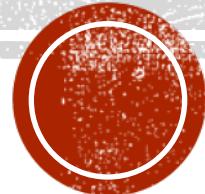


# **COMP-447: Neural Networks and Deep Learning**

## **LECTURE 6**

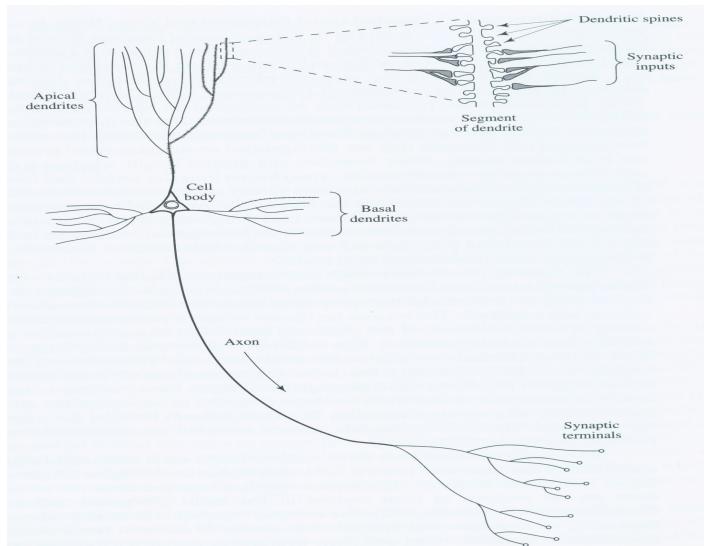
### **Radial Basis Function Networks**

Dr Michalis Agathocleous



# PREVIOUS LECTURE

- Temporality
- Recurrent Neural Networks



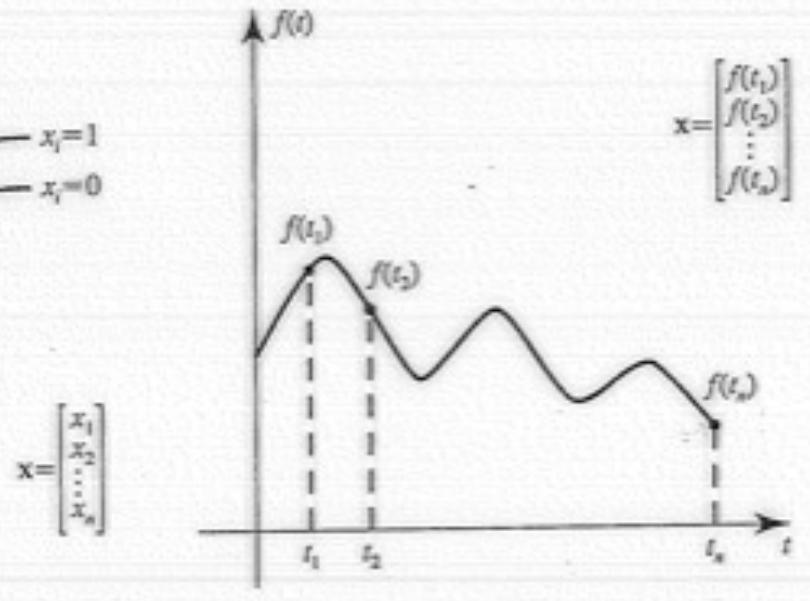
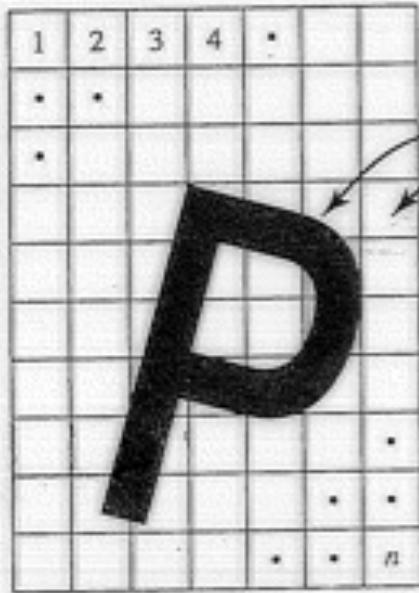
# SEQUENTIAL DATA

- **Sequential data (or time series) refers to data that appear in a specific order.**
  - The order defines a time axis, that differentiates this data from other cases we have seen so far
- Examples
  - The price of a stock (or of many stocks) over time
  - Environmental data (pressure, temperature, precipitation etc) over time
  - The sequence of queries in a search engine, or the frequency of a query over time
  - The words in a document as they appear in order
  - A DNA sequence of nucleotides
  - Event occurrences in a log over time
  - Etc...
- **Time series: usually we assume that we have a vector of numeric values that change over time.**

# TIME REPRESENTATION IN NEURAL NETWORKS

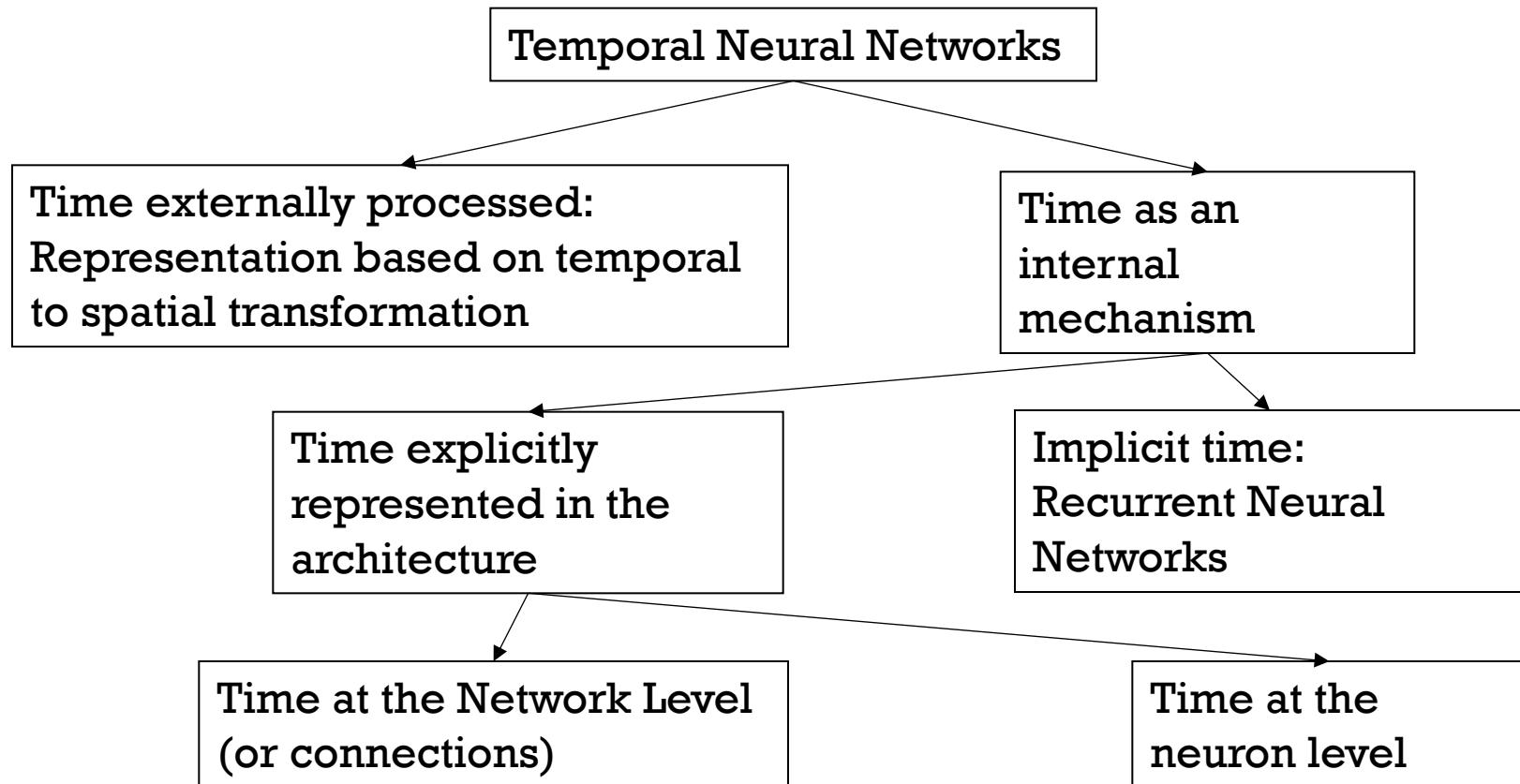
- Neural networks described so far deal with patterns which are **static** (lacking of temporal context)
- Real world living systems exists in **dynamic environments**; most of the problems **are intrinsically temporal in nature** and a number of practical **applications require neural networks that respond to a sequence of patterns**
- **A Neural Network should produce a particular response to a particular sequence of inputs**
- Example: Motion Detection, Speech recognition, robot trajectory generation, time series prediction, optical character recognition by sequence following etc.

