

Soft Computing: Neural Networks

---

# Neural Networks

(Chapter 9)

**Kai Goebel, Bill Cheetham**  
**GE Corporate Research & Development**  
goebel@cs.rpi.edu  
cheetham@cs.rpi.edu

1

Soft Computing: Neural Networks

---

## Outline

- Introduction
- Categories
- Hopfield Net
- Perceptron
  - Single Layer
  - Multi Layer

2

Soft Computing: Neural Networks

## Introduction

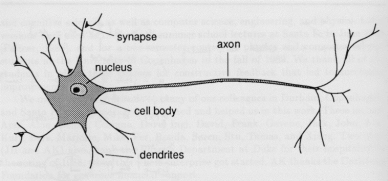
Human brain is superior to digital computer at many tasks

- + e.g., processing of visual information
- + robust and fault tolerant (nerve cells in the brain die every day)
- + flexible; adjusts to new environment
- + can deal with information that is sparse, imprecise, noisy, inconsistent
- + highly parallel
- + small, compact, dissipates very little power
- slower in primarily (simple) arithmetic operations

3

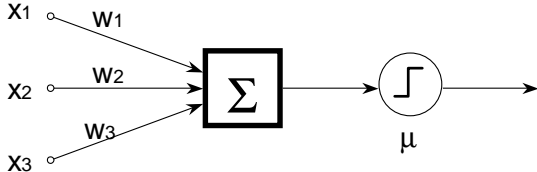
Soft Computing: Neural Networks

## Neurons



**McCulloch & Pitts (1943)**

- simple model of neuron as a binary threshold unit
- uses step function to “fire” when threshold  $\mu$  is surpassed



4

	<p>Soft Computing: Neural Networks</p> <h2>Real Neurons</h2>
5	<h3>Real Neurons</h3> <ul style="list-style-type: none"> <li>- use not even approximately threshold devices</li> <li>- it is assumed they use a non-linear summation method</li> <li>- produce a sequence of pulses (not a single output level)</li> <li>- do not have the same fixed delay (<math>t \rightarrow t+1</math>)</li> <li>- are not updated synchronously</li> <li>- amount of transmitter substance varies unpredictably</li> </ul>

	<p>Soft Computing: Neural Networks</p> <h2>Issues</h2>
6	<p>What does that leave us with?</p> <p>What is the best architecture? (layers, connections, activation functions, updating, # units?)</p> <p>How can it be programmed? (can it learn, # examples needed, time to learn, amount of supervision, real-time learning)</p> <p>What can it do? (how many tasks, how well, how fast, how robust, level of generalization)</p>

Soft Computing: Neural Networks

## Neural Nets: Categorization

**Supervised Learning**

- Multilayer perceptrons
- Radial basis function networks
- Modular neural networks
- LVQ (learning vector quantization)

**Reinforcement Learning**

- Temporal Difference Learning
- Q-Learning

**Unsupervised Learning**

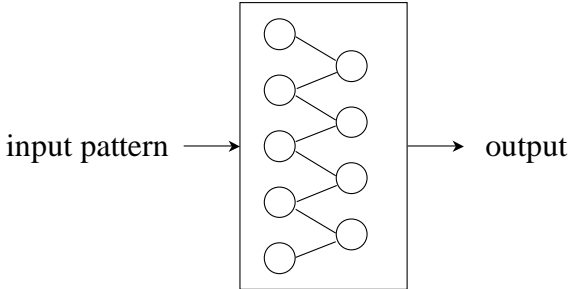
- Competitive learning networks
- Kohonen self-organizing networks
- ART (adaptive resonant theory)

