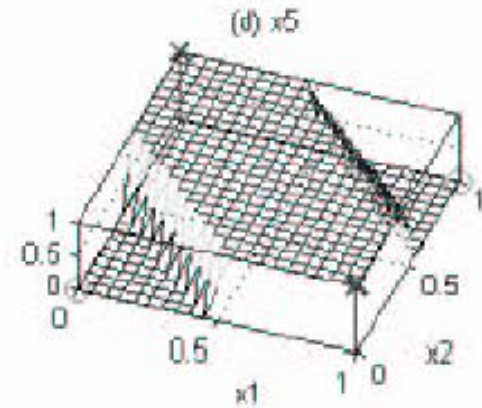
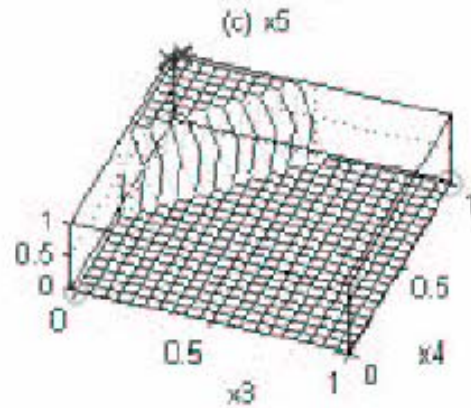
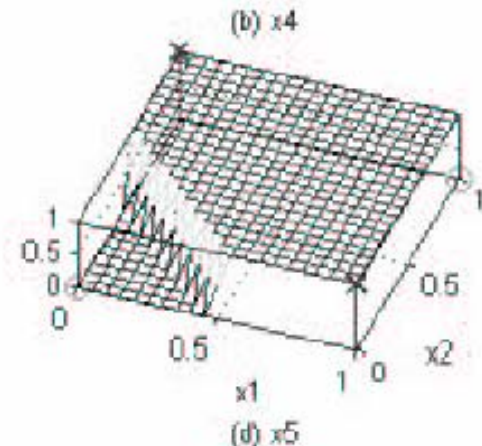
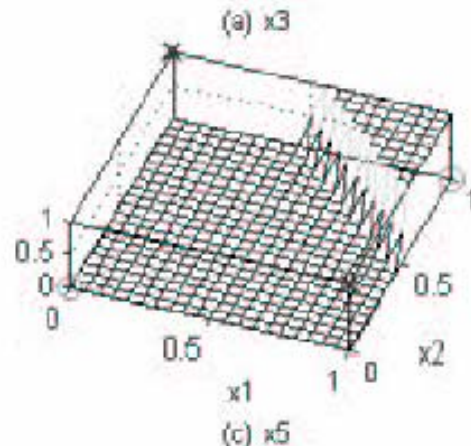
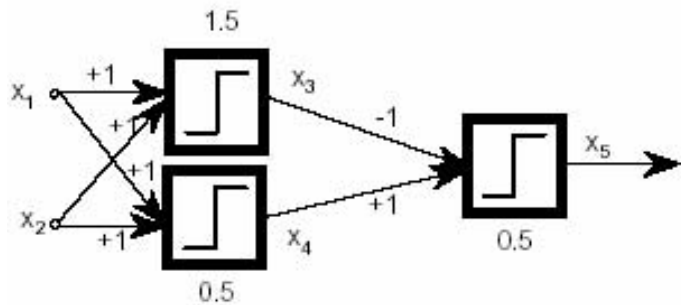
(b)



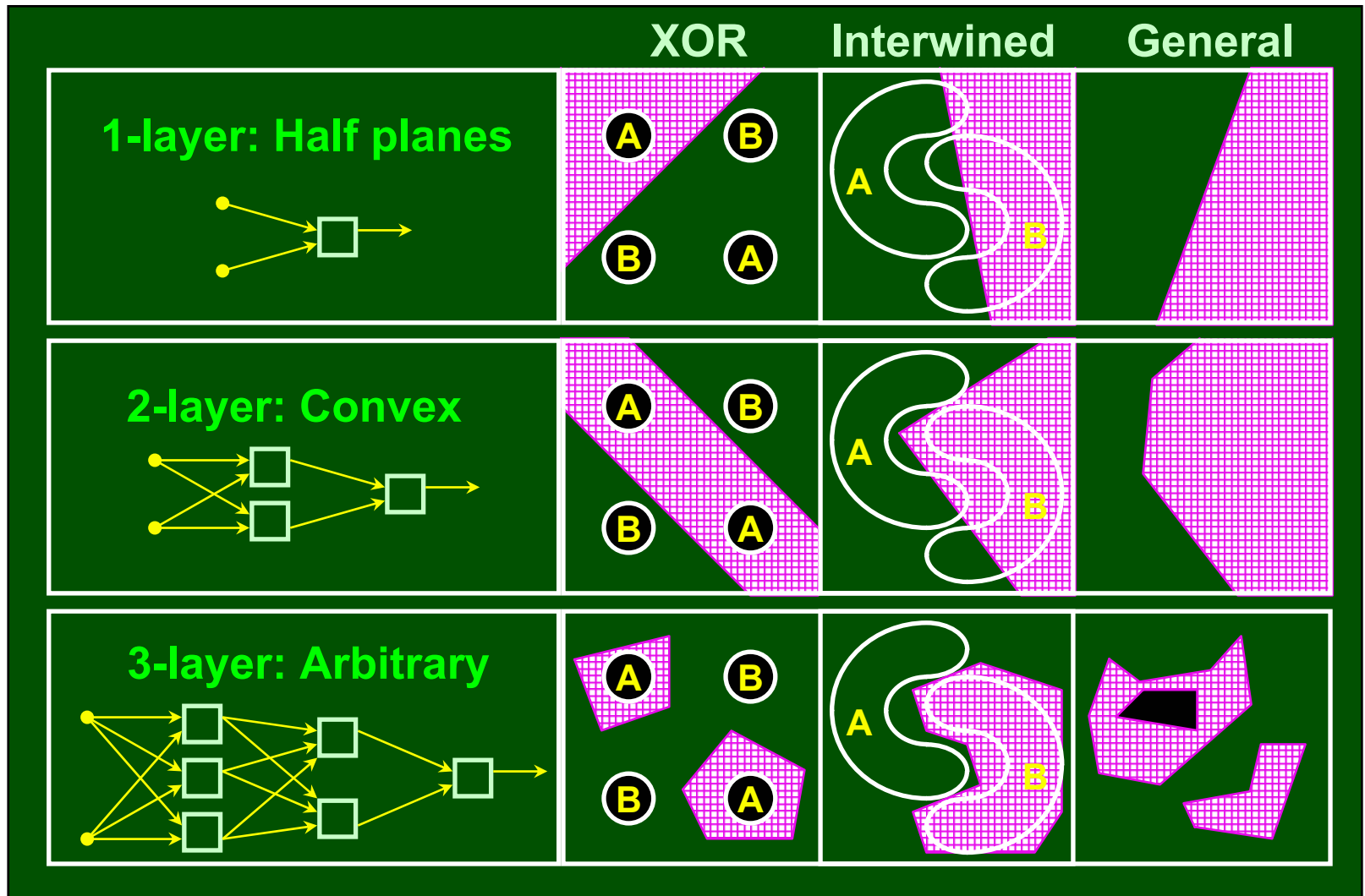
(c)

- (a) Decision boundary constructed by hidden neuron 3;
- (b) Decision boundary constructed by hidden neuron 4;
- (c) Decision boundaries constructed by the complete three-layer network