# TIME REPRESENTATION IN NEURAL NETWORKS

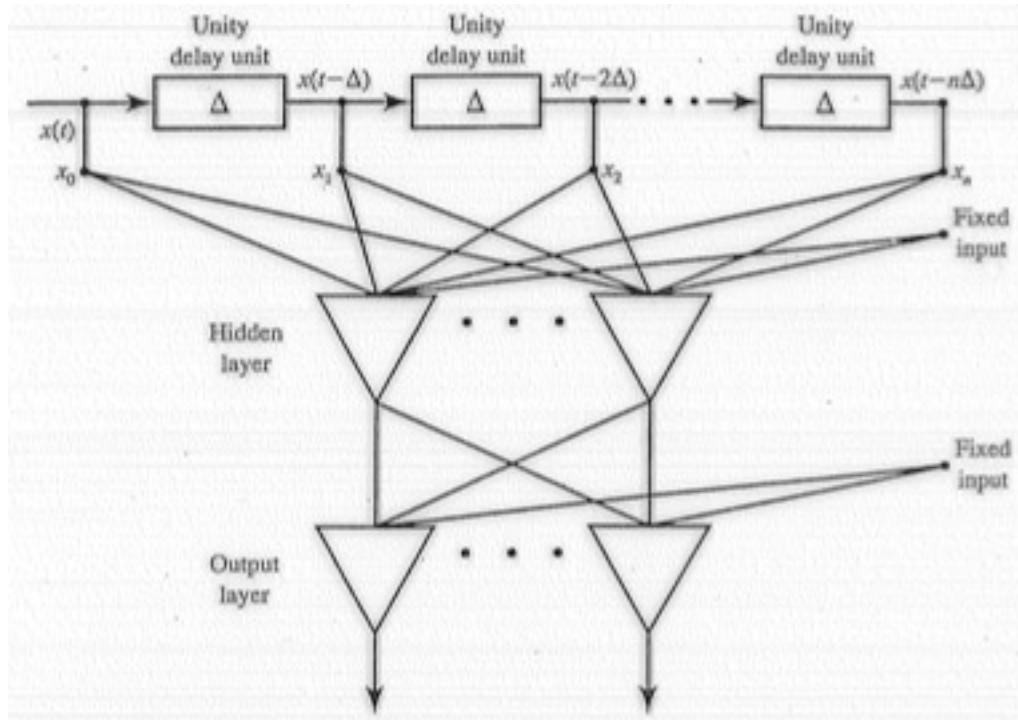


Two simple ways of coding patterns into pattern vectors: (a) special object and (b) temporal objects (waveforms)

# TEMPORAL NEURAL NETWORK CATEGORIES



# TIME DELAY NEURAL NETWORKS (TDNN)



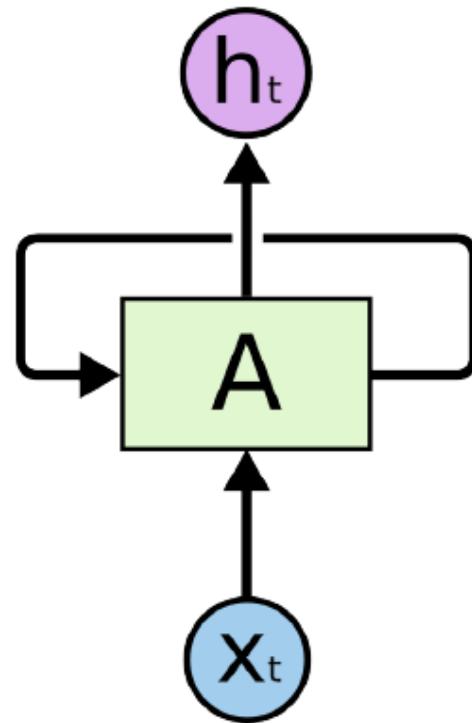
Inputs to the network  
are:  $X_i = X(t-i\Delta)$   
( $i=0, 1, \dots, n$ )

Note:  $\Delta$  is equal to the sampling period

# RECURRENT NEURAL NETWORKS

- **Current activation of the network can depend upon input history of a system and not just current input.**
- RNNs utilize internal memory and can be applied to problems involving **temporal context**. They calculate the **next state vector from the current state vector and the input vector**. The current output can either be calculated by using only the next state vector, or calculated from both the current state and input to the network.
- **Jordan** and **Elman** networks which are of recurrent backpropagation type, involve loops of activation flow.
- They are trained with a modified version of Backpropagation learning algorithm, called **Backpropagation Through Time learning algorithm**.

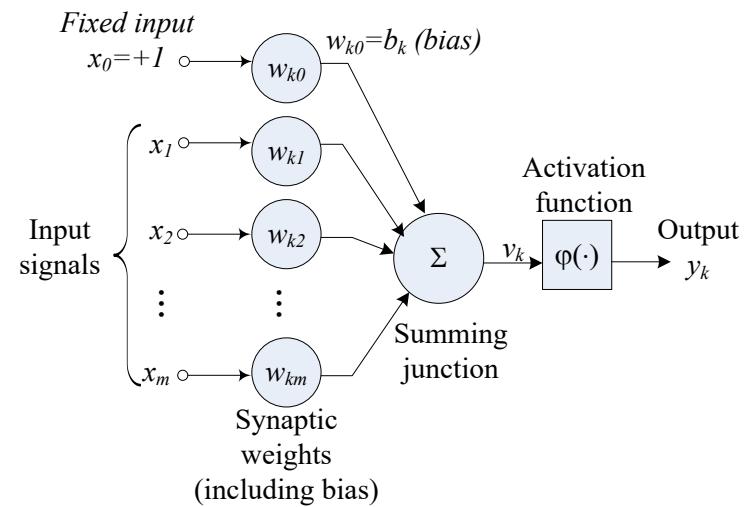
# BASIC STRUCTURE



- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# TODAY

- Radial Basis Function Networks



# INTRODUCTION

Much of these notes come from

1. Neural Networks and Learning Machines (3rd ed.) by Simon O. Haykin, Prentice Hall, Englewood Cliffs, NJ, 2008.
2. Neural Computing: an introduction by R. Beale and T. Jackson, 1990

# RADIAL BASIS FUNCTION NETWORKS

# INTERPOLATION AND APPROXIMATION

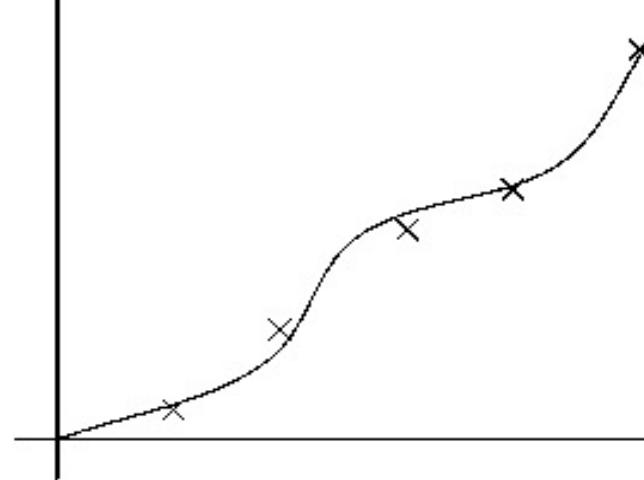
- RBF networks have their origins in techniques for performing exact *interpolation* of a set of input-output samples (What is the difference between regression and classification?)
- (exact) *interpolation* – after training NN gives the exactly correct outputs for all training data
- *approximation*– NN relates well the training input vectors to the desired outputs but does not necessary produce the exactly correct outputs
- Vacations example: how much you enjoy vacations as a function of their duration

Ex.	Duration (days)	Rating (1-10)
1	4	5
2	7	9
3	9	2
4	12	6

- We would like to create an interpolation and approximation networks that predict the ratings for other durations
- Given data is called sample input-output pairs (training data)

# APPROXIMATION THEORY

- Fitting of lines or surfaces to a set of points in a 2 or more dimensional space:
  - In 2-D space - lines are fitted to points.
  - In 3-D space - planes are fitted to points.
  - In n-D space - hyperplanes are fitted to points.
- Curve fitting is essentially the process that occurs as MLPs learn their training data.
- If there are K weights and biases in the network, the learning process thus traces out a point in  $(K+1)$  dimensional weight/error space.