7

Soft Computing: Neural Networks

## Supervised Neural Networks

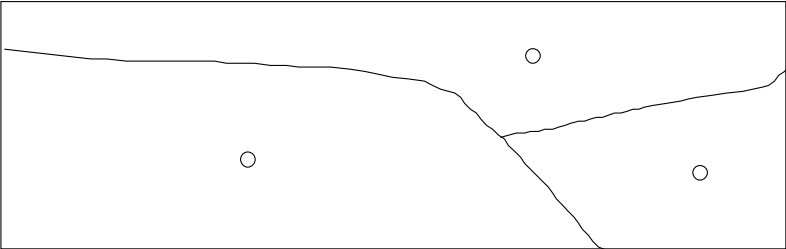
**Requirement:**  
known input-output relations



The diagram illustrates a supervised neural network. On the left, the text "input pattern" is followed by an arrow pointing into a rectangular box. Inside the box, there is a network of nodes: a vertical column of five circles on the left, and a vertical column of three circles on the right. Lines connect each of the five left nodes to the top-left node of the right column, and each of the three top nodes of the left column to the top-right node of the right column. An arrow points from the right side of the box to the text "output".

8

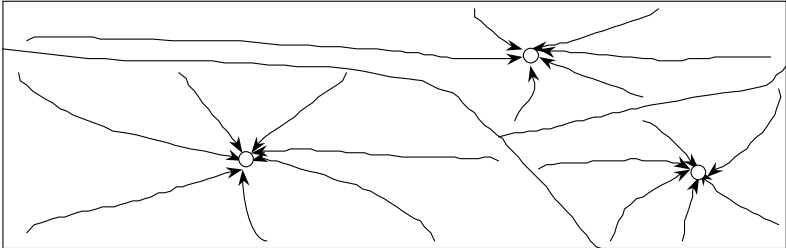
	<p>Soft Computing: Neural Networks</p> <h2>Hopfield Model</h2>
9	<p>Associative Memory is considered the “fruit fly” of this field.</p> <p>It illustrates in the simplest possible manner the way that collective computation can work.</p> <p>Store a set of patterns in such a way that when presented with a new pattern, the network responds by producing the closest stored pattern.</p> <p>Conventional approach: store a list of patterns, compute the Hamming distance, find the smallest, et voila!</p>

	<p>Soft Computing: Neural Networks</p> <h2>Hopfield Network Operation</h2>
10	<p>Picture is pattern; stored as attractor in the configuration space.</p> <p>From arbitrary starting points, one attractor will be found</p> 

Soft Computing: Neural Networks

## Hopfield Network Operation

Picture is pattern; stored as attractor in the configuration space.  
 From arbitrary starting points, one attractor will be found

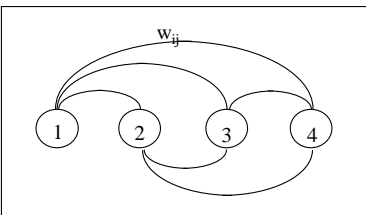


11

Soft Computing: Neural Networks

## Hopfield Architecture

- Recurrent Network
- Symmetric Architecture
- Evaluation until no more changes are observed  
 i.e., network settles into local minimum config.



- local minimum corresponds to “energy function”

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} x_i x_j$$

12

Soft Computing: Neural Networks

## Hopfield Network Equations

The operative equation, i.e., the network output at each step is

$$y_i = \operatorname{sgn}\left(\sum_j w_{ij}x_j\right)$$

where

$$\operatorname{sgn}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$

13

Soft Computing: Neural Networks

## Learning in Hopfield Models

Learning Rule:

$$w_{ij}^{(n+1)} = w_{ij}^{(n)} + x_i x_j$$

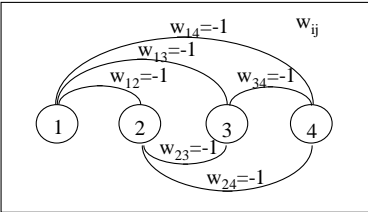
$$w_{ii} = 0$$

14

Soft Computing: Neural Networks

## Hopfield Example

Learn  $x=[1 \ 1 \ -1 \ -1]$   
 which gives us the weight matrix

$$w = \begin{bmatrix} 0 & 1 & -1 & -1 \\ 1 & 0 & -1 & -1 \\ -1 & -1 & 0 & 1 \\ -1 & -1 & 1 & 0 \end{bmatrix}$$


