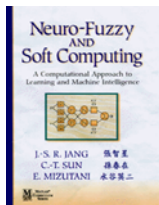


## Chapter 9: Supervised Learning Neural Networks

- Introduction (9.1)
- Perceptrons (9.2)
- Adaline (9.3)
- Backpropagation Multilayer Perceptrons (9.4)
- Radial Basis Function Networks (9.5)



Jyh-Shing Roger Jang et al., *Neuro-Fuzzy and Soft Computing: A Computational Approach to Learning and Machine Intelligence*, First Edition, Prentice Hall, 1997

### Introduction (9.1)

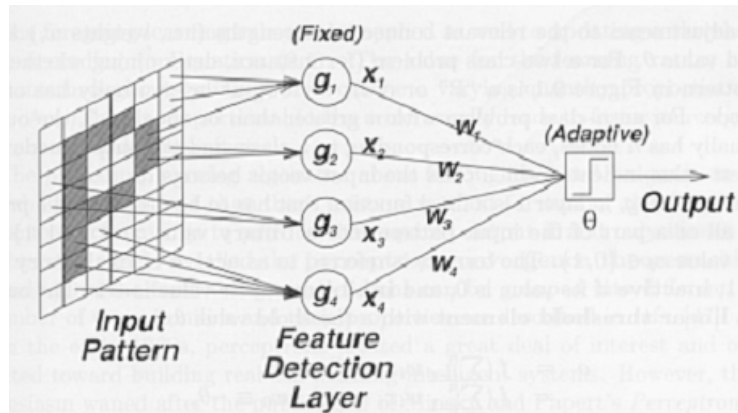
- Artificial Neural Networks (NN) have been studied since 1950
- Minsky & Papert in their report of perceptron (Rosenblatt) expressed pessimism over multilayer systems, the interest in NN dwindled in the 1970's
- The work of Rumelhart, Hinton, Williams & others, in learning algorithms created a resurgence of the lost interest in the field of NN

## Introduction (9.1) (cont.)

- Several NN have been proposed & investigated in recent years
  - Supervised versus unsupervised
  - Architectures (feedforward vs. recurrent)
  - Implementation (software vs. hardware)
  - Operations (biologically inspired vs. psychologically inspired)
- In this chapter, we will focus on modeling problems with desired input-output data set, so the resulting networks must have adjustable parameters that are updated by a supervised learning rule

## Perceptrons (9.2)

- Architecture & learning rule
  - The perceptron was derived from a biological brain neuron model introduced by Mc Culloch & Pitts in 1943
  - Rosenblatt designed the perceptron with a view toward explaining & modeling pattern recognition abilities of biological visual systems
  - The following figure illustrate a two-class problem that consists of determining whether the input pattern is a “p” or not



The perceptron

## Perceptrons (9.2) (cont.)

$$\begin{aligned}
 f\left(\sum_{i=1}^{i=n} w_i x_i - \theta\right) &= \text{Output} \\
 &= f\left(\sum_{i=1}^{i=n} w_i x_i + w_0\right), w_0 \equiv -\theta \\
 &= f\left(\sum_{i=1}^{i=n} w_i x_i\right), x_0 \equiv 1
 \end{aligned}$$

- A signal  $x_i$  is binary, it could be active (or excitatory) if its value is 1, inactive if its value is 0 and inhibitory if its value is -1

## Perceptrons (9.2) (cont.)

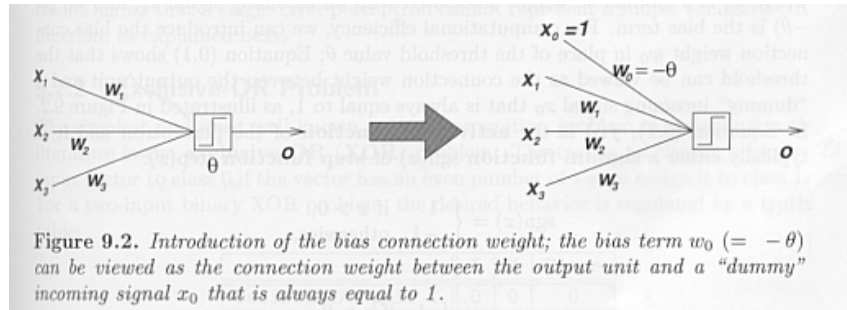
- Architecture & learning rule (cont.)
  - The output unit is a linear threshold element with a threshold value  $\theta$
  - $w_i$  is a modifiable weight associated with an incoming signal  $x_i$
  - The threshold  $\theta = w_0$  can be viewed as the connection weight between the output unit & a dummy incoming signal  $x_0 = 1$