# Decision boundaries



# MLP Decision Boundaries



# Gradient-based Optimization

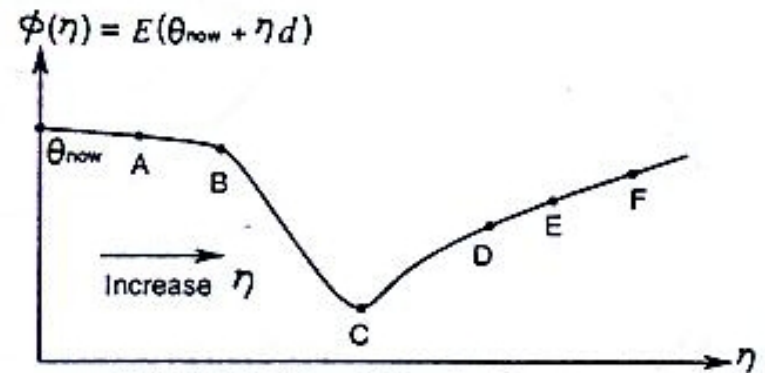
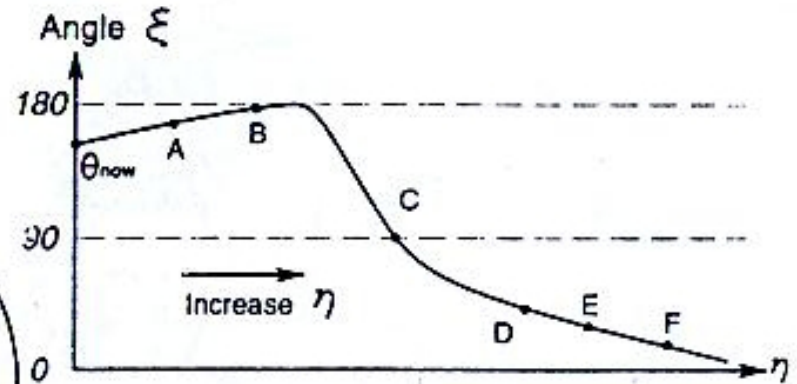
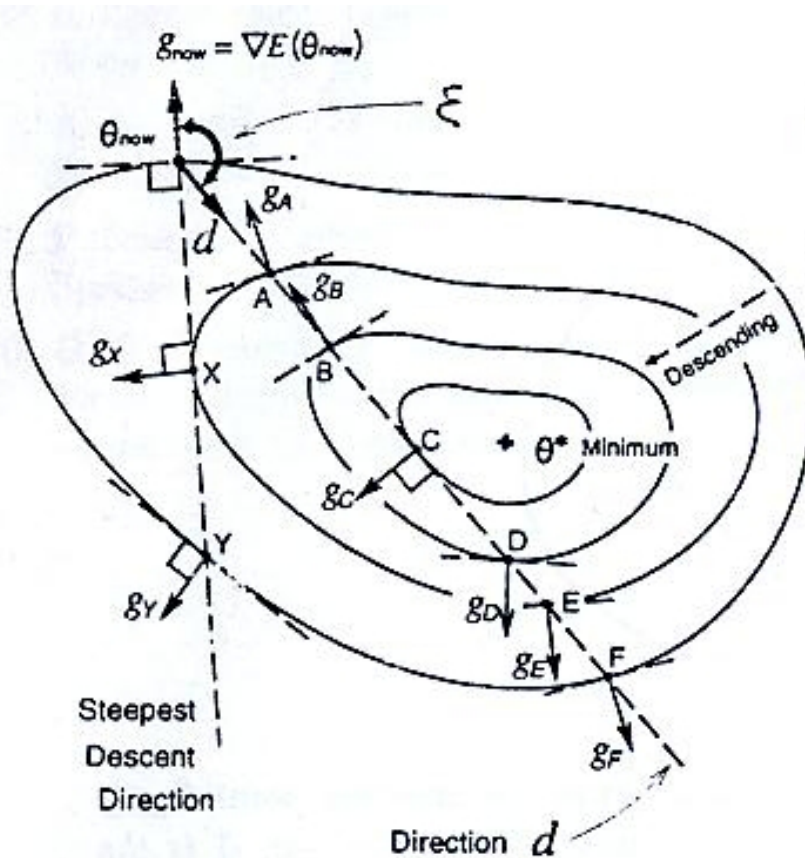
Finding a (possibly local) minimum point  $\theta^*$  that minimizes  $E(\theta)$

Steepest descent direction:

(g: gradient (derivate) of  $E(\theta)$ )

$$d = -g$$

$$g(\theta) = \frac{\partial E(\theta)}{\partial \theta}$$



A Sectional View of the Contours in Direction  $d$

# Steepest descent - Back-propagation

- **Minimizing the error  $E_p$  of the  $p^{\text{th}}$  training data**

$$e_k(p) = y_{d,k}(p) - y_k(p)$$

$$E_p = \sum_{k=1}^N (y_{d,k}(p) - y_k(p))^2$$

$$\frac{\partial E_p}{\partial y_k(p)} = -2 \cdot (y_{d,k}(p) - y_k(p))$$

- **Calculate the gradient for all the tunable variables**

$$g(\theta) = \frac{\partial E(\theta)}{\partial \theta}$$

- **Steepest descent**

$$\theta_{next} = \theta_{now} - \eta \cdot g(\theta)$$

# Steepest descent - Back-propagation

- For the neurons in the output layer:

$$\frac{\partial E_p}{\partial x_k(p)} = \frac{\partial E_p}{\partial y_k(p)} \cdot \frac{\partial y_k(p)}{\partial x_k(p)} = -2 \cdot (y_{d,k}(p) - y_k(p)) \cdot y_k(p) \cdot (1 - y_k(p))$$

$$x_k(p) = \sum_{j=1}^m x_{jk}(p) \cdot w_{jk}(p)$$

$$Y^{\text{sigmoid}} = \frac{1}{1 + e^{-X}}$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(y)}{\partial y} \cdot \frac{\partial g(x)}{\partial x} \quad \text{with } y = f(x)$$

$$\frac{\partial E_p}{\partial w_{jk}(p)} = \frac{\partial E_p}{\partial y_k(p)} \cdot \frac{\partial y_k(p)}{\partial x_k(p)} \cdot \frac{\partial x_k(p)}{\partial w_{jk}(p)} = -2 \cdot (y_{d,k}(p) - y_k(p)) \cdot y_k(p) \cdot (1 - y_k(p)) \cdot x_{jk}(p)$$

$$\frac{\partial x_k(p)}{\partial w_{jk}(p)} = x_{jk}(p)$$

$$e_k(p) = y_{d,k}(p) - y_k(p)$$

$$\delta_k(p) = y_k(p) \cdot [1 - y_k(p)] \cdot e_k(p)$$

$$w_{\text{next}} = w_{\text{now}} - \eta \cdot g(w)$$

$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p)$$

# Steepest descent - Back-propagation

- For the neurons in the **hidden layer**:

**Ordered derivates** (not partial, because of the dependences):

$$\frac{\partial f(x, g(x))}{\partial x} = \frac{\partial f(x, y)}{\partial x} \Big|_{y=f(x)} + \frac{\partial f(x, y)}{\partial y} \Big|_{y=f(x)} \cdot \frac{\partial g(x)}{\partial x}$$

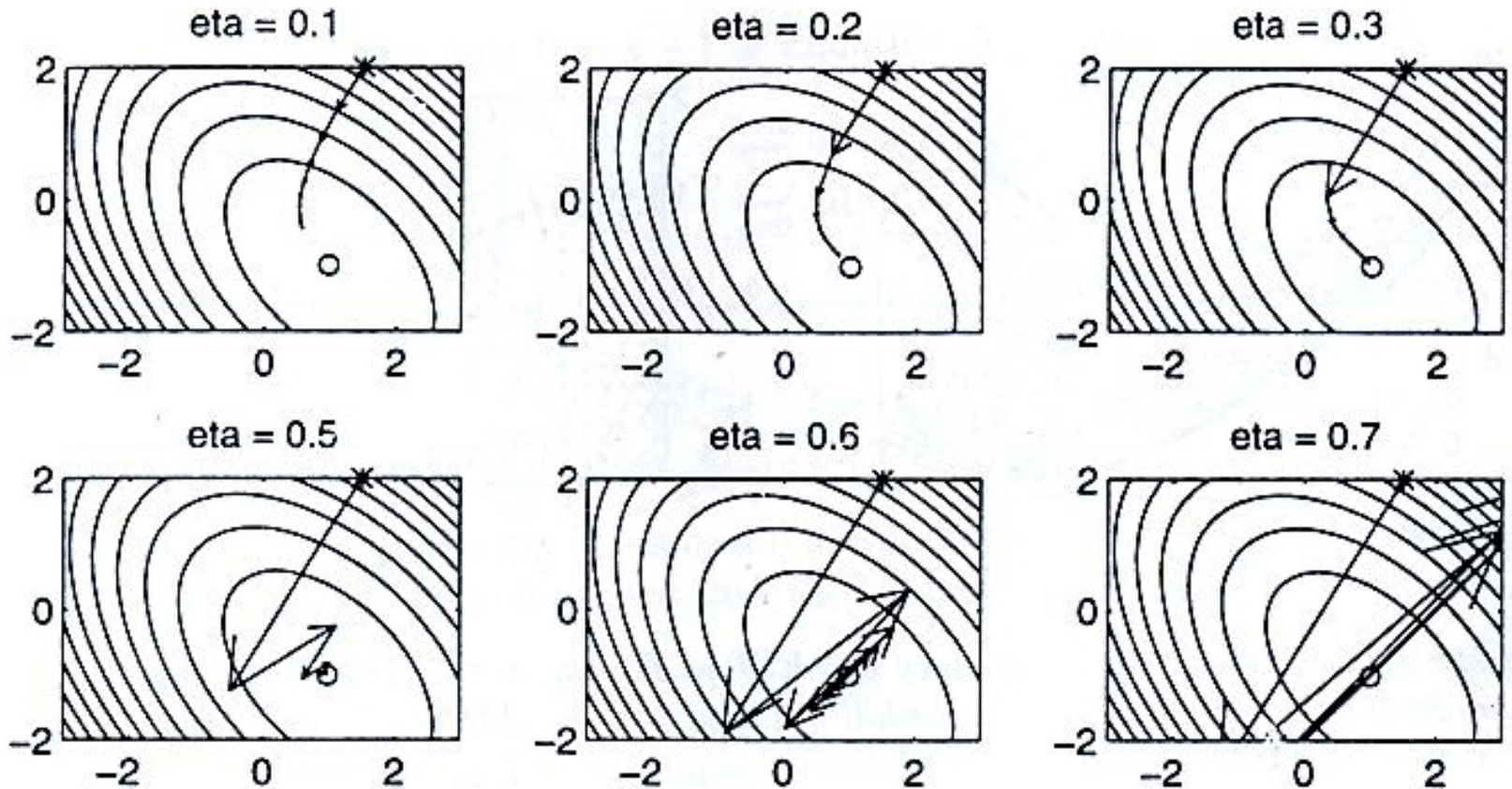
$$\delta_j(p) = y_j(p) \cdot [1 - y_j(p)] \cdot \sum_{k=1}^l \delta_k(p) w_{jk}(p)$$