Now let's check the slightly corrupted pattern  
 $p=[1 \ 1 \ -1 \ 1]$   
 which will restore the pattern found close  
 $y=[1 \ 1 \ -1 \ -1]$   
 with an energy level of  $E=-6$

15

Soft Computing: Neural Networks

## Hopfield Example

Learn second pattern  $x=[-1 \ -1 \ 1 \ 1]$   
 which gives us the new weight matrix

$$w = \begin{bmatrix} 0 & 2 & -2 & -2 \\ 2 & 0 & -2 & -2 \\ -2 & -2 & 0 & 2 \\ -2 & -2 & 2 & 0 \end{bmatrix}$$

Now let's check the slightly corrupted pattern  
 $p=[-1 \ -1 \ -1 \ 1]$   
 which will restore the pattern  
 $y=[-1 \ -1 \ 1 \ 1]$   
 with an energy level of  $E=-12$

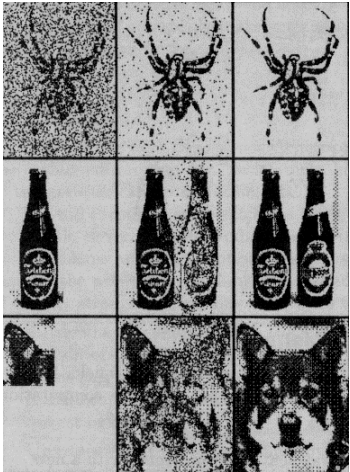
16



Soft Computing: Neural Networks

## More Complex Hopfield Examples

Reconstruction of Images



binary images are 130x180 pixels

17

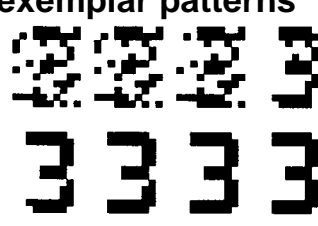
Soft Computing: Neural Networks

## Hopfield Book Example

Character Recognition

0 1 2 3  
4 6 9

eight exemplar patterns



output pattern for noisy "3" input

18

	<p>Soft Computing: Neural Networks</p> <h2>Hopfield: Issues</h2>
19	<ul style="list-style-type: none"> <li>- Other memories can get lost</li> <li>- Memories are created that were not supposed to be there</li> <li>- crosstalk: if there are many memories, they might interfere</li> <li>- no emphasis on learning; rather handcrafting to get desired properties</li> <li>- goes towards optimization</li> </ul>

	<p>Soft Computing: Neural Networks</p> <h2>Perceptrons</h2>
20	<ul style="list-style-type: none"> <li>-Rosenblatt: 1950s</li> <li>-Input patterns represented is binary</li> <li>-Single layer network can be trained easily</li> <li>-Output <math>o</math> is computed by           <math display="block">o = f\left(\sum_{i=1}^n w_i x_i - \theta\right)</math> </li> </ul> <p>where</p> <ul style="list-style-type: none"> <li><math>w_i</math> is a (modifiable) weight</li> <li><math>x_i</math> is the input signal</li> <li><math>\theta</math> is some threshold (weight of constant input)</li> <li><math>f(\cdot)</math> is the activation function <math>f(x) = \text{sgn}(x) = \begin{cases} 1 &amp; \text{if } x &gt; 0 \\ 0 &amp; \text{otherwise} \end{cases}</math></li> </ul>