## Perceptrons (9.2) (cont.)

- Architecture & learning rule (cont.)
  - $f(\cdot)$  is the activation function of the perceptron & is either a signum function  $\text{sgn}(x)$  or a step function  $\text{step}(x)$

$$\text{sgn}(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{otherwise} \end{cases}$$

$$\text{step}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$



## Perceptrons (9.2) (cont.)

- Algorithm (Single-layer perceptron)

1. Start with a set of random connection weights
2. Select an input vector  $x$  from the training data set  
If the perceptron provides a wrong response then modify all connection weights  $w_i$  to  $w_i = \eta t_i x_i$   
where:  $t_i$  is a target output  
 $\eta$  is a learning rate
1. Test the weight convergence: if converge stop else go to 1

This learning algorithm is based on gradient descent

## Perceptrons (cont.)

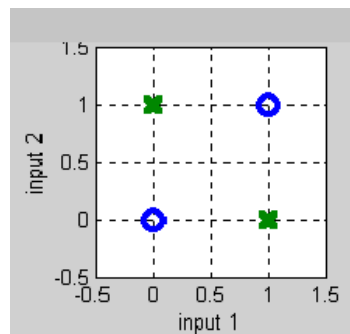
- Exclusive-OR problem (XOR)

Goal: classify a binary input vector to class 0 if the vector has an even number of 1's, otherwise assign it to class 1

	X	Y	Class
Desired i/o pair 1	0	0	0
Desired i/o pair 2	0	1	1
Desired i/o pair 3	1	0	1
Desired i/o pair 4	1	1	0

## Perceptrons (9.2) (cont.)

- Exclusive-OR problem (XOR) (cont.)
  - From the following figure, we can say that the XOR problem is not linearly separable



### Perceptrons (9.2) (cont.)

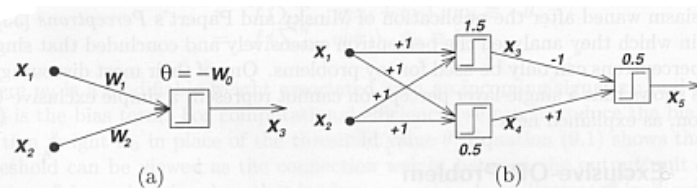
- Using a single-layer perceptron and the step function to solve this problem requires satisfying the following inequalities

$$\begin{aligned}
 0 * w_1 + 0 * w_2 + w_0 &\leq 0 \Leftrightarrow w_0 \leq 0 \\
 0 * w_1 + 1 * w_2 + w_0 &> 0 \Leftrightarrow w_0 > -w_2 \\
 1 * w_1 + 0 * w_2 + w_0 &> 0 \Leftrightarrow w_0 \leq -w_1 \\
 1 * w_1 + 1 * w_2 + w_0 &\leq 0 \Leftrightarrow w_0 \leq -w_1 - w_2
 \end{aligned}$$

This self of inequalities is self-contradictory  
 $\Rightarrow$  Minsky & Pappert criticism of perceptron was justified in part by this result!

### Perceptrons (9.2) (cont.)

- The XOR problem can be solved using a two-layer perceptron illustrated by the following figure



**Figure 9.4.** Perceptrons for the two-input exclusive-OR problem: (a) the single-layer perceptron, and (b) the two-layer perceptron. Both use the step function as the activation function for each node.