$$\Delta w_{ij}(p) = \alpha \cdot x_i(p) \cdot \delta_j(p)$$

# Steepest descent step size

- The effect of the step size  $\eta$

$$\theta_{next} = \theta_{now} - \eta \cdot g(\theta)$$



# Accelerated learning

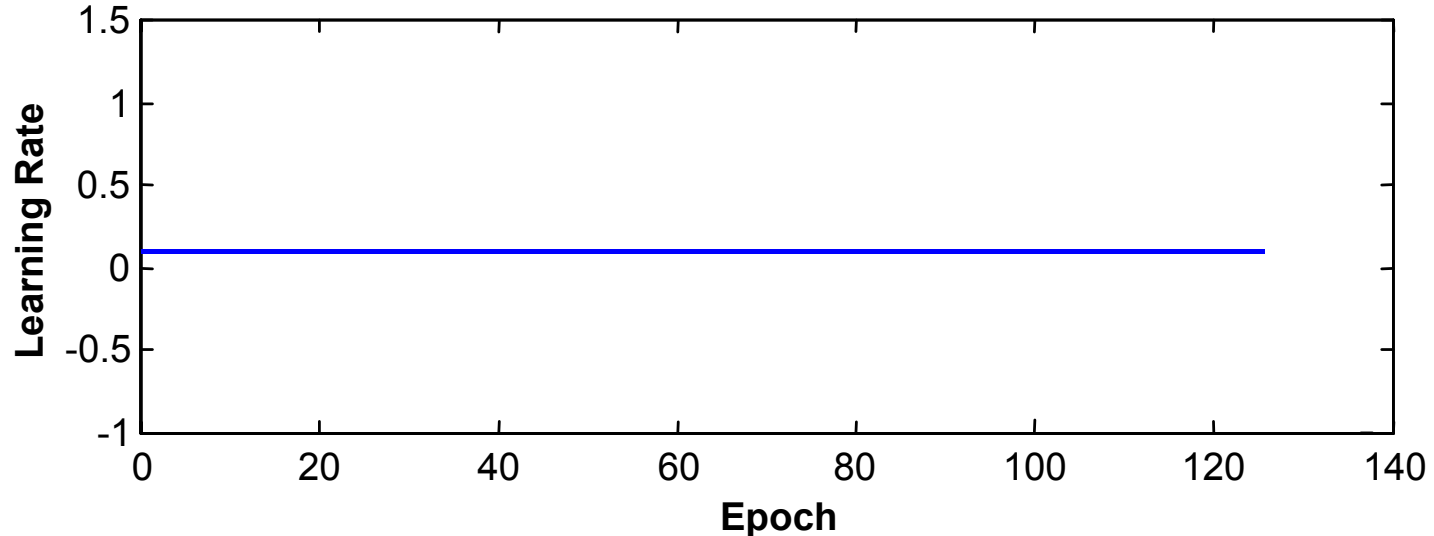
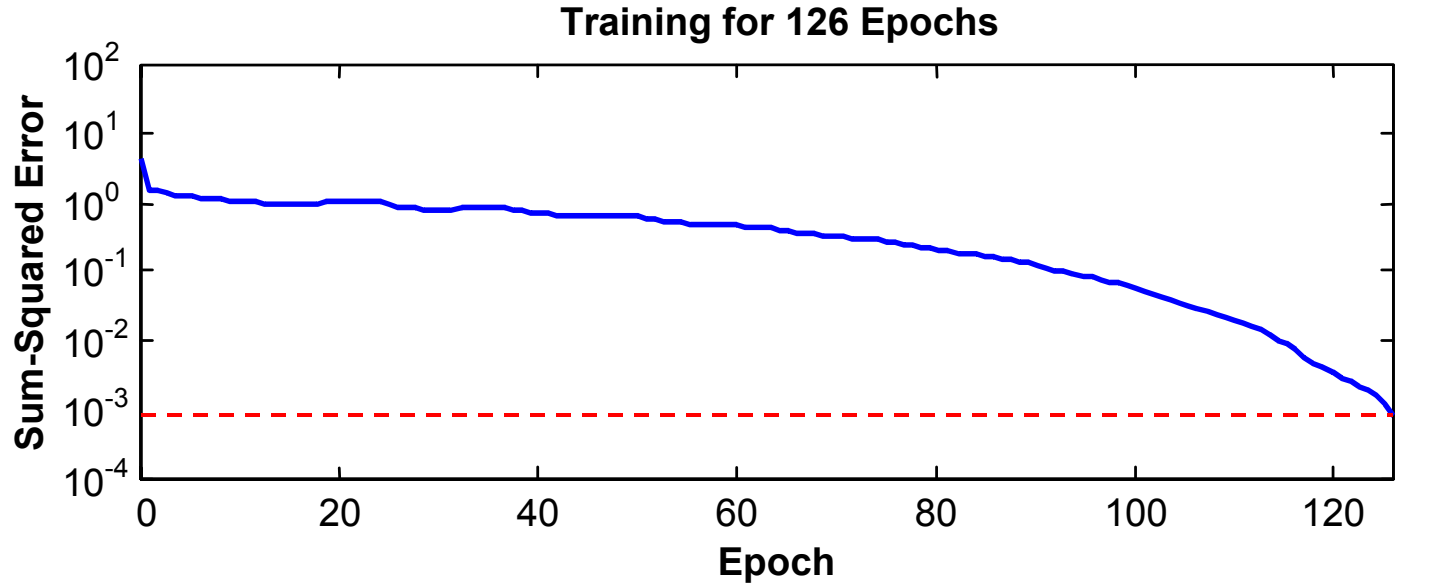
- We also can accelerate training by including a **momentum term** in the delta rule:

$$\Delta w_{jk}(p) = \beta \cdot \Delta w_{jk}(p-1) + \alpha \cdot y_j(p) \cdot \delta_k(p)$$

where  $\beta$  is a positive number ( $0 \leq \beta < 1$ ) called the **momentum constant**. Typically, the momentum constant is set to 0.95.

- This equation is called the **generalised delta rule**.

# Learning with momentum for operation *Exclusive-OR*



# Learning with adaptive learning rate

To accelerate the convergence and yet avoid the danger of instability, we can apply two heuristics:

## Heuristic 1

If the change of the sum of squared errors has the same algebraic sign for several consequent epochs, then the learning rate parameter,  $\alpha$ , should be increased.