Soft Computing: Neural Networks

## Single-Layer Perceptrons

**Network architecture**

$y = \text{signum}(\sum w_i x_i + w_0)$

21

Soft Computing: Neural Networks

## Single-Layer Perceptron

**Example: Gender classification (according to Jang)**

**Network Arch.**

$y = \text{signum}(hw_1 + vw_2 + w_0)$

$= \begin{cases} -1 & \text{if female} \\ 1 & \text{if male} \end{cases}$

**Training data**

22

	Soft Computing: Neural Networks
	<b>Perceptron</b>
	<p><b>Learning:</b></p> <p><b>select an input vector</b></p> <p><b>if the response is incorrect, modify all weights</b></p> $\Delta w_i = \eta t_i x_i$ <p><b>where</b></p> <p><b><math>t_i</math> is a target output</b></p> <p><b><math>\eta</math> is the learning rate</b></p> <p><b>If a set of weights for converged state exists, then a method for tuning towards convergence exists</b></p> <p><b>(Rosenblatt, 1962)</b></p>
23	

	Soft Computing: Neural Networks
	<b>ADALINE</b>
	<p><b>single layer network (conceived by Widrow and Hoff)</b></p> <p><b>output is weighted linear combination of weights</b></p> $o = \sum_{i=1}^n w_i x_i - w_0$ <p><b>error is described as</b></p> $E_p = (t_p - o_p)^2 \quad \text{(for pattern p)}$ <p><b>where</b></p> <p><b><math>t_p</math> is the target output</b></p> <p><b><math>o_p</math> is the actual output</b></p>
24	

	<p>Soft Computing: Neural Networks</p> <h2>ADALINE</h2>
25	<p>To decrease the error, the derivative wrt the weights is taken</p> $\frac{\partial E_p}{\partial w_i} = -2(t_p - o_p)x_i$ <p>The delta rule is:</p> $\Delta_p w_i = \eta(t_p - o_p)x_i$ <p>Intuitive appeal:</p> <ul style="list-style-type: none"> <li>if <math>t_p &gt; o_p</math>, boost <math>o_p</math> by increasing <math>w_i x_i</math></li> <li>increase <math>w_i</math> if <math>x_i</math> is positive</li> <li>decrease <math>w_i</math> if <math>x_i</math> is negative</li> </ul>

	<p>Soft Computing: Neural Networks</p> <h2>ADALINE and MADALINE</h2>
26	<ul style="list-style-type: none"> <li>+ Simplicity of learning procedure</li> <li>+ Distributed learning; can be performed locally at node level</li> <li>+ on-line (pattern by pattern) learning</li> <li>+ connect several ADALINEs to MADALINEs to deal with XOR problem</li> <li>+ were used for noise cancellation, adaptive inverse control</li> <li>- only one layer; no suitable training method for multi-layer perceptron ... why?</li> </ul>

Soft Computing: Neural Networks

## XOR

Minsky and Papert reported a severe shortcoming of single layer perceptrons, the XOR problem...

x1	x2	output
0	0	0
0	1	1
1	0	1
1	1	0

not linearly separable

27

Soft Computing: Neural Networks

## XOR

Minsky and Papert reported a severe shortcoming of single layer perceptrons, the XOR problem...