## Perceptrons (9.2) (cont.)

$(x_1 = 0, x_2 = 0 \Rightarrow 0)$

results at the hidden layer

$$0 * (+1) + 0 * (+1) = 0 < 1.5 \Rightarrow x_3 = 0$$

$$0 * (+1) + 0 * (+1) = 0 < 0.5 \Rightarrow x_4 = 0$$

results at the output layer

$$0 * (-1) + 0 * (+1) = 0 < 0.5 \Rightarrow x_5 = \text{output} = 0$$

$(x_1 = 0, x_2 = 1 \Rightarrow 1)$

results at the hidden layer

$$0 * (+1) + 1 * (+1) = 1 < 1.5 \Rightarrow x_3 = 0$$

$$1 * (+1) + 0 * (+1) = 1 > 0.5 \Rightarrow x_4 = 1$$

results at the output layer

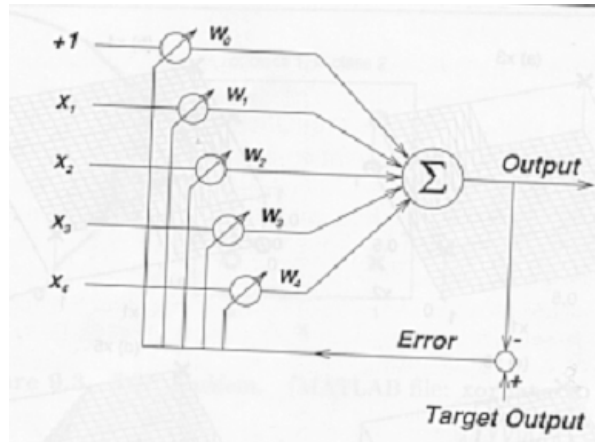
$$0 * (-1) + 1 * (+1) = +1 > 0.5 \Rightarrow x_5 = \text{output} = 1$$

In summary, multilayer perceptrons can solve nonlinearly separable problems and are thus more powerful than the single-layer perceptrons

## ADALINE (9.3)

- Developed by Widrow & Hoff, this model represents a classical example of the simplest intelligent self-learning system that can adapt itself to achieve a given modeling task





Adaline (Adaptive linear element)

### ADALINE (9.3) (cont.)

$$\text{output} = \sum_{i=1}^n w_i x_i + w_0$$

- One possible implementation of ADALINE is the following:
  - The input signals are voltages
  - The weights  $w_i$  are conductances of controllable resistors
  - The network output is the summation of the currents caused by the input voltages
  - The problem consists of finding a suitable set of conductances such that the input-output behavior of ADALINE is close to a set of desired input-output data points

## ADALINE (9.3) (cont.)

- The ADALINE model can be solved using a linear least-square method,  $(n + 1)$  linear parameters in order to minimize the error

$$\sum_{p=1}^m (t_p - o_p)^2 \quad (m \text{ training data})$$

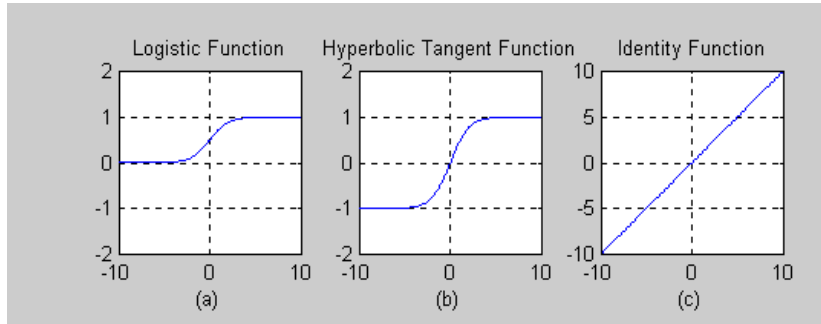
- However, since this method can be slow (requires too many calculations!) if  $n$  is large, therefore Widrow & Hoff fell back on gradient descent

$$\text{if } E_p = (t_p - o_p)^2 \Rightarrow \frac{\partial E_p}{\partial w_i} = -2(t_p - o_p) * x_i \quad \begin{array}{l} \text{(Least-mean square} \\ \text{(LMS) learning} \\ \text{procedure)} \end{array}$$

$$\text{which provides : } w_{\text{next}}^p = w_{\text{now}}^p + \eta \underbrace{(t_p - o_p) * x_i}_g$$

## Backpropagation Multilayer Perceptrons (9.4)

- There was a change in 1985 of the reformulation of the backpropagation training method by Rumelhart
- The sigum and the step functions are not differentiable, the use of logistic (hyperbolic) functions contribute for a better learning scheme
  - Logistic:  $f(x) = 1 / (1 + e^{-x})$
  - Hyperbolic tangent:  $f(x) = \tanh(x/2) = (1 - e^{-x}) / (1 + e^{-x})$
  - Identity:  $f(x) = x$
- The sigum function is approximated by the hyperbolic tangent function & the step function is approximated by the logistic function



Activation functions for backpropagation MLPs

## Backpropagation Multilayer Perceptrons (9.4) (cont.)

- Backpropagation learning rule

Principle:

The net input  $\bar{x}$  of a node is defined as the weighted sum of the incoming signals plus a bias term:

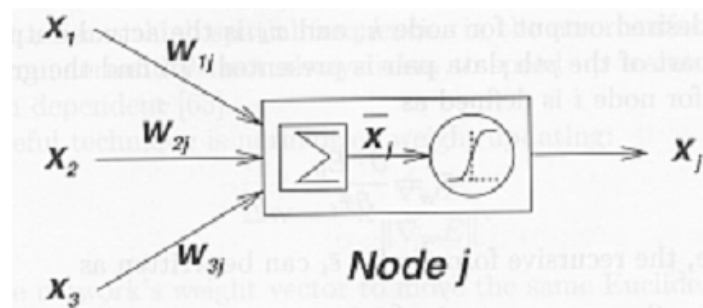
$$\bar{x}_j = \sum_i w_{ij} x_i + w_j \quad (\text{Logistic function})$$

$$x_j = f(\bar{x}_j) = \frac{1}{1 + \exp(-\bar{x}_j)}$$

Where:  $x_i$  = output of node  $i$  at any of the previous layers

$w_{ij}$  = weight associated with the link connecting nodes  $i$  &  $j$

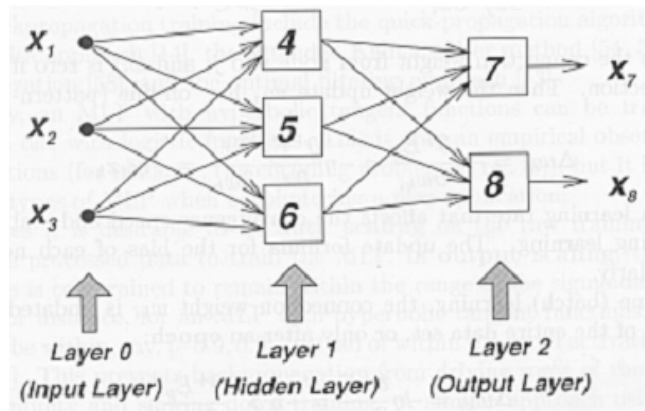
$w_j$  = bias of node  $j$



Node j of a backpropagation MLP

### Backpropagation Multilayer Perceptrons (9.4) (cont.)

- The following figure shows a two-layer backpropagation MLP with 3 inputs in the input layer, 3 neurons in the hidden layer & 2 output neurons in the output layer



A 3-3-2 backpropagation MLP

### Backpropagation Multilayer Perceptrons (9.4) (cont.)

- The square error measure for the  $p$ -th input-output pair is defined as:

$$E_p = \sum_k (d_k - x_k)^2$$

where:  $d_k$  = desired output for node  $k$

$x_k$  = actual output for node  $k$  when the  $p$ -th data pair is presented

- Since it is a propagation scheme, an error term  $\bar{\epsilon}_i$  for node  $i$  is needed:

$$\bar{\epsilon}_i = \frac{\partial^+ E_p}{\partial \bar{x}_i}$$

Using a chain rule derivation, we obtain:

$$w_{k,i}^{\text{next}} = \underbrace{w_{k,i}^{\text{now}} - \eta \nabla_{w_{k,i}} E}_{\text{steepest descent}} = w_{k,i}^{\text{now}} - \eta \bar{\epsilon}_i x_k$$

## Backpropagation Multilayer Perceptrons (9.4) (cont.)

- Algorithm (Stochastic backpropagation)

```

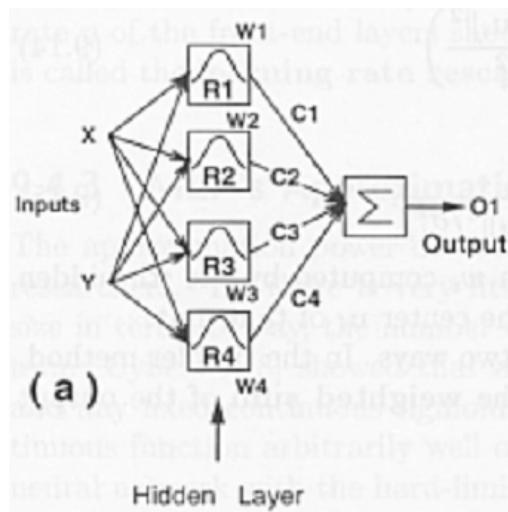
Begin initialize number-of-hidden-units,
w, criterion  $\theta, \eta, m$  (training data size)
  Do  $m \leftarrow m + 1$ 
       $x^m \leftarrow$  randomly chosen pattern
           $w_{k,I} \leftarrow w_{k,I} + \eta \bar{\epsilon}_i x_k$ 
      Until  $||\nabla E(w)|| < \theta$ 
  Return  $w$ 
End.