# THE LAYERS OF AN RBF NETWORK

- Input layer

- the nodes distribute the network input vector to all the nodes in the hidden layer.

- Hidden layer

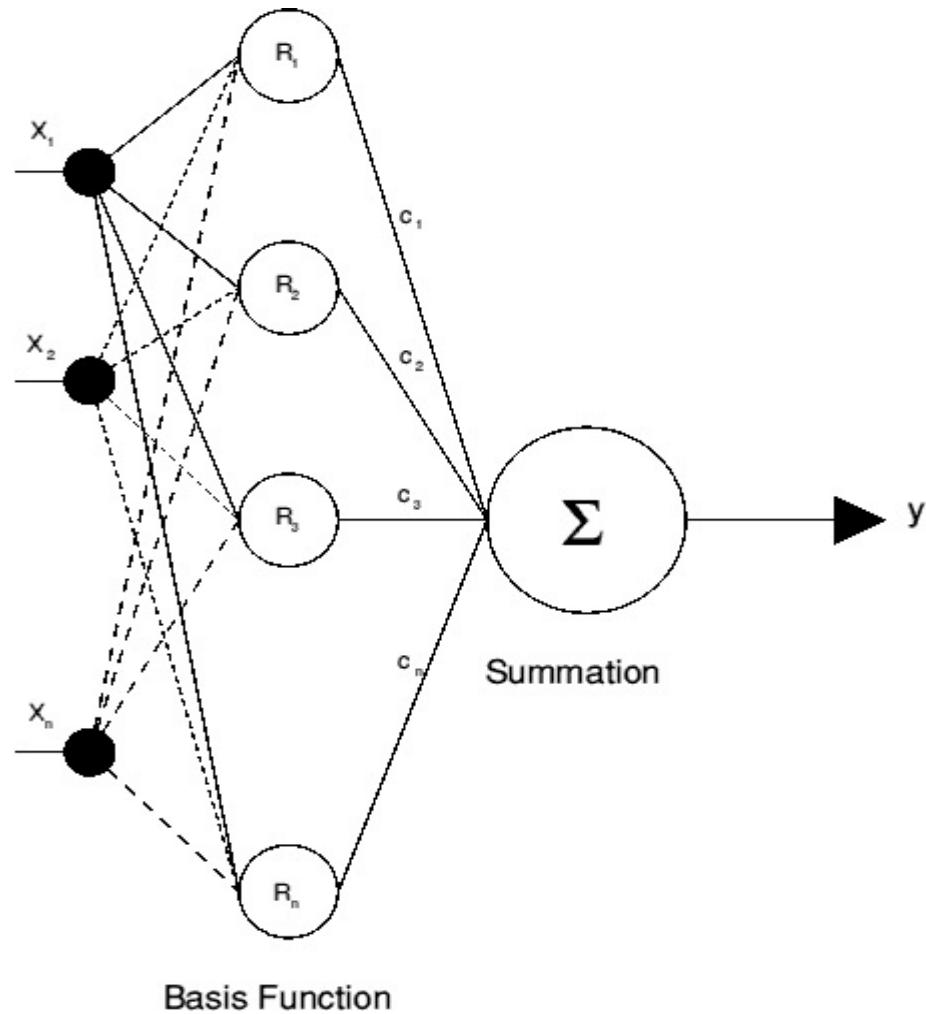
- each node contains a **centre** and a function called the **Basis Function**. It performs a non-linear mapping from the input space into a higher dimensional space in which the patterns become linearly separable.
- [https://www.youtube.com/watch?v=vMmG\\_7JcfIc](https://www.youtube.com/watch?v=vMmG_7JcfIc)
- <https://www.youtube.com/watch?v=3liCbRZPrZA>

- Output layer

- multiplies the output of each hidden layer node by a coefficient and sums the resulting values.

- Network output

- the result of the summation at the output layer.



Fully connected layer – all neurons of the previous layer are connected with all of the next layer

# MAIN IDEAS

Process performed in the hidden layer

- \* Idea: patterns in the input space form clusters
- If the centres of the clusters are known then the distance from the cluster centre can be measured
- \* the distance measure is made non-linear – so: if a pattern is in an area that is close to the cluster centre, it gives a value close to 1; beyond this area the value drops dramatically.
- \* Notion: this area is **radially symmetrical** around the cluster centre – so: The non-linear function becomes known as RBF

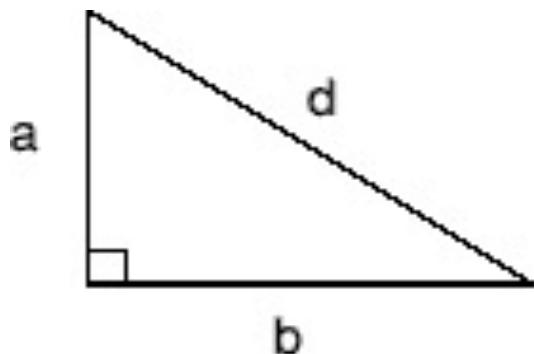
# WHAT DO RBFS CONSIST OF?

- A set of n-dimensional centres.
  - A set of scalar coefficients.
  - An input layer.
  - A hidden layer.
  - An output layer.
- 
- ***The Parameters***
    - ***A Centre ( $R$ )***
      - A vector of the same dimensionality as the network's input vector
      - This is associated with each node in the hidden layer.
    - ***A Coefficient ( $c$ )***
      - A scalar value which weights the connection between the hidden layer and the output layer.

# BASIS FUNCTIONS

- Equivalent to the activation function of an MLP.
- Takes in a scalar number and returns a scalar number which is the output of the hidden layer node.
- Commonly used basis functions:
  - Gaussian:  $\phi(r) = \exp(-r^2/2\sigma^2)$
  - *where  $\sigma$  represents the width of the Gaussian.*
- ***What happens in a hidden layer node?***
  - Euclidean distance between the input vector and the centre vector is found.
  - Distance produced is passed through the basis function to produce the output of the hidden layer node.

# EUCLIDEAN DISTANCE



$$d = \sqrt{a^2 + b^2}$$

For a n-Dimensional vector  $x = [x_1 \dots x_n]$ :

$$\|x\| = \text{Length}_{\text{EUC}} = \sqrt{\sum_{i=1}^N x_i^2}$$

# AN RBF HIDDEN NODE

- $\mathbf{x} = [x_1, x_2, \dots, x_N]$  = input vector
- $\mathbf{R} = [R_1, R_2, \dots, R_N]$  = centre vector

$$\|\mathbf{x} - \mathbf{R}\| = d_{\text{EUC}} = \sqrt{\sum_{i=1}^N (x_i - R_i)^2}$$

- *Euclidean distance between the vectors  $\mathbf{X}$  and  $\mathbf{R}$*
- Centre vector is subtracted from the input vector.
- The resulting distance is a scalar number which is used as the input to the basis function, to produce the hidden layer node's output ( $o$ ).

$$o = f(\|\mathbf{x} - \mathbf{R}\|)$$

# CALCULATING THE NETWORK OUTPUT

- The outputs of the hidden layer nodes are passed to the output layer. The output layer weights each hidden layer output with a coefficient and sums the resultant values.

$$y = \sum_{i=1}^m c_i o_i$$

- Where  $m$  is the number of nodes in the hidden layer,  $o_i$  represents the output of the  $i_{th}$  hidden layer node and  $c_i$  is the coefficient weighting the connection between the  $i_{th}$  hidden node and the output layer node.

- Output equation:

$$y = \sum_{i=1}^m c_i \phi(|x - R_i|)$$

- Where  $R_i$  is the centre belonging to hidden layer node  $i$ .