x1	x2	output
0	0	0
0	1	1
1	0	1
1	1	0

not linearly separable



$$0w_1 + 0w_2 + w_0 \leq 0 \Leftrightarrow w_0 \leq 0$$

$$0w_1 + 1w_2 + w_0 > 0 \Leftrightarrow w_2 > -w_0$$

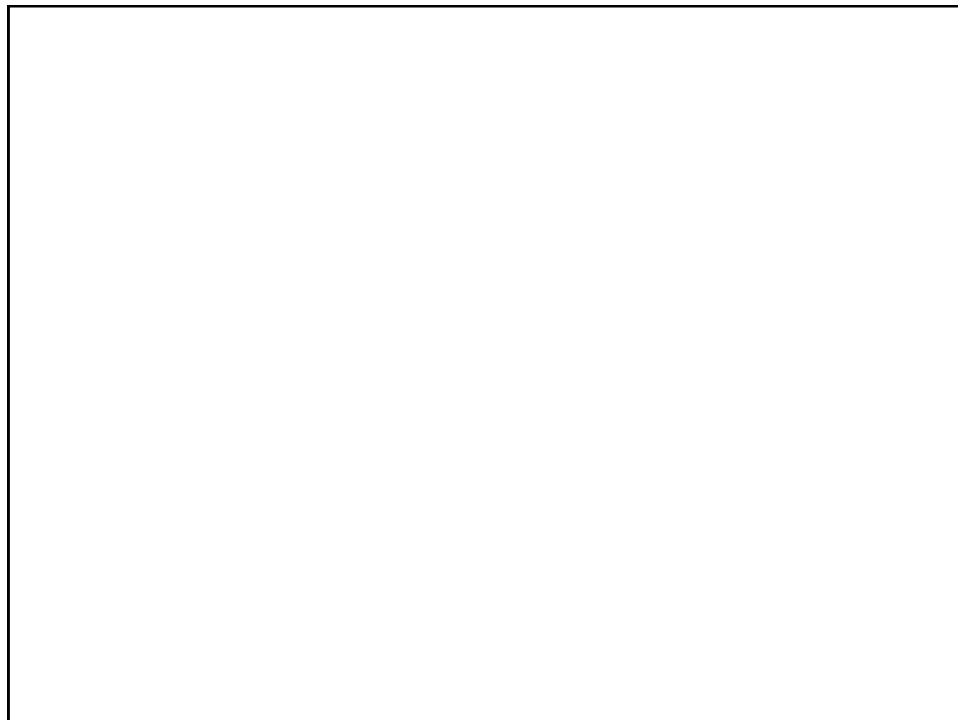
$$1w_1 + 0w_2 + w_0 > 0 \Leftrightarrow w_1 > -w_0$$

$$1w_1 + 1w_2 + w_0 \leq 0 \Leftrightarrow w_1 + w_2 \leq -w_0$$

28

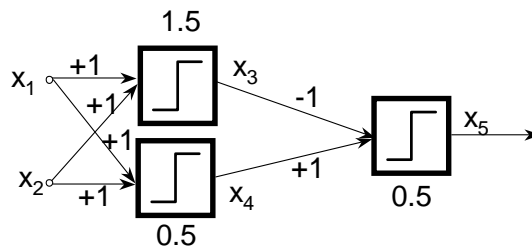
	<p>Soft Computing: Neural Networks</p> <h2>Enter the Dark Ages of NNs</h2> 
	<p><b>...which (together with a lack or proper training techniques for multi-layer perceptrons) all but killed interest in neural nets in the 70s and early 80s.</b></p>

29



## Multilayer Perceptrons

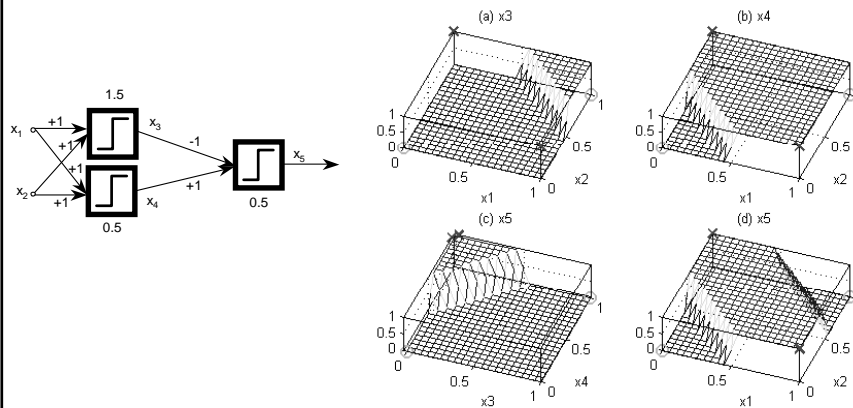
### Two-layer perceptron



31

## Two-Layer Perceptron: XOR

### Node output as surface of their two inputs



32

note location of "o" and "x"



Soft Computing: Neural Networks

## Multilayer Perceptrons (MLPs)

**Network architecture**

**Learning rule:**

- Steepest descent (Backprop)
- Conjugate gradient method
- All optim. methods using first derivative
- Derivative-free optim.