```

## Radial Basis Function Networks (9.5)

- Architectures & Learning Methods

- Inspired by research in regions of the cerebral cortex & the visual cortex, RBFNs have been proposed by Moody & Darken in 1988 as a supervised learning neural networks
- The activation level of the  $i$ th receptive field unit is:
 
$$w_i = R_i(x) = R_i(||x - u_i|| / \sigma_i), i = 1, 2, \dots, H$$
  - $x$  is a multidimensional input vector
  - $u_i$  is a vector with same dimension as  $x$
  - $H$  is the number of radial basis functions called also receptive field units
  - $R_i(\cdot)$  is the  $i$ th radial basis function with a single maximum at the origin



Single-output RBFN that uses weighted sum

## Radial Basis Function Networks (9.5) (cont.)

- Architectures & Learning Methods (cont.)

- $R_i(\cdot)$  is either a Gaussian function

$$R_i^G(x) = \exp\left(-\frac{\|x - u_i\|^2}{2\sigma_i^2}\right)$$

or a logistic function

$$R_i^L(x) = \frac{1}{1 + \exp[\|x - u_i\|^2 / \sigma_i^2]}$$

if  $x = u_i \Rightarrow R_i^G = 1$  (Maximum) &  $R_i^L = 1/2$  (Maximum)

## Radial Basis Function Networks (9.5) (cont.)

- Architectures & Learning Methods (cont.)

– The output of an RBFN

- $$d(x) = \sum_{i=1}^{i=H} c_i w_i = \sum_{i=1}^{i=H} c_i R_i(x) \quad (\text{weighted sum})$$

where  $c_i$  = output value associated with the  $i$ th receptive field

- $$d(x) = \frac{\sum_{i=1}^{i=H} c_i w_i}{\sum_{i=1}^{i=H} w_i} = \frac{\sum_{i=1}^{i=H} c_i R_i(x)}{\sum_{i=1}^{i=H} R_i(x)} \quad (\text{weighted average})$$

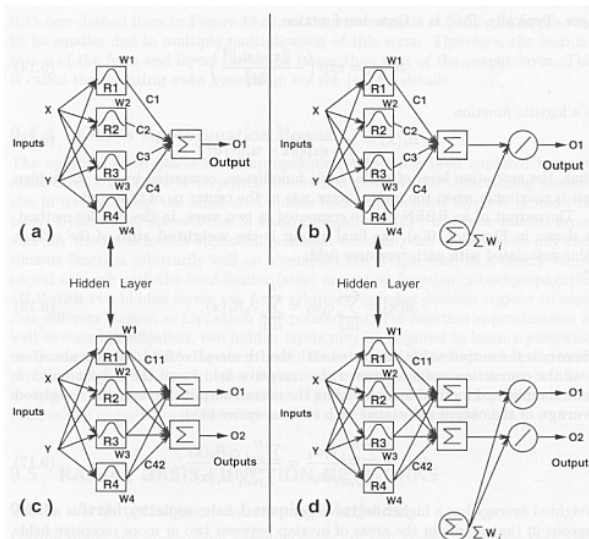


Figure 9.10. Four RBFNs that possess four basis functions: (a) single-output RBFN that uses weighted sum; (b) single-output RBFN that uses weighted average; (c) two-output RBFN that uses weighted sum; (d) two-output RBFN that uses weighted average. The network in (d) is equivalent to Figure 13.3 (upper right). [Note that in (b) and (d), four connections to the lower summation unit are omitted for simplicity.]



## Radial Basis Function Networks (9.5) (cont.)

- Architectures & Learning Methods (cont.)
  - Moody-Darke's RBFN may be extended by assigning a linear function to the output function of each receptive field
 
$$\mathbf{c}_i = \mathbf{a}_i^T \mathbf{x} + \mathbf{b}$$

( $\mathbf{a}_i$  is a parameter vector &  $\mathbf{b}_i$  is a scalar parameter)
  - Supervised adjustments of the center & shape of the receptive field (or radial basis) functions may improve RBFNs approximation capacity
  - Several learning algorithms have been proposed to identify the parameters ( $u_i$ ,  $\sigma_i$  &  $c_i$ ) of an RBFN