# RBFs AND PERCEPTRONS

- RBFs have similarities to perceptrons, but are more powerful than single layer perceptrons ...
  - ... As a RBF network can learn problems which a single layer perceptron cannot
    - i.e. 'Hard' problems which require a Multi-Layer Perceptron.
- Perceptrons weight and sum their inputs before passing them through an activation function:

$$y = \sigma \sum_1^N x_i w_i$$

- RBFs pass each input through a basis function before weighting and summing:

$$y = \sum_1^m c_i \phi(|x - R_i|)$$

# COMPARISON WITH RBFS AND MLPS

- **Differences:**

<u>Feature</u>	<u>MLP</u>	<u>RBF</u>
▪ <i>Variables:</i>	Weights & Thresholds	Centres & Coefficients
▪ <i>Speed of Training:</i>	Slow	Fast (er)
▪ <i>Computation:</i>	Computationally expensive	Less computation
▪ <i>Accuracy:</i>	Very good	Not as good
▪ <i>Mathematical Description</i>	Difficult	Easy

- **Similarities:**

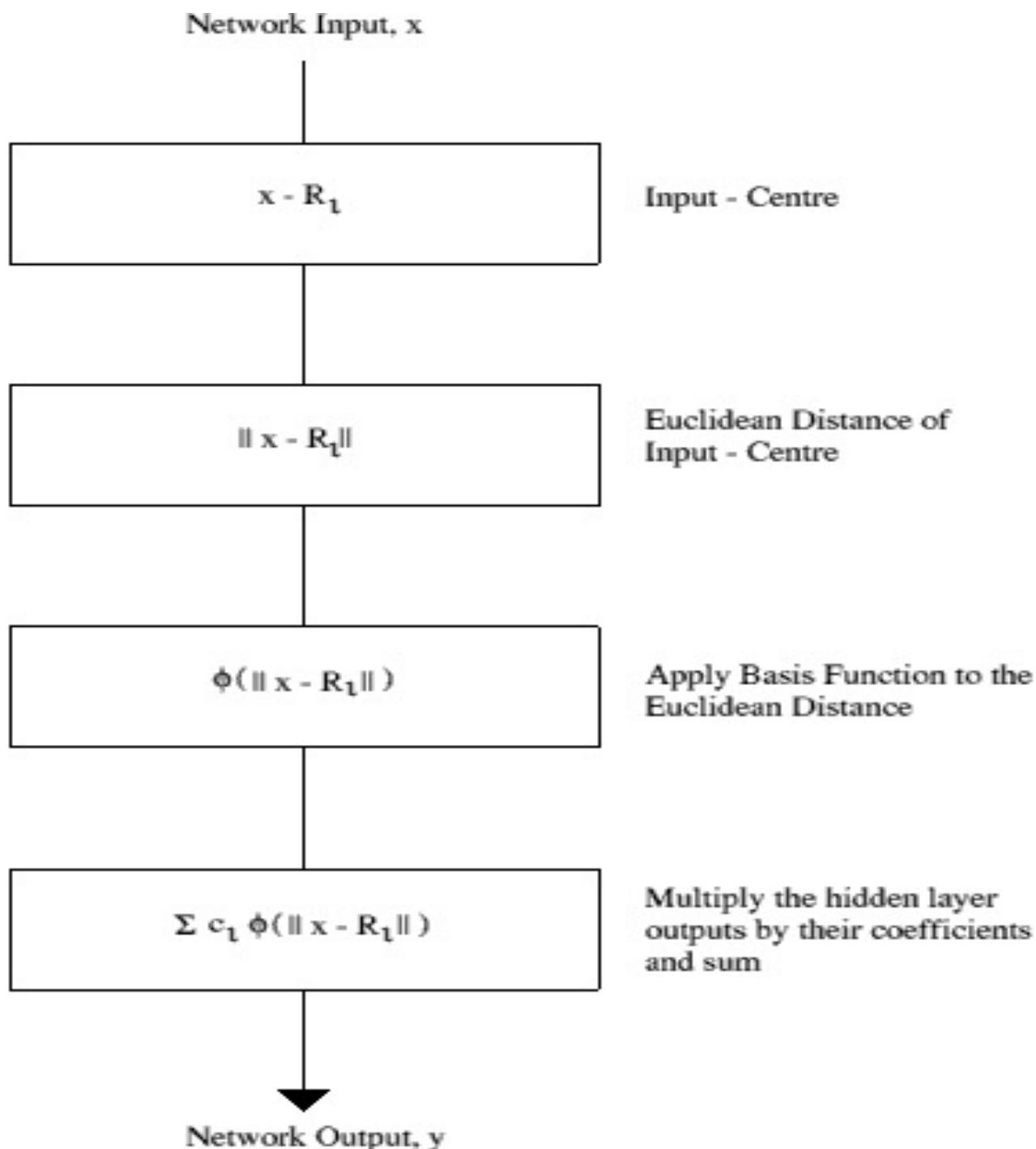
- Learn from a set of training data.
- Ability to generalise from training data.
- Store information in a distributed manner.
- Can be implemented using parallel processing.

# AN EXAMPLE: RBF TO EVALUATE XOR

- The Centres:

- The two centres are chosen as:
  - $R_1 = [1 \ 1]^T$
  - $R_2 = [0 \ 0]^T$

- The weights & bias:
  - $c_1 = c_2 = -1.97$
  - $b = +2.45$



# CALCULATING THE OUTPUT

- The input
  - Let  $X = [0 \ 1]$ ; For XOR the desired output is  $[+1]$
- Distance of  $X$  from  $R_1 [1 \ 1]$ 
  - $\text{SQRT } ((0 - 1)^2 + (1 - 1)^2) = +1$
  - Output  $R_1 = \text{Gaussian } (+1) = e^{-1} = 0.368$
- Distance of  $X$  from  $R_2 [0 \ 0]$ 
  - $\text{SQRT } ((0 - 0)^2 + (1 - 0)^2) = +1$
  - Output  $R_2 = \text{Gaussian } (+1) = e^{-1} = 0.368$
- Output =  $(0.368 * (-1.97)) + (0.368 * (-1.97)) + (+2.45)$   
 $= 1.00008$

# DESIGNING AN RBF

- **How to choose centres**

- Centres chosen to represent input training data set.
- The correct choice of centres is critical for good RBF performance.

- **The number of centres**

- Too many or too few leads to inaccuracy.
- Currently there are no rigorous formal methods to determine the optimum number of centres.

- **Methods for choosing centres:**

- *Simple distribution*
  - Uniform distribution over data space
  - Gaussian distribution over data space
- *Distribution related to the data distribution*
  - Eg. Use of clustering algorithms

# TRAINING THE RBF

$$y = \sum_{i=1}^n c_i \phi(\|x - R_i\|)$$

- Before training commences, the centres,  $R$ , are chosen. The only parameters which need to be trained are the coefficients  $c$ . All the other RBF values are known.
- The input and output training sets provide pairs of values for  $x$  and  $y$ . For each input/output,  $x/y$ , pair, an equation is produced:

$$y = c_1 \phi(x, R_1) + c_2 \phi(x, R_2) + \dots + c_n \phi(x, R_n)$$

- Where:  $\phi(x, R) = \phi(\|x - R\|)$

# CALCULATING THE CENTRE COEFFICIENTS

- As there is usually more than one pattern in the training set a group of equations is produced:

$$y_1 = c_1 \phi(x_1, R_1) + c_2 \phi(x_1, R_2) + \dots + c_n \phi(x_1, R_n)$$
$$y_2 = c_1 \phi(x_2, R_1) + c_2 \phi(x_2, R_2) + \dots + c_n \phi(x_2, R_n)$$
$$y_m = c_1 \phi(x_m, R_1) + c_2 \phi(x_m, R_2) + \dots + c_n \phi(x_m, R_n)$$

- There are usually more training pairs than there are unknown coefficients, so we have an over determined set of equations.

# ... IN VECTOR NOTATION

- Each equation can be expressed in vector notation, using vector dot product to represent the summation, producing the following set of equations:

$$y_1 = \phi(x_1, \underline{R}) \cdot \underline{c}$$

$$y_2 = \phi(x_2, \underline{R}) \cdot \underline{c}$$

$$y_m = \phi(x_m, \underline{R}) \cdot \underline{c}$$

- These equations can be represented more simply in matrix form by:

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{bmatrix} = \begin{bmatrix} \phi(x_1, \underline{R}) \\ \phi(x_2, \underline{R}) \\ \dots \\ \phi(x_m, \underline{R}) \end{bmatrix} \cdot \begin{bmatrix} c_1 \\ c_2 \\ \dots \\ c_n \end{bmatrix}$$
$$\underline{y} = \underline{A} \times \underline{c}$$