33

Soft Computing: Neural Networks

## Multi-Layer Perceptrons

**-Recall the output**

$$o = f\left(\sum_{i=1}^n w_i x_i - \theta\right)$$

**-and the squared error measure**

$$E_p = (t_p - o_p)^2 \quad \text{which is amended to} \quad {}^p E_k = \sum_{k=1}^n ({}^p t_k - {}^p o_k)^2$$

**-and the activation function**

$$f(x) = \frac{1}{1 + e^{-x}} \quad \text{or} \quad f(x) = \frac{1 - e^{-x}}{1 + e^{-x}} = \tanh\left(\frac{x}{2}\right) \quad \text{or} \quad f(x) = x$$



**then the learning rule for each node can be derived using the chain rule...**

34

Soft Computing: Neural Networks

## Multi-Layer Perceptrons

**-Recall the output**

$$o = f\left(\sum_{i=1}^n w_i x_i - \theta\right)$$

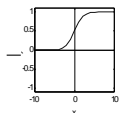
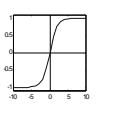
**-and the squared error measure**

$$E_p = (t_p - o_p)^2 \quad \text{which is amended to} \quad {}^p E_k = \sum_{k=1}^n ({}^p t_k - {}^p o_k)^2$$

**-and the activation function**

$$f(x) = \frac{1}{1 + e^{-x}} \quad \text{or} \quad f(x) = \frac{1 - e^{-x}}{1 + e^{-x}} = \tanh\left(\frac{x}{2}\right) \quad \text{or} \quad f(x) = x$$

**squashing functions**

35

Soft Computing: Neural Networks

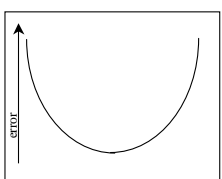
## Backpropagation

**make incremental change in the direction  $\partial E / \partial w$  to decrease the error.**

**The learning rule for each node can be derived using the chain rule...**

**...to propagate the error back through a multi-layer perceptron.**

$$\frac{\partial E}{\partial \text{parameters}}$$

$$\Delta w_{ki} = -\eta \sum_p \frac{\partial E_p}{\partial w_{ki}}$$


36

## Back-prop procedure

1. Initialize weights to small random values
2. Choose a pattern and apply it to input layer
3. Propagate the signal forward through the network
4. Compute the deltas for the output layer
5. Compute the deltas for the preceding layers by propagating the error backwards
6. Update all weights
7. Go back to step 2 and repeat for next pattern
8. Repeat until error rate is acceptable

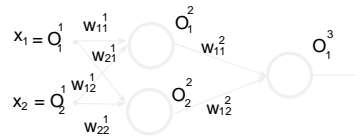
37

## Step 1

1. Initialize weights to small random values

Example:

Weight	Value
$w_{11}^1$	-0.4
$w_{12}^1$	-0.3
$w_{21}^1$	-0.1
$w_{22}^1$	-0.1
$w_{11}^2$	-0.6
$w_{12}^2$	0.3



Threshold	Value
$t_1^2$	0.2
$t_2^2$	0.5
$t_1^3$	-0.1

38

Soft Computing: Neural Networks

## Step 2

2. Choose a pattern and apply it to input layer

x1	x1	Target output
0	0	0
0	1	1
1	0	1
1	1	0

39

Soft Computing: Neural Networks

## Step 3

Propagate signal forward through the network

$$O_i^m = g \left( \sum_j W_{ij}^m O_j^{m-1} \right)$$

until all outputs have been calculated  
 For m=0 (input layer), the output is the pattern.  
**Example:**

$t_1^2 = 0.2$   
 $x_1 = O_1^1 = 0$     $W_{11}^1 = -0.4$     $O_1^2 = 1 / (1 + \exp(-2 * \beta * (0 * (-0.4) + 1 * (-0.3) + 0.2)))$   
 $x_2 = O_2^1 = 1$     $W_{21}^1 = -0.3$