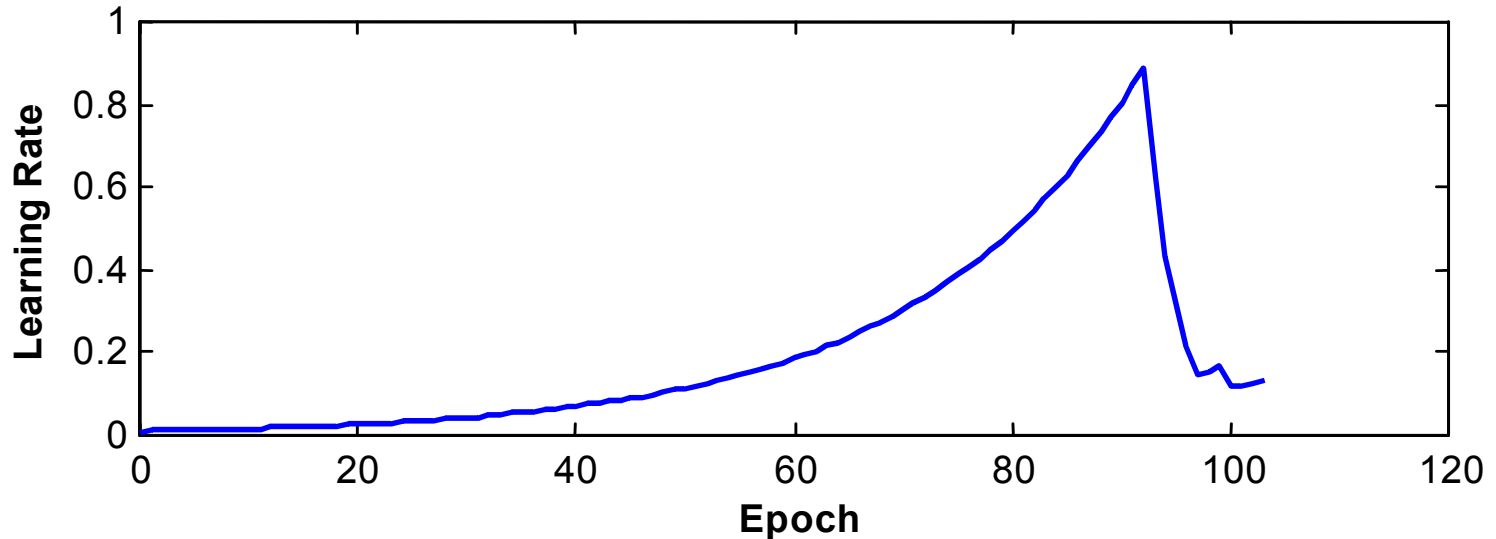
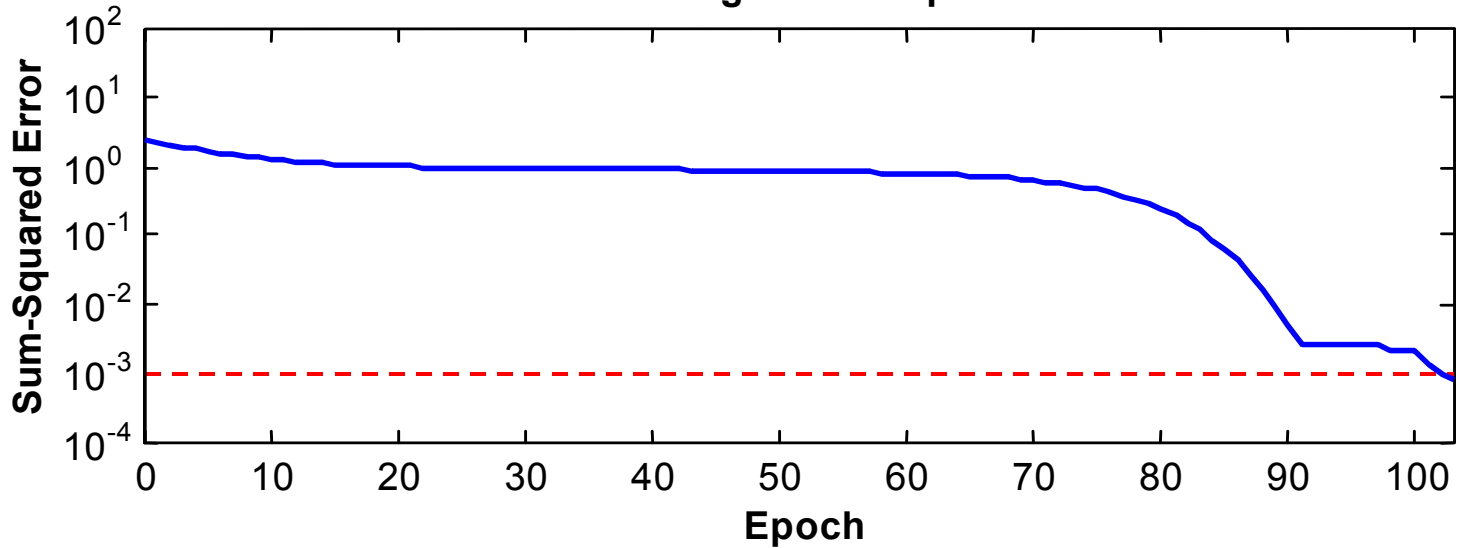
## Heuristic 2

If the algebraic sign of the change of the sum of squared errors alternates for several consequent epochs, then the learning rate parameter,  $\alpha$ , should be decreased.

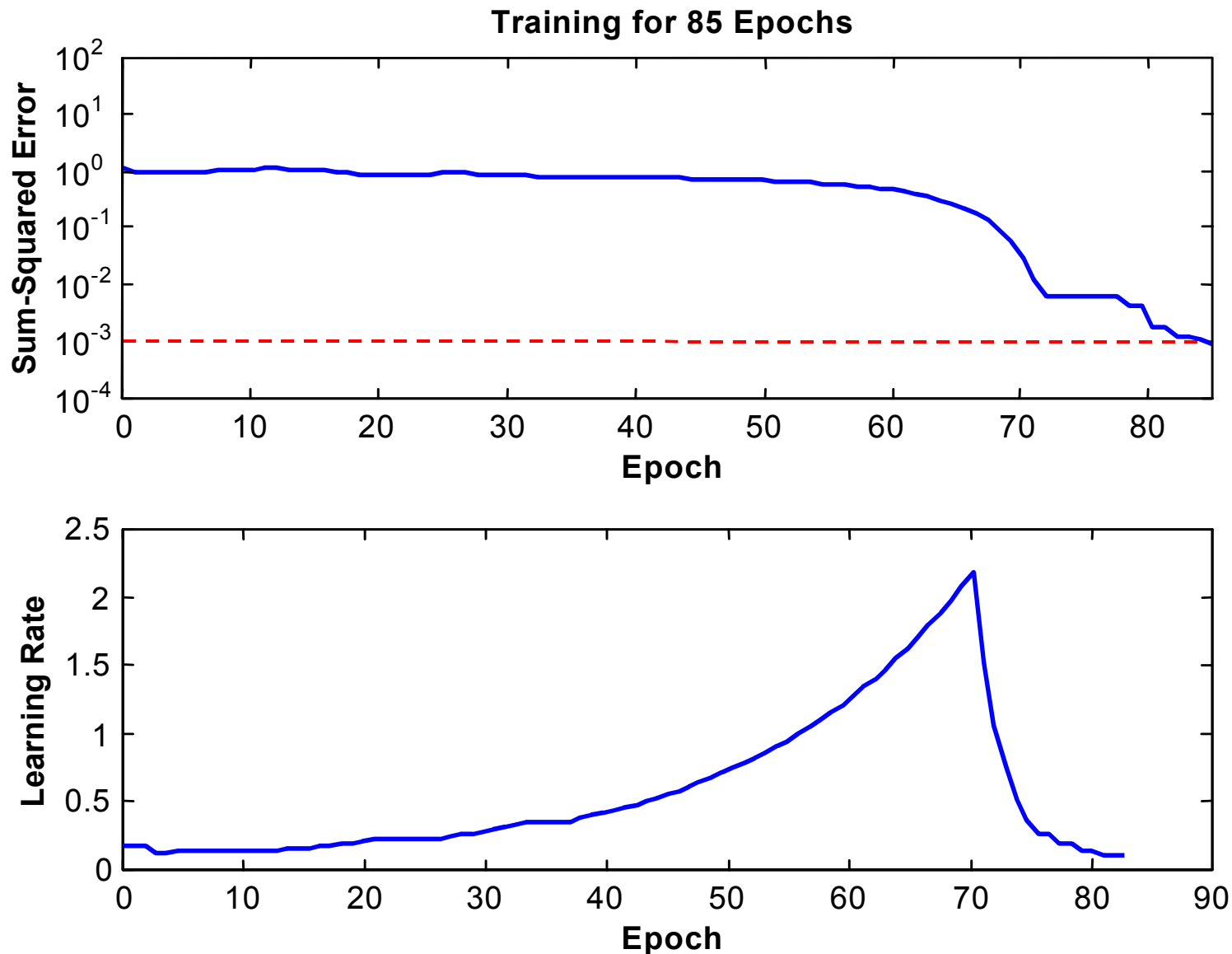
- Adapting the learning rate requires some changes in the back-propagation algorithm.
- If the sum of squared errors at the current epoch exceeds the previous value by more than a predefined ratio (typically 1.04), the learning rate parameter is decreased (typically by multiplying by 0.7) and new weights and thresholds are calculated.
- If the error is less than the previous one, the learning rate is increased (typically by multiplying by 1.05).