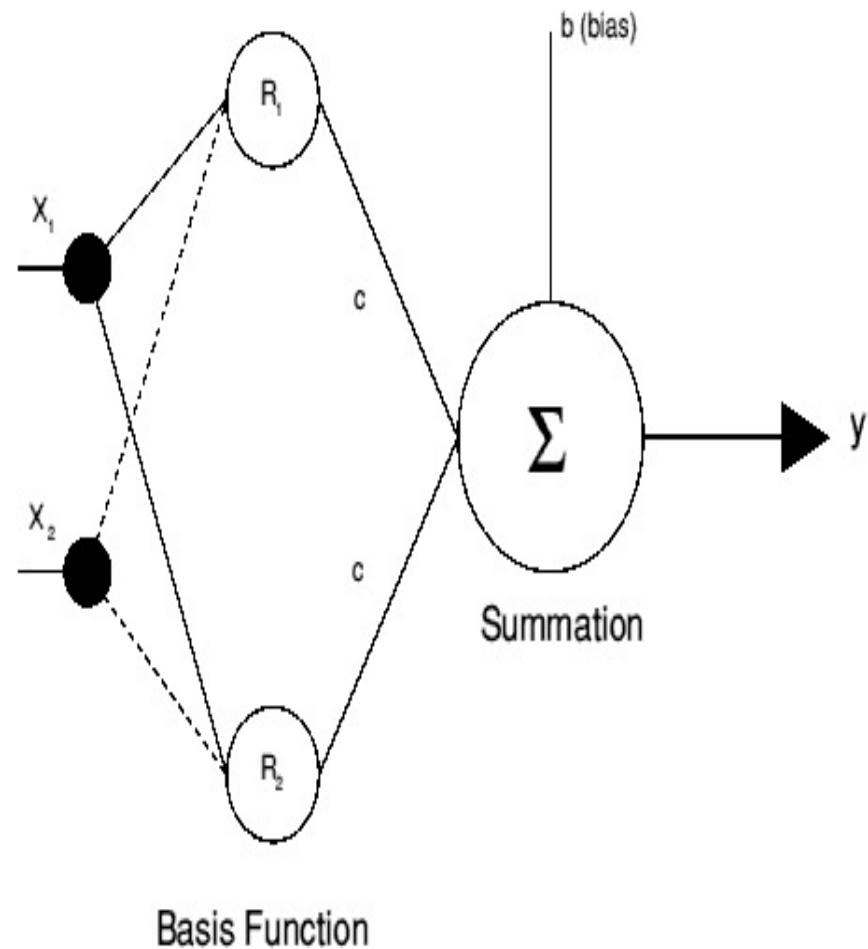
# SOLVING FOR C

$$\underline{c} = \underline{y} \times \underline{A}^{-1}$$

- Solving the above equation involves finding the inverse of  $A$ .
- There are a variety of techniques that can be used to produce this inverse and hence solve for the coefficients.
- The most widely used technique is **singular value decomposition**; this technique can deal with cases where the training patterns are closely related to each other.

# TRAINING WITH FIXED CENTRES: AN EXAMPLE

- A simple RBF which has:
  - 2 inputs,
  - 1 output,
  - 2 centres,
  - 2 unknown coefficients (one weight and one bias).



- The XOR problem revisited:

- $0\ 0 \Rightarrow 0$
- $0\ 1 \Rightarrow 1$
- $1\ 0 \Rightarrow 1$
- $0\ 0 \Rightarrow 0$

# THE CENTRES

- For the characterisation of the output unit we will use the following information:
  - The output unit will employ ***weight sharing***. This is justified by the underlying symmetry of the XOR problem.
  - The output unit includes a ***bias term***. This is justified by the fact that the desired output values of the XOR function have a non zero mean.

## ▪ The Centres

- The two centres are chosen as:
  - $R_1 = [1 \ 1]^T$
  - $R_2 = [0 \ 0]^T$

# BASIS FUNCTION

- Use a Gaussian function:

$$G(\|x - R\|) = e^{-\|x-R\|^2} = e^{-\sum_i (x_i - R_i)^2}$$

- Hence the output function is:  $y(x) = \sum_{i=1}^2 c G(\|x - R\|) + b$

- Hence learning in this RBF is the process of finding the weight coefficient,  $c$ , and the bias term,  $b$ . Four equations can be produced, from the 4 training patterns:

$$0.0 = c G(\|[0\ 0] - [1\ 1]\|) + c G(\|[0\ 0] - [0\ 0]\|) + b$$

$$1.0 = c G(\|[0\ 1] - [1\ 1]\|) + c G(\|[0\ 1] - [0\ 0]\|) + b$$

$$1.0 = c G(\|[1\ 0] - [1\ 1]\|) + c G(\|[1\ 0] - [0\ 0]\|) + b$$

$$0.0 = c G(\|[1\ 1] - [1\ 1]\|) + c G(\|[1\ 1] - [0\ 0]\|) + b$$

# SOLVING FOR $C$

## Calculating the basis function value

$$\begin{aligned} G([0 \ 0] - [1 \ 1]) &= e^{-(\sqrt((0-1)^2 + (0-1)^2))^2} \\ &= e^{-((-1)^2 + (-1)^2)} \\ &= e^{-\sqrt{2}} \\ &= 0.2431 \end{aligned}$$

$$0 = c(0.2431) + c(1.0000) + b$$

$$1 = c(0.3678) + c(0.3678) + b$$

$$1 = c(0.3678) + c(0.3678) + b$$

$$0 = c(1.0000) + c(0.2431) + b$$

Let  $y = [0 \ 1 \ 1 \ 0]^T$  and  $c = [c \ c \ b]^T$

$$\text{Let } A = \begin{vmatrix} 0.2431 & 1.0000 & 1 \\ 0.3678 & 0.3678 & 1 \\ 0.3678 & 0.3678 & 1 \\ 1.0000 & 0.2431 & 1 \end{vmatrix}$$

Then  $Ac = y$ . But this problem is **over-determined** which is why A is not a square matrix.

$$\text{and } c = \begin{vmatrix} -1.97 \\ -1.97 \\ +2.45 \end{vmatrix}$$

These values can now be used with new input patterns to generate new outputs as required.

# TRAINING ALGORITHM FOR RBFS WITH FIXED CENTRES

1. Choose appropriate centres  $\underline{R}$ .
2. Choose appropriate  $\sigma$ .
3. Calculate the output  $\underline{y} = \underline{A} \times \underline{c} = \Phi(\underline{x}_n, \underline{R}) \times \underline{c}$
4. Find the optimum weight values with  $\underline{c} = \underline{y} \times \underline{A}^{-1}$

# RBF NETWORKS

- As we have seen so far in training of RBF Network the centres  $R$  are fixed. We choose random  $R_s$  to cover any  $x$ .
- **Training RBF Networks with non Fixed Centres**
  - We use Stochastic Gradient Descent Approach
  - In this case we can change/train coefficients  $c$ , centres  $R$  and variances  $\sigma$ .

# TRAINING RBF NETWORKS WITH NON FIXED CENTRES

- The instantaneous total error is calculated by:

$$J[k] = \frac{1}{2} [e[k]]^2 = \frac{1}{2} [d[k] - \sum_{k=1}^N c_k [k] \Phi(x[k], R_k[k])]^2$$

- Given that  $\Phi$  is the Gaussian Function, the instantaneous total error becomes:

$$J[k] = \frac{1}{2} [d[k] - \sum_{k=1}^N c_k [k] e^M]^2$$

$$M = \left( \frac{\|x[k], R_k[k]\|^2}{\sigma_k^2 [k]} \right)$$

# TRAINING RBF NETWORKS WITH NON FIXED CENTRES

- The functions of Gradient Descent for  $c$ ,  $R$  and  $\sigma$  are given by:

$$c_k[k+1] = c_k[k] - n_c \frac{\partial J[k]}{\partial c_k} \mid c_k = c_k[k]$$

$$R_k[k+1] = R_k[k] - n_R \frac{\partial J[k]}{\partial R_k} \mid R_k = R_k[k]$$

$$\sigma_k[k+1] = \sigma_k[k] - n_\sigma \frac{\partial J[k]}{\partial \sigma_k} \mid \sigma_k = \sigma_k[k]$$

# TRAINING RBF NETWORKS WITH NON FIXED CENTRES

## Training Algorithms:

1. Choose appropriate centres  $c$ .
2. Choose appropriate  $\sigma$ .
3. Initialize the coefficients  $c$  with small random values.
4. Give an input to the network and calculate the network's output:

$$y[k] = \sum_{k=1}^N c_k \Phi(x[k], R_k[k], \sigma_k[k])$$

# TRAINING RBF NETWORKS WITH NON FIXED CENTRES

5. Calculate new values for  $c$ ,  $R$  and  $\sigma$ :

$$c[k+1] = c[k] + c[k] + n_c e[k] \Psi[k]$$

$$R_k[k+1] = R_k[k] + n_R \frac{e[k] c_k[k]}{\sigma_k^2[k]} \Phi(x[k], R_k[k], \sigma_k[k]) [x[k] - R_k[k]]$$

$$\sigma_k[k+1] = \sigma_k[k] + n_\sigma \frac{e[k] c_k[k]}{\sigma_k^2[k]} \Phi(x[k], R_k[k], \sigma_k[k]) ||x[k] - R_k[k]||^2$$

where

$$C[k] = [\Phi(x[k], R_1[k], \sigma_1[k]), \Phi(x[k], R_2[k], \sigma_2[k]), \dots, \Phi(x[k], R_N[k], \sigma_N[k])]$$

$$e[k] = d[k] - y[k]$$

6. Stop if the algorithm has converged, otherwise go back to 4.

# TRAINING RBF NETWORKS WITH NON FIXED CENTRES

- This training algorithm is **much easier than Backpropagation** because we have only one layer to train.
- The learning rates for  $c$ ,  $R$  and  $\sigma$  can be different between each other.