40

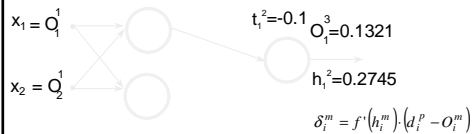
## Step 4

Compute the deltas for the output layer

$$\delta_i^M = g'(h_i^M) [d_i^p - O_i^m]$$

by comparing the actual output **O**  
with the target output **t**  
for the pattern **p** considered

**Example:**



41

## Step 5

Compute the deltas for the preceding layers by  
propagating the errors backwards

$$\delta_i^{m-1} = g'(h_i^{m-1}) \sum_j^n w_{ij}^m \delta_j^m$$

for  $m=M, M-1, M-2, \dots$

until a delta has been calculated for every unit

42

Soft Computing: Neural Networks

## Step 6

Use

$$\Delta w_{ij}^m = \eta \delta_i^m O_j^{m-1}$$

to update all connections to

$$w_{ij}^{new} = w_{ij}^{old} + \Delta w_{ij}$$

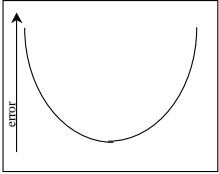
43

Soft Computing: Neural Networks

## Step Size, Initial Weights

Step size:

- too big
- too small
- variable:



- compute error
- backpropagate
- compute error again
- if error bigger, reduce step size (0.5)
- otherwise, increase a little (1.1)

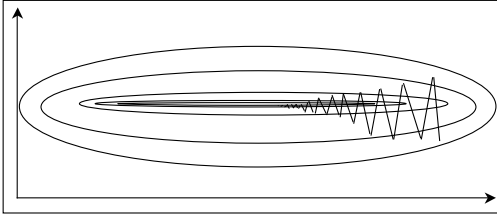
Initial weights: randomize ( $\neq 0$ )

44

Soft Computing: Neural Networks

## Momentum

If error minimum in long narrow valley, then updating can happen to zig-zag down the valley



smoothes weight updating  
can speed learning up

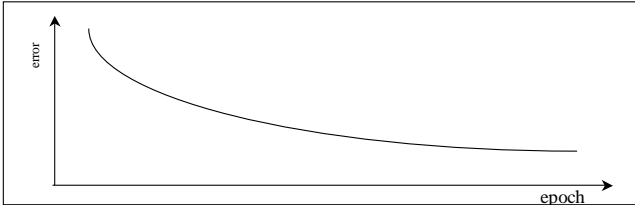
$$\Delta w = -\eta \nabla_w E + \alpha \Delta w_{prev}$$

45

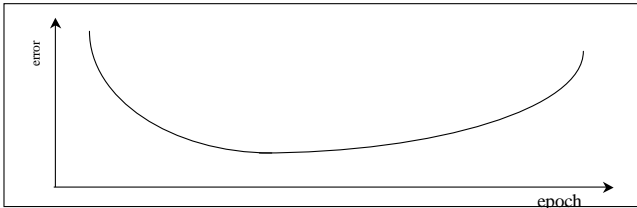
Soft Computing: Neural Networks

## Overfitting

Error on learning cases



Error on validation cases

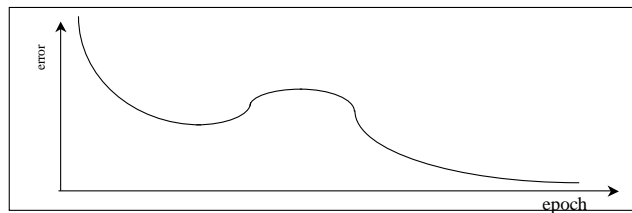


trained things that are accidental and unimportant

46

## Local Minima

There is no guarantee that the algorithm converges to a global minimum



- check with different initial conditions (different weights, etc.)
- perturb the system (data) with noise to improve result

47

## Architectures and other Techniques

### Normalize weights

move weights same Euclidean distance  
each epoch

### Data scaling

Input scaling: allows weights to have  
same order of magnitude

Output scaling: let target go between  $\pm 0.9$   
to avoid saturation

What number of nodes per layer?