# Learning with adaptive learning rate

Training for 103 Epochs



# Learning with momentum and adaptive learning rate



# The Hopfield Network

- Neural networks were designed on analogy with the brain. The brain's memory, however, works by association.
- For example, we can recognise a familiar face even in an unfamiliar environment within 100-200 ms.
- We can also recall a complete sensory experience, including sounds and scenes, when we hear only a few bars of music.
- The brain routinely associates one thing with another.

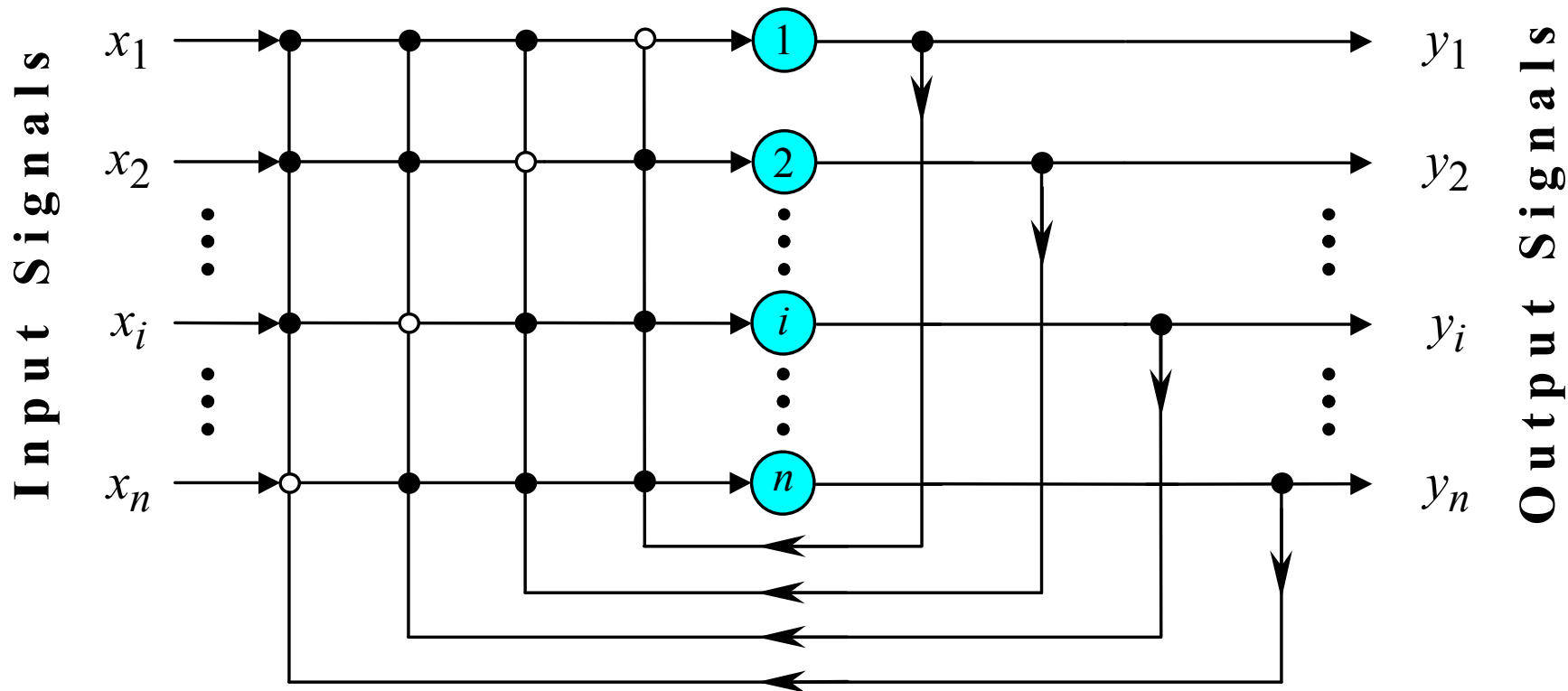
# The Hopfield Network

- Multilayer neural networks trained with the back-propagation algorithm are used for pattern recognition problems. However, to emulate the human memory's associative characteristics we need a different type of network: a **recurrent neural network**.
- A recurrent neural network has feedback loops from its outputs to its inputs. The presence of such loops has a profound impact on the learning capability of the network.

# The Hopfield Network

- The stability of recurrent networks intrigued several researchers in the 1960s and 1970s. However, none was able to predict which network would be stable, and some researchers were pessimistic about finding a solution at all.
- The problem was solved only in 1982, when **John Hopfield** formulated the physical principle of storing information in a dynamically stable network.

# Single-layer n-neuron Hopfield network



# The Hopfield Network

- The Hopfield network uses **McCulloch and Pitts neurons** with the *sign activation function* as its computing element:

$$Y^{sign} = \begin{cases} +1, & \text{if } X > 0 \\ -1, & \text{if } X < 0 \\ Y, & \text{if } X = 0 \end{cases}$$

# The Hopfield Network

- The current state of the Hopfield network is determined by the current outputs of all neurons,

$$y_1, y_2, \dots, y_n.$$

Thus, for a single-layer  $n$ -neuron network, the state can be defined by the **state vector** as:

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

# The Hopfield Network

- In the Hopfield network, synaptic weights between neurons are usually represented in matrix form as follows:

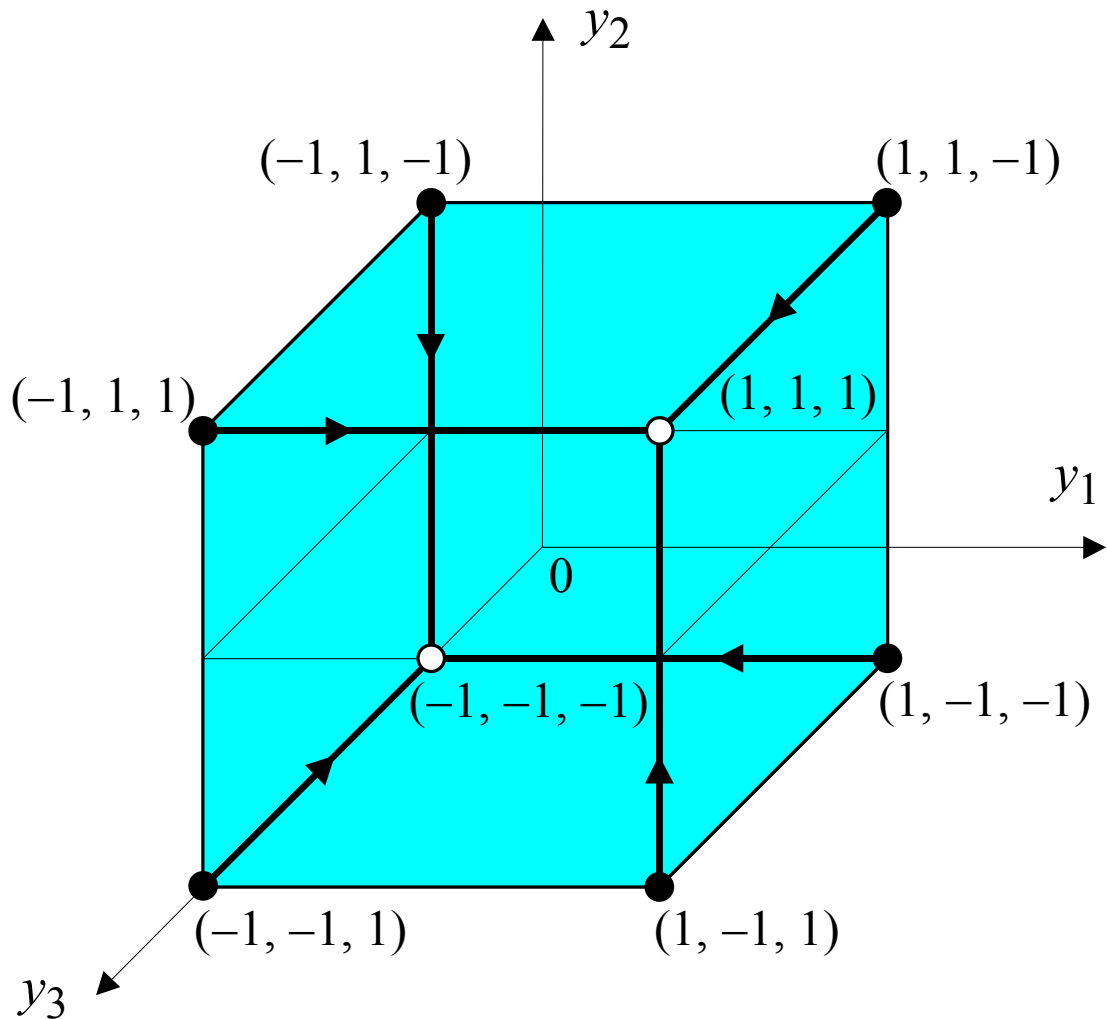
$$\mathbf{W} = \sum_{m=1}^M \mathbf{Y}_m \mathbf{Y}_m^T - M \mathbf{I}$$

$$\mathbf{Y}(p) = \text{sign}\{\mathbf{W} \cdot \mathbf{X}(p)\}$$

$$\mathbf{X}(p+1) = \mathbf{Y}(p)$$

- where  $M$  is the number of states to be memorised by the network,  $\mathbf{Y}_m$  is the  $n$ -dimensional binary vector,  $\mathbf{I}$  is  $n \times n$  identity matrix, and superscript  $T$  denotes a matrix transposition.