# MORE ON RBF NETWORKS AND MLPS - SIMILARITIES

**Similarities - both MLP and RBF NNs are:**

- examples of non-linear, layered, feed-forward networks
- universal approximators
  - Universality of RBF NNs: Any function can be approximated to arbitrary small error by an RBF NN given there are enough RBF neurons
  - Hence, there always exist an RBF network capable of accurately mimicking a specified MLP, and vice versa

# MORE ON RBF NETWORKS AND MLPS – DIFFERENCES

- An RBF network has a simple architecture (2 layers: nonlinear and linear) with single hidden layer; MLP may have 1 or 2 hidden layers and an active output layer
- MLPs may have a complex pattern of connectivity (not necessarily fully connected), RBF is a fully connected network.
- The hidden layer of an RBF net is non-linear, and the output layer is linear. The hidden and output layers of an MLP used as a classifier are nonlinear.
- Hidden and output neurons in a MLP share a common neuron model (nonlinear transformation of linear combination). Hidden neurons in RBF nets are quite different and serve a different purpose than the output neurons.

# RBF NETWORKS AND MLPS – DIFFERENCES (CONT)

- **Computation at hidden neurons**

- \* **The argument of the activation function of a hidden neuron in RBF nets computes the Euclidean distance between the input vector and the centre of the neuron.**
- \* **The activation function of a hidden neuron in MLP computes the inner product of the input vector and the weight vector of that neuron.**

- **Representation formed (in the space of hidden neurons activations with respect to the input space)**

- \* **An MLP is said to form a distributed representation since for a given input vector many hidden neurons will typically contribute to the output value.**
- \* **RBF nets form a local representation as for a given input vector only a few hidden neurons will have significant activation and contribute to the output.**

# RBF NETWORKS AND MLPS – DIFFERENCES (CONT)

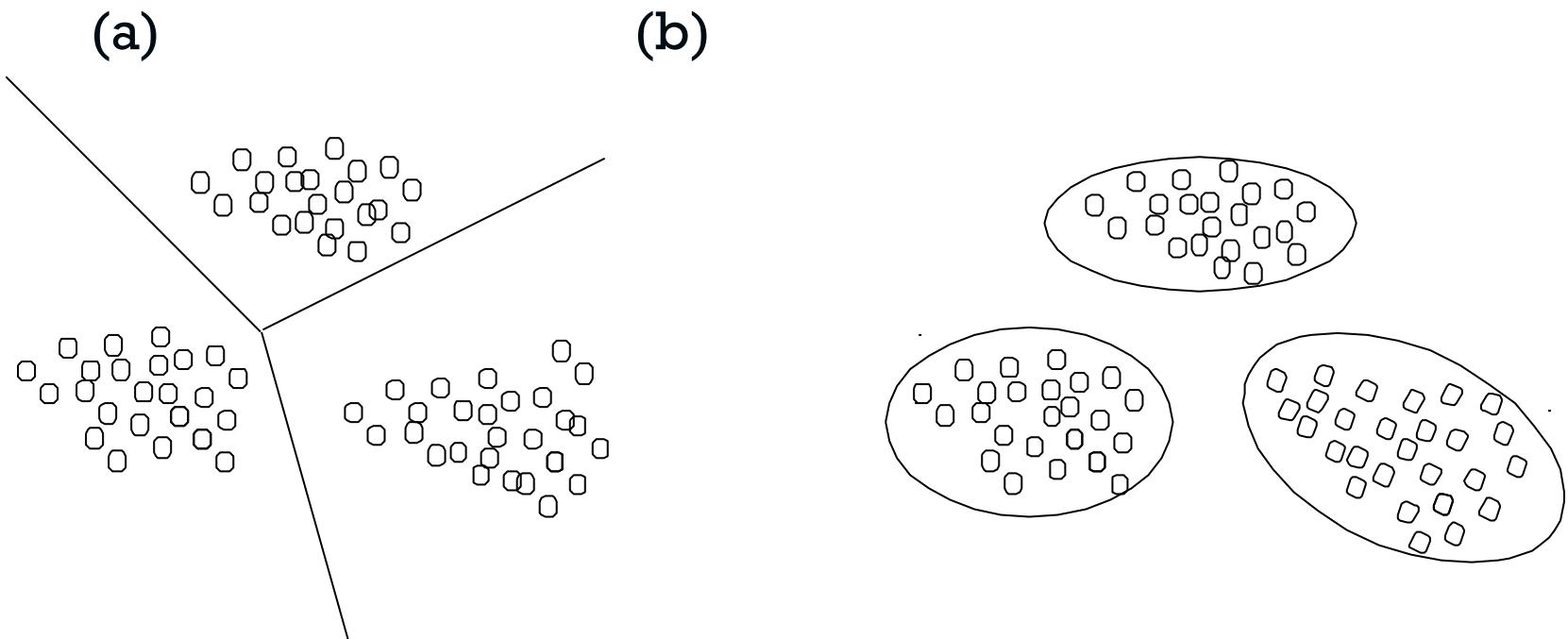
## Differences in how the weights are determined

- In MLP - usually at the same time as a part of a single global training strategy involving supervised learning.
- A RBF net is typically trained in 2 stages with the RBFs being determined first by **unsupervised technique using** the given data alone, and the weights between the hidden and output layer being found by fast linear supervised methods.

# RBF NETWORKS AND MLPS – DIFFERENCES (CONT)

## Differences in the separation of the data

- Hidden neurons in MLP define hyper-planes in the input space (a)
- RBF nets fit each class using kernel functions (b)



# RBF NETWORKS AND MLPS – DIFFERENCES (CONT)

## Number of neurons

- \* MLPs typically have less number of neurons than RBF NNs. This is because they respond to large regions of the input space.
- \* RBF nets may require many RBF neurons to cover the input space (larger number of input vectors and larger ranges of input vectors => more RBF neurons are needed.) However, they are capable of fast learning and have reduced sensitivity to the order of presentation of training data.

## RADIAL BASIS FUNCTIONS (continue)

**Vacations example (revisited): how much do you enjoy vacations as a function of their duration?**

**Ex. Duration (days) Rating (1-10)**

1	4	5
2	7	9
3	9	2
4	12	6

**We would like to create an interpolation and approximation networks that predict the ratings for other durations**

**Interpolation RBF Network – use of RBFs with fixed centres:**

- Can be solved by a RBF network with  $k$  neurons in the hidden layer and 1 output neuron
  - In this case we could make each neuron in the hidden layer to respond to one sample input vector
  - Centres would therefore be the input vectors in this case;  $\sigma$  is the width of the Gaussians; if small – each input vector will have only local influence, if wide – global influence
  - The output layer adds the weighted outputs of the hidden layer

will

# Equations for the Vacation Example

Given  $\sigma$ , we can compute the weights

\* A system of linear equations

\* Which are the unknowns? How many are they?

\* How many equations?

$$y(\mathbf{x}_1) = w_1 e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_1\|^2}{2\sigma^2}} + w_2 e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}} + w_3 e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_3\|^2}{2\sigma^2}} + w_4 e^{-\frac{\|\mathbf{x}_1 - \mathbf{x}_4\|^2}{2\sigma^2}}$$

$$y(\mathbf{x}_2) = w_1 e^{-\frac{\|\mathbf{x}_2 - \mathbf{x}_1\|^2}{2\sigma^2}} + w_2 e^{-\frac{\|\mathbf{x}_2 - \mathbf{x}_2\|^2}{2\sigma^2}} + w_3 e^{-\frac{\|\mathbf{x}_2 - \mathbf{x}_3\|^2}{2\sigma^2}} + w_4 e^{-\frac{\|\mathbf{x}_2 - \mathbf{x}_4\|^2}{2\sigma^2}}$$