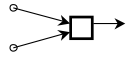
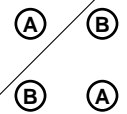
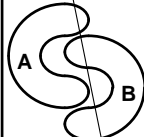
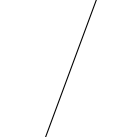
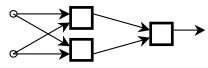
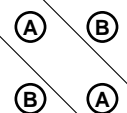
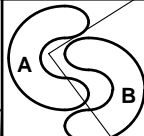
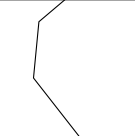
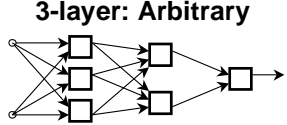
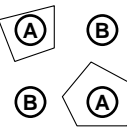
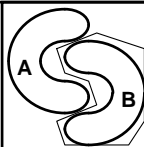
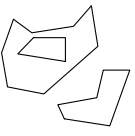
How many layers?

48



Soft Computing: Neural Networks

## MLP Decision Boundaries

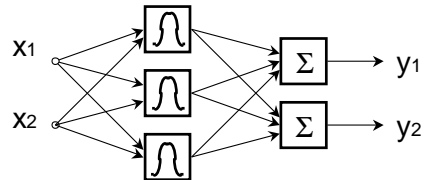
	XOR	Interwined	General
<b>1-layer: Half planes</b> 			
<b>2-layer: Convex</b> 			
<b>3-layer: Arbitrary</b> 			

49

Soft Computing: Neural Networks

## Radial Basis Function (RBF) Networks

**Network architecture**



**Each node is described by a bell shaped function**

$$a_i = \exp\left(-\frac{\|x - c_i\|^2}{2\sigma_i^2}\right)$$

**where**

**$c_i$  is the center of the curve**

50

Soft Computing: Neural Networks

## RBF

---

**Output:**  
 weighted sum  
 weighted average  
 linear combination

**Location of Center:**  
 Use (fuzzy) k-means clustering

**Size of Variance:**  
 Use knn-classifier and take average distance

51

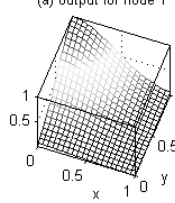
Soft Computing: Neural Networks

## XOR, revisited

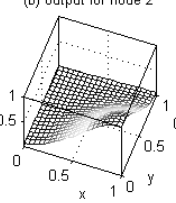
---

x	y	output
0	0	0
0	1	1
1	0	1
1	1	0

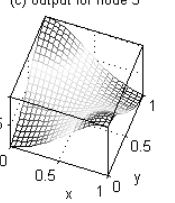
(a) output for node 1

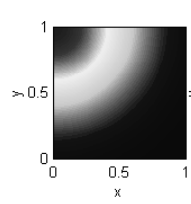
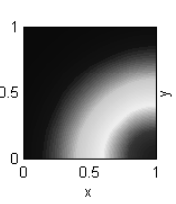
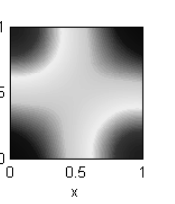


(b) output for node 2



(c) output for node 3



52

Soft Computing: Neural Networks

## RBF and FIS

Consider the radial basis functions:

and a linear combination of the output variables

$$y_i = \vec{a}_i \vec{o} + b_i$$

then the response is equivalent to ...

53

Soft Computing: Neural Networks

## Modular Networks

Task decomposition  
Local Experts  
Fuse information

54

Soft Computing: Neural Networks

last slide

55

The slide features a title bar at the top containing the text "Soft Computing: Neural Networks". A horizontal arrow points from the left side of the slide towards the right. The text "last slide" is centered in the main area. A small number "55" is located in the bottom-left corner of the slide frame.