# Possible states for the three-neuron Hopfield network



# The Hopfield Network

- The stable state-vertex is determined by the weight matrix  $\mathbf{W}$ , the current input vector  $\mathbf{X}$ , and the threshold matrix  $\theta$ . If the input vector is partially incorrect or incomplete, the initial state will converge into the stable state-vertex after a few iterations.
- Suppose, for instance, that our network is required to memorise two opposite states,  $(1, 1, 1)$  and  $(-1, -1, -1)$ . Thus,

$$\mathbf{Y}_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{Y}_2 = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \quad \text{or} \quad \mathbf{Y}_1^T = [1 \quad 1 \quad 1] \quad \mathbf{Y}_2^T = [-1 \quad -1 \quad -1]$$

where  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$  are the three-dimensional vectors.

– The  $3 \times 3$  identity matrix  $\mathbf{I}$  is

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

– Thus, we can now determine the weight matrix as follows:

$$\mathbf{W} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & -1 & -1 \end{bmatrix} - 2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 2 & 2 \\ 2 & 0 & 2 \\ 2 & 2 & 0 \end{bmatrix}$$

– Next, the network is tested by the sequence of input vectors,  $\mathbf{X}_1$  and  $\mathbf{X}_2$ , which are equal to the output (or target) vectors  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$ , respectively.

– First, we activate the Hopfield network by applying the input vector  $\mathbf{X}$ . Then, we calculate the actual output vector  $\mathbf{Y}$ , and finally, we compare the result with the initial input vector  $\mathbf{X}$ .

$$\mathbf{Y}_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{Y}_2 = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

$$\mathbf{Y}(p) = \text{sign} \{ \mathbf{W} \cdot \mathbf{X}_1 \} = \text{sign} \left\{ \begin{bmatrix} 0 & 2 & 2 \\ 2 & 0 & 2 \\ 2 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \mathbf{Y}_1 = \mathbf{X}_1$$

$$\mathbf{Y}(p) = \text{sign} \{ \mathbf{W} \cdot \mathbf{X}_2 \} = \text{sign} \left\{ \begin{bmatrix} 0 & 2 & 2 \\ 2 & 0 & 2 \\ 2 & 2 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \right\} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = \mathbf{Y}_2 = \mathbf{X}_2$$

# Hopfield network

$$\mathbf{W} = \sum_{m=1}^M \mathbf{Y}_m \mathbf{Y}_m^T - M \mathbf{I}$$

$$\mathbf{Y}(p) = \mathbf{W} \cdot \mathbf{X}(p) \quad \mathbf{X}(p+1) = \mathbf{Y}(p)$$

$$\mathbf{Y}(p+1) = \mathbf{W} \cdot \mathbf{Y}(p) = \left[ \sum_{m=1}^M \mathbf{Y}_m \cdot \mathbf{Y}_m^T - M \cdot \mathbf{I} \right] \cdot \mathbf{Y}(p)$$

$$\mathbf{Y}(p) = \sum_{m=1}^M \mathbf{Y}_m \cdot \mathbf{Y}_m^T \cdot \mathbf{Y}(p) - M \cdot \mathbf{I} \cdot \mathbf{Y}(p)$$

$$\mathbf{Y}_m^T \cdot \mathbf{Y}_m = M \quad \text{Max! E.g. } \mathbf{Y}_2^T \cdot \mathbf{Y}_2 = [-1 \quad -1 \quad -1] \cdot \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = 3$$

# Hopfield network

- The remaining six states are all unstable. However, stable states (also called **fundamental memories**) are capable of attracting states that are close to them.
- The fundamental memory  $(1, 1, 1)$  attracts unstable states  $(-1, 1, 1)$ ,  $(1, -1, 1)$  and  $(1, 1, -1)$ . Each of these unstable states represents a single error, compared to the fundamental memory  $(1, 1, 1)$ .
- The fundamental memory  $(-1, -1, -1)$  attracts unstable states  $(-1, -1, 1)$ ,  $(-1, 1, -1)$  and  $(1, -1, -1)$ .
- Thus, the Hopfield network can act as an **error correction network**.

# Storage capacity of the Hopfield network

- **Storage capacity** is or the largest number of fundamental memories that can be stored and retrieved correctly.
- The maximum number of fundamental memories  $M_{max}$  that can be stored in the  $n$ -neuron recurrent network is limited by

$$M_{max} = 0.15 n$$

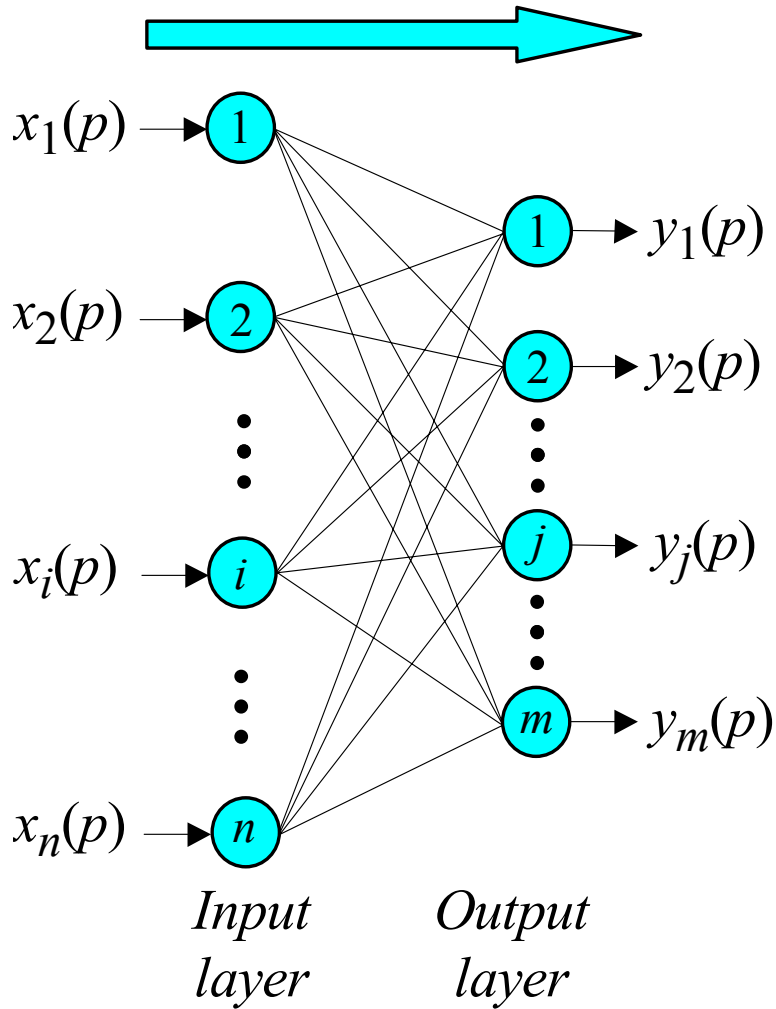
# Bidirectional associative memory (BAM)

- The Hopfield network represents an **autoassociative** type of memory – it can retrieve a corrupted or incomplete memory but cannot associate this memory with another different memory.
- Human memory is essentially **associative**. One thing may remind us of another, and that of another, and so on. We use a chain of mental associations to recover a lost memory. If we forget where we left an umbrella, we try to recall where we last had it, what we were doing, and who we were talking to. We attempt to establish a chain of associations, and thereby to restore a lost memory.