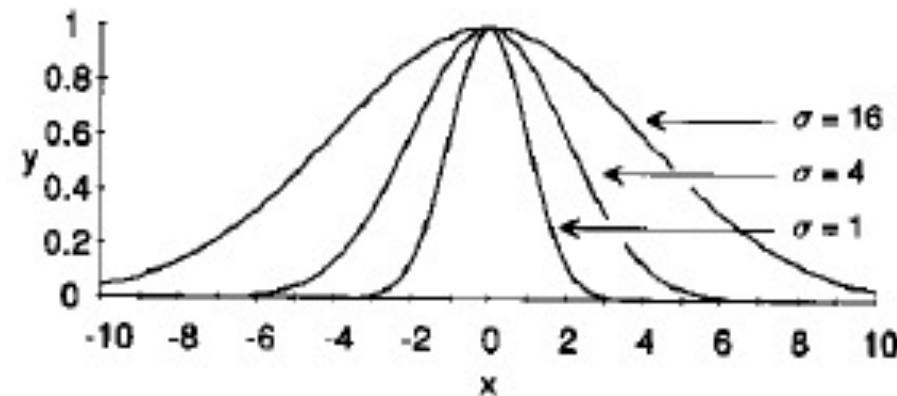
$$y(\mathbf{x}_3) = w_1 e^{-\frac{\|\mathbf{x}_3 - \mathbf{x}_1\|^2}{2\sigma^2}} + w_2 e^{-\frac{\|\mathbf{x}_3 - \mathbf{x}_2\|^2}{2\sigma^2}} + w_3 e^{-\frac{\|\mathbf{x}_3 - \mathbf{x}_3\|^2}{2\sigma^2}} + w_4 e^{-\frac{\|\mathbf{x}_3 - \mathbf{x}_4\|^2}{2\sigma^2}}$$

$$y(\mathbf{x}_4) = w_1 e^{-\frac{\|\mathbf{x}_4 - \mathbf{x}_1\|^2}{2\sigma^2}} + w_2 e^{-\frac{\|\mathbf{x}_4 - \mathbf{x}_2\|^2}{2\sigma^2}} + w_3 e^{-\frac{\|\mathbf{x}_4 - \mathbf{x}_3\|^2}{2\sigma^2}} + w_4 e^{-\frac{\|\mathbf{x}_4 - \mathbf{x}_4\|^2}{2\sigma^2}}$$

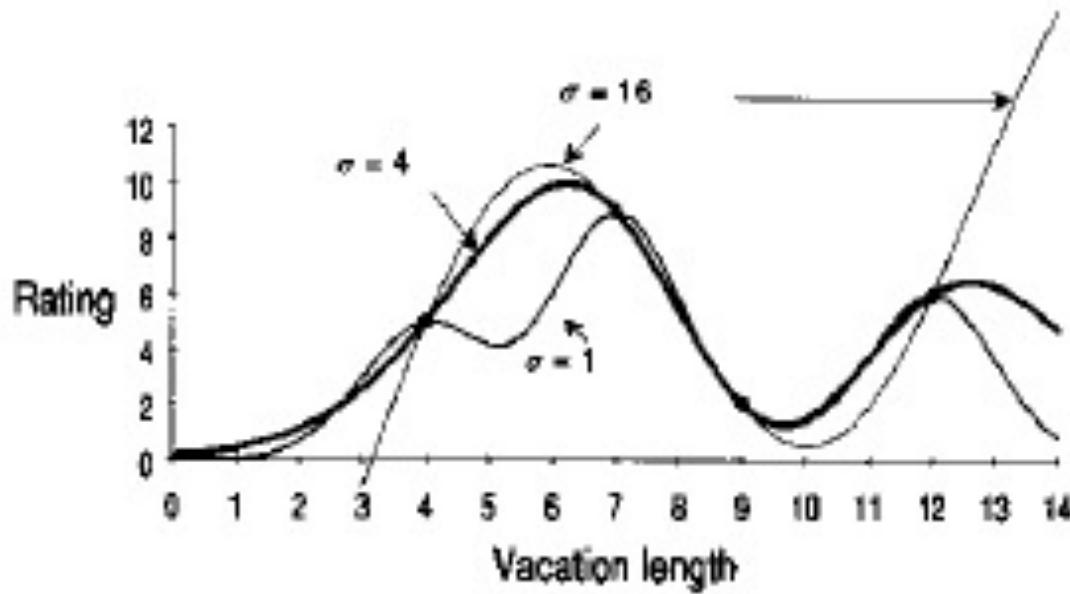
# Vacation Example - Solving the Equations

\* Solutions for 3 different  $\sigma$ :

$\sigma$	w1	w2	w3	w4
1	4.90	8.84	0.73	5.99
4	0.87	13.93	-9.20	8.37
16	-76.50	236.49	-237.77	87.55



\* The vacation rating approximation functions



- Small  $\sigma \Rightarrow$  bumpy interpolations
- Large  $\sigma \Rightarrow$  extend the influence of each sample
- Too large  $\sigma \Rightarrow$  too global influence of each sample

## **Summary: How to Create an Interpolation RBF Net**

- For each given sample (training set example), create a neuron in the first layer (i.e. the hidden layer) centred on the sample input vector
- Create a system of linear equations
  - Choose  $\sigma$
  - Compute the Euclidean distance between the sample input and each of the centres
    - Compute the Gaussian functions of each distance
    - Multiply the Gaussian function by the corresponding weight/coefficient of the neuron
    - Equate the sample output with the sum of the weighted Gaussian functions of distance
- Solve the system equations for the weights/coefficients
- There is no training of the weights but solving of equations
- The network can be used to predict the output of new examples

# Creation of an Approximation RBF net

- If there are **many samples**, the number of neurons in the hidden layer of the interpolation nets will be high
- Solution: Approximate the unknown function with a smaller number of neurons that are **somewhat representative**
  - The simplest way – pick up a few random samples to serve as centres
- Such networks are called approximation nets
  - As the number of neurons in the hidden layer is lower than the number of samples, a set of weights/coefficients such that the network **produces the exact correct output for all samples does not exist**
    - However, a set of weights/coefficients that result in a **reasonable approximation for all samples can be found**
    - One way is to use the gradient descent to minimize the error – *RBF networks with adjustable centres*

# How to Create RBF Approximation Nets

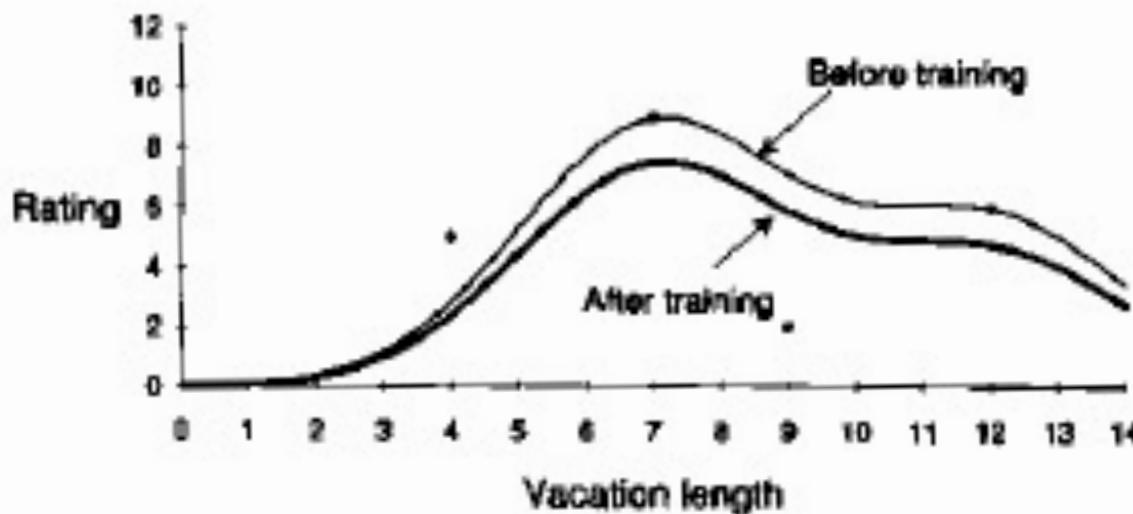
- For a few samples, create an approximation RBF net
- Select learning rates ( $\eta_{w/c}$ ,  $\eta_R$ ,  $\eta_\sigma$ )
- Until performance is satisfactory (error below a minimum threshold or max number of epochs reached):
  - For all sample inputs
  - Compute the resulting outputs
  - Compute the adjustment of weights/coefficients, the centres and the standard deviations
  - Add up the parameter ( $c/w$ ,  $R$ ,  $\sigma$ ) changes for all the sample inputs and update the parameters (batch mode) or update the parameters after each example (online update)