## Radial Basis Function Networks (9.5) (cont.)

- Functional Equivalence to FIS

$$\mathbf{c}_i = \vec{\mathbf{a}}_i \cdot \vec{\mathbf{x}} + \mathbf{b}_i$$

- The extended RBFN response given by the weighted sum or the weighted average is identical to the response produced by the first-order Sugeno fuzzy inference system provided that the membership functions, the radial basis function are chosen correctly

$$\mathbf{d}(\mathbf{x}) = \sum_{i=1}^{i=H} (\vec{\mathbf{a}}_i \vec{\mathbf{x}} + \mathbf{b}_i) \mathbf{w}_i(\mathbf{x}) = \sum_{i=1}^{i=H} (\vec{\mathbf{u}}_i \vec{\mathbf{x}} + \mathbf{v}_i)$$

$$\text{where : } \vec{\mathbf{u}}_i = [\mathbf{u}_i^1, \mathbf{u}_i^2, \dots, \mathbf{u}_i^m]^T, \vec{\mathbf{x}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]^T$$

## Radial Basis Function Networks (9.5) (cont.)

- Functional Equivalence to FIS (CONT.)
  - While the RBFN consists of radial basis functions, the FIS comprises a certain number of membership functions
  - The FIS & the RBFN were developed on different bases, they are rooted in the same soil

## Radial Basis Function Networks (9.5)(cont.)

- Functional Equivalence to FIS (cont.)
  - Condition of equivalence between FIS & RBFN
    - RBFN & FIS use both of them the same aggregation method (weighted average & weighted sum)
    - The number of receptive field units in RBFN is equal to the number of fuzzy if-then rules in the FIS
    - Each radial basis function of the RBFN is equal to a multidimensional composite MF of the premise part of a fuzzy rule in the FIS
    - Corresponding radial basis & fuzzy rule should have the same response function

## Radial Basis Function Networks (9.5) (cont.)

- Interpolation & approximation RBFN
  - The interpolation case: each RBF is assigned to each training pattern

Goal: Estimate a function  $d(\cdot)$  that yields exact desired outputs for all training data

- Our goal consists of finding  $c_i$  ( $i = 1, 2, \dots, n$ ) ( $n = H$ ) such that  $d(x_i) = o_i =$  desired output

since  $w_i = R_i (\|x - u_i\|) = \exp [- (x - u_i)^2 / (2 \sigma_i^2 )]$

Therefore, starting with  $x_i$  as centers for the RBFNs, we can write:

$$d(x) = \sum_{i=1}^n c_i \exp \left[ - \frac{(x - x_i)^2}{2\sigma_i^2} \right]$$

## Radial Basis Function Networks (9.5) (cont.)

- Interpolation & approximation RBFN (cont.)
  - **The interpolation case (cont.)**
    - For given  $\sigma_i$  ( $i = 1, \dots, n$ ), we obtain the following  $n$  simultaneous linear equations with respect to unknown weights  $c_i$  ( $i = 1, 2, \dots, n$ )

## Radial Basis Function Networks (9.5) (cont.)

– The interpolation case (cont.)

$$\text{First pattern } d(x_1) = \sum_{i=1}^n c_i \exp \left[ -\frac{(x_1 - x_i)^2}{2\sigma_i^2} \right] = d_1$$

$$\text{Second pattern } d(x_2) = \sum_{i=1}^n c_i \exp \left[ -\frac{(x_2 - x_i)^2}{2\sigma_i^2} \right] = d_2$$

⋮

$$\text{nth pattern } d(x_n) = \sum_{i=1}^n c_i \exp \left[ -\frac{(x_n - x_i)^2}{2\sigma_i^2} \right] = d_n$$

$$\mathbf{D} = \mathbf{GC} \quad \text{where } \mathbf{D} = [d_1, d_2, \dots, d_n]^T, \mathbf{C} = [c_1, c_2, \dots, c_n]^T, \\ \text{and } \mathbf{G} = \text{exp onential function values}$$

⇓

$$\mathbf{C} = \mathbf{G}^{-1}\mathbf{D} \quad (\text{optimal weights if } \mathbf{G} \text{ is nonsingular})$$

## Radial Basis Function Networks (9.5) (cont.)

- Interpolation & approximation RBFN (cont.)

– **Approximation RBFN**

- This corresponds to the case when there are fewer basis functions than there are available training patterns
- In this case, the matrix  $\mathbf{G}$  is rectangular & the least square methods are commonly used in order to find the vector  $\mathbf{C}$