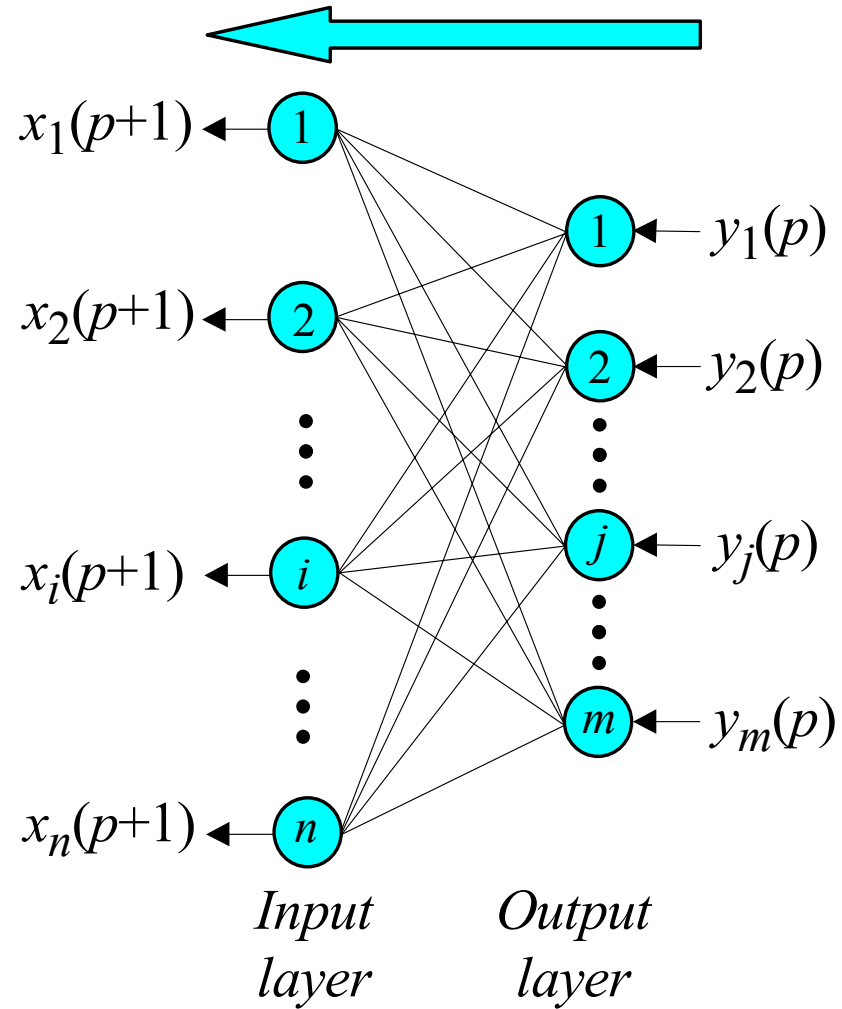
# Bidirectional associative memory (BAM)

- To associate one memory with another, we need a recurrent neural network capable of accepting an input pattern on one set of neurons and producing a related, but different, output pattern on another set of neurons.
- **Bidirectional associative memory (BAM)**, first proposed by **Bart Kosko**, is a heteroassociative network. It associates patterns from one set, set  $A$ , to patterns from another set, set  $B$ , and vice versa. Like a Hopfield network, the BAM can generalise and also produce correct outputs despite corrupted or incomplete inputs.

# BAM operation



(a) Forward direction.



(b) Backward direction.

# Bidirectional associative memory (BAM)

The basic idea behind the BAM is to store pattern pairs so that when  $n$ -dimensional vector  $\mathbf{X}$  from set  $A$  is presented as input, the BAM recalls  $m$ -dimensional vector  $\mathbf{Y}$  from set  $B$ , but when  $\mathbf{Y}$  is presented as input, the BAM recalls  $\mathbf{X}$ .

# Bidirectional associative memory (BAM)

- To develop the BAM, we need to create a correlation matrix for each pattern pair we want to store. The correlation matrix is the matrix product of the input vector  $\mathbf{X}$ , and the transpose of the output vector  $\mathbf{Y}^T$ . The BAM weight matrix is the sum of all correlation matrices, that is,

$$\mathbf{W} = \sum_{m=1}^M \mathbf{X}_m \mathbf{Y}_m^T$$

- where  $M$  is the number of pattern pairs to be stored in the BAM.

# Bidirectional associative memory (BAM)

$$\mathbf{W} = \sum_{m=1}^M \mathbf{X}_m \mathbf{Y}_m^T$$

$$\mathbf{Y}(p+1) = \mathbf{W}^T \cdot \mathbf{X}(p) = \sum_{m=1}^M \mathbf{Y}_m \cdot \mathbf{X}_m^T \cdot \mathbf{X}(p)$$

$$\mathbf{X}(p+1) = \mathbf{W} \cdot \mathbf{Y}(p) = \sum_{m=1}^M \mathbf{X}_m \cdot \mathbf{Y}_m^T \cdot \mathbf{Y}(p)$$

$$\mathbf{X}_m^T \cdot \mathbf{X}_m = M \quad \text{Max!} \quad \text{E.g.} \quad \mathbf{X}_2^T \cdot \mathbf{X}_2 = [-1 \quad -1 \quad -1] \cdot \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} = 3$$
$$\mathbf{Y}_m^T \cdot \mathbf{Y}_m = M$$

# Stability and storage capacity of the BAM

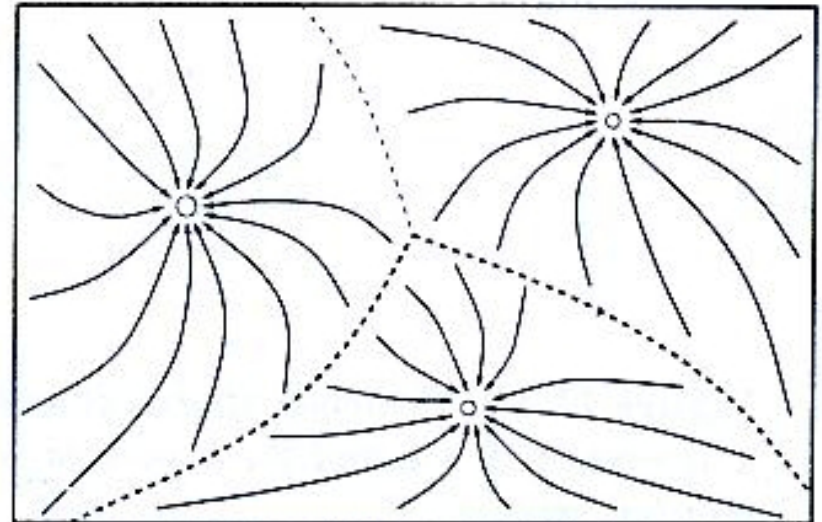
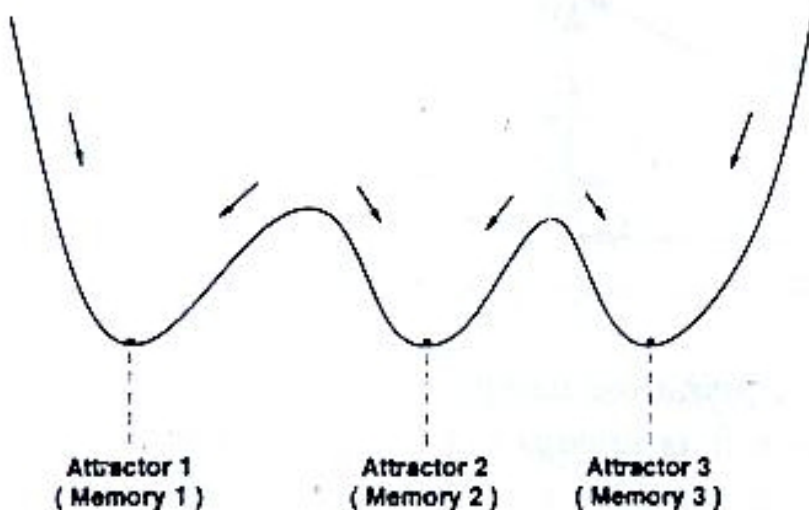
- The BAM is **unconditionally stable**. This means that any set of associations can be learned without risk of instability.
- The maximum number of associations to be stored in the BAM should not exceed the number of neurons in the smaller layer.
- The more serious problem with the BAM is **incorrect convergence**. The BAM may not always produce the closest association. In fact, a stable association may be only slightly related to the initial input vector.

# Hopfield network

- Recording learning system
- Associative memory network
- Content addressable
- Auto-associative

# Hopfield network

- Any physical system whose dynamics in phase space is dominated by a substantial number of locally stable states to which it is attracted can therefore be regarded as a **general content-addressable memory**.
- The physical system will be potentially useful memory if, in addition, any prescribed set of states can readily be made the stable states of the system.