# Approximation RBF Net for the Vacation Example

## Adjustment of the weights/coefficients only

- 2 neurons only in the hidden layer
- initialized to example 2 (7 days) and 4 (12 days) (i.e., these are taken as the centres for the two hidden neurons)
- Learning rate  $\eta_{w/c} = 0.1$ ,  $\sigma = 4$
- The weights after 100 epochs:

	w1	w2	R1	R2
Initial values	8.75	5.61	7	12
Final values	7.33	4.47	7	12

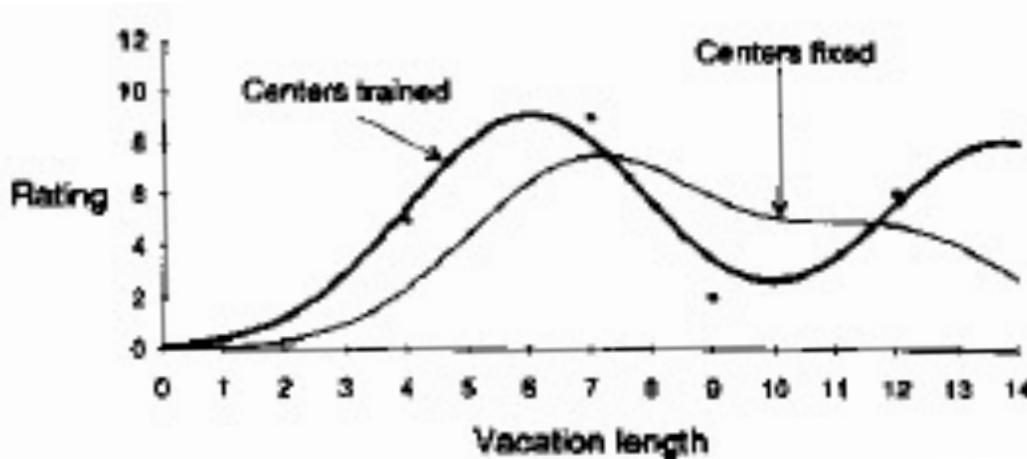


# Approximation RBF Net for the Vacation Example

## Adjustment of the weights/coefficients and centres

- \* Not only weights/coefficients but also the centres and standard deviations can be adjusted
- \* How much (for the centres)? Find the derivatives of the error function with respect to the centre coordinates (for centre adjustment, see relevant equation from previous lecture)
- \* The weights after 100 epochs:

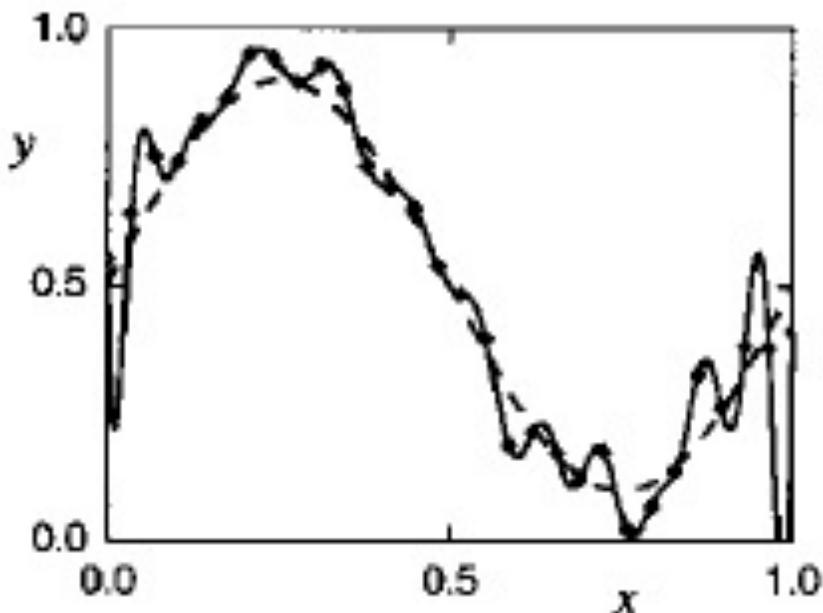
	w1	w2	R1	R2
Initial values	8.75	5.61	7	12
Final values	9.13	8.06	6	13.72



- \* Adjustment of both weights and centres provides a better approximation than adjustment of either of them alone

# RBFs for Interpolation and Approximation – Summary Example

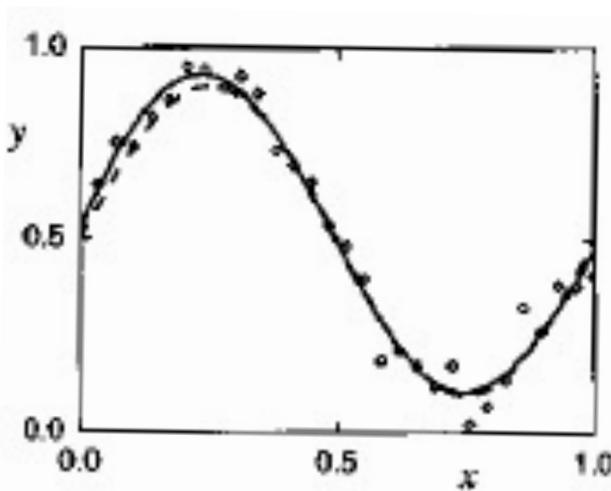
- Example taken from Bishop (1995), “NNs for Pattern Recognition”
- A set of 30 data points was generated by sampling the function  $y = 0.5 + 0.4\sin(2\pi x)$ , shown by the dashed curve, and adding Gaussian noise with standard deviation 0.05



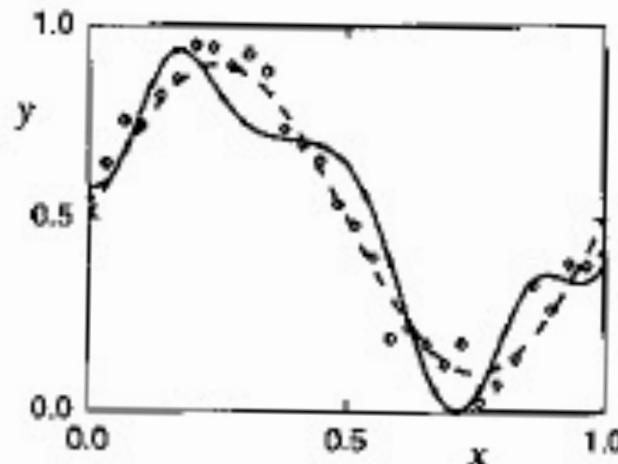
- RBF network for exact interpolation has been created
- Solid curve - the interpolation function which resulted using  $\sigma = 0.067$  ( $\approx$  twice the average spacing between the data points)
- The weights between the RBF hidden and output layers were found by solving the system of equations

## Summary Example (continue)

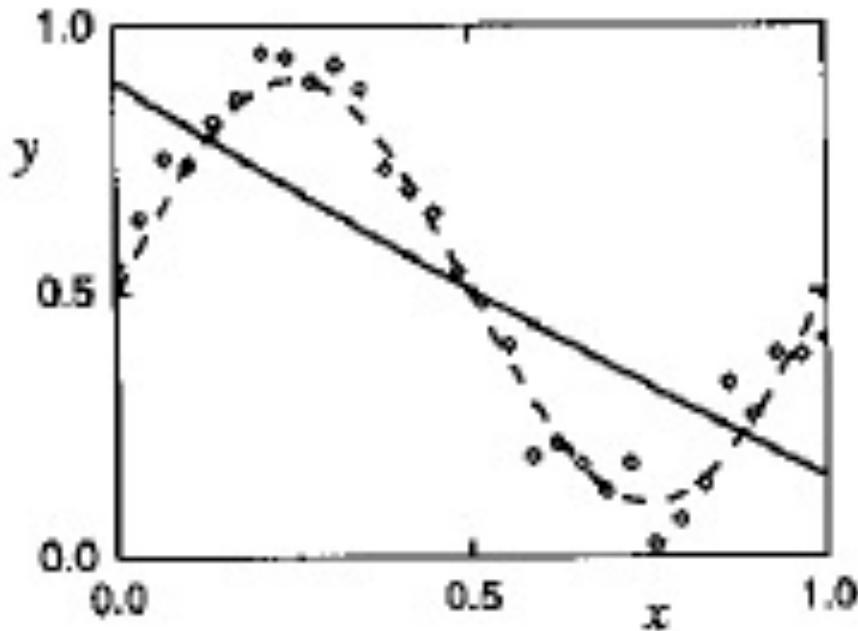
- Approximation RBF networks with 3 different widths have been created



- RBF hidden neurons = 5 (significantly smaller than the number of data points)
- RBF centres were set to a random subset of the given data points
- $\sigma = 0.4$  ( $\approx$  twice the average spacing between centres)
- The weights were found by training
- Good fit
- As above, but  $\sigma = 0.08$
- The resulting network function is insufficiently