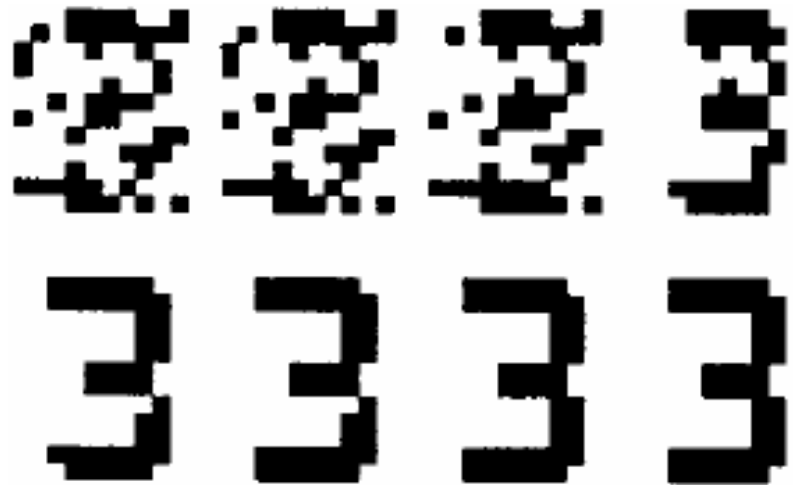


# Hopfield examples

- Character recognition  
eight exemplar patterns

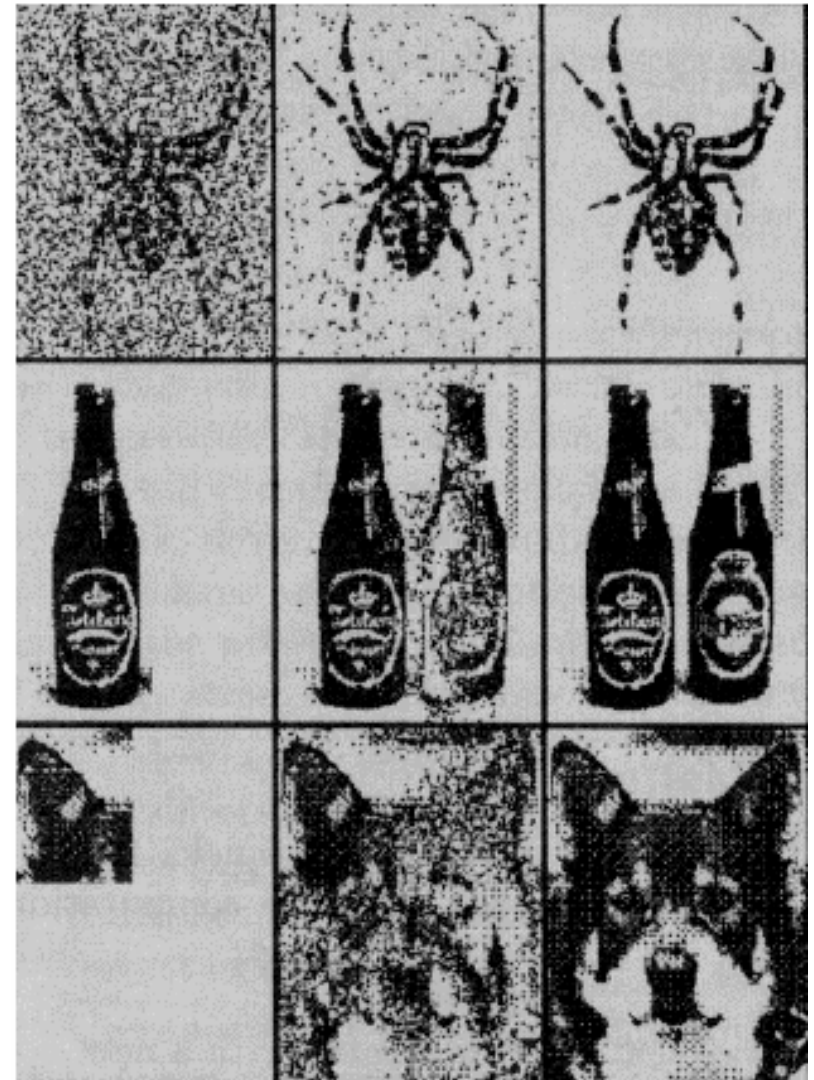


- Output pattern for noisy “3” input



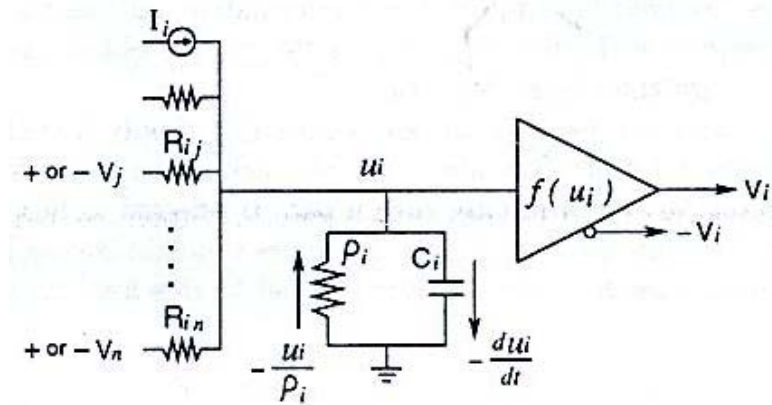
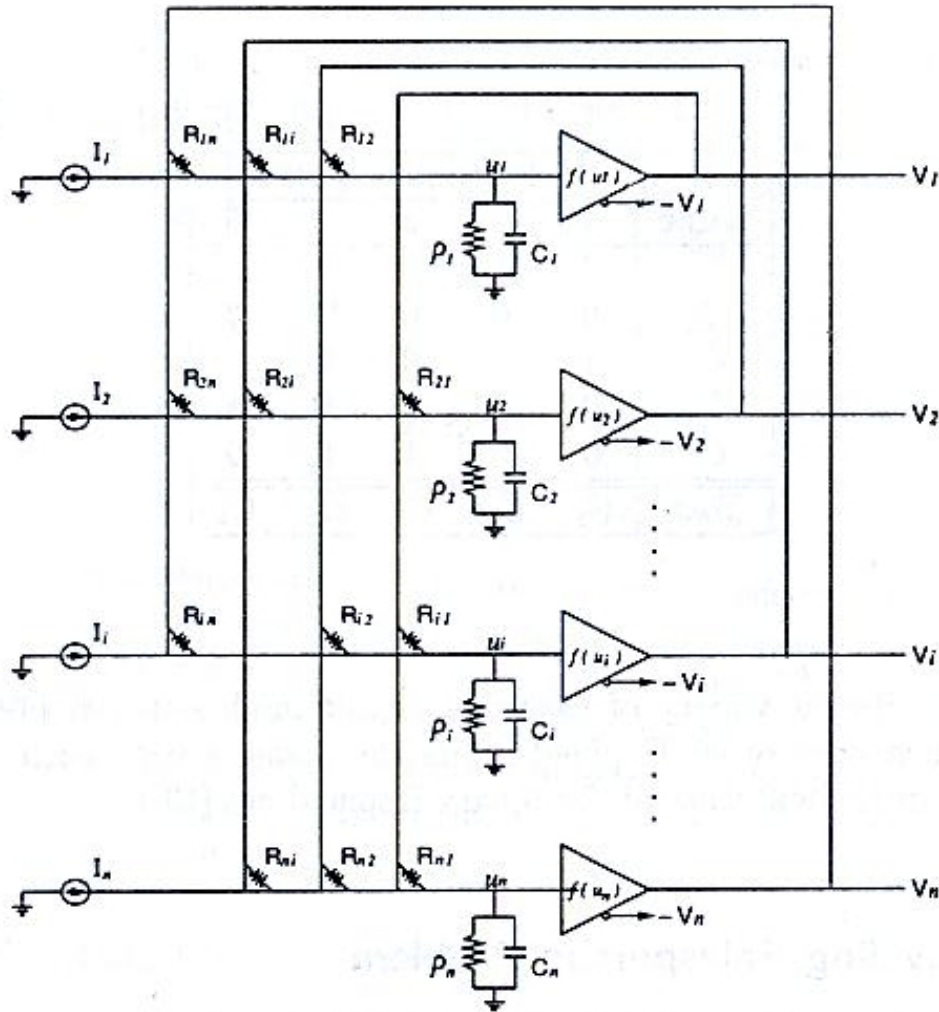
# Hopfield examples

- **Reconstruction of Images**  
binary images are 130x180 pixels



# Continuous-Valued Hopfield network

- Electrical implementation



# Ajánlott irodalom

- **The slides of this lecture are partially based on the books:**

**Kóczy T. László és Tikk Domonkos: *Fuzzy rendszerek*,  
Typotex Kiadó, 2000, ISBN 963-9132-55-1**

**J.-S. R. Jang, C.-T. Sun, E. Mizutani: *Neuro-Fuzzy and Soft  
Computing*, Prentice Hall, 1997, ISBN 0-13-261066-3**

**Michael Negnevitsky: *Artificial Intelligence: A Guide to  
Intelligent Systems*, Addison Wesley, Pearson Education  
Limited, 2002, ISBN 0201-71159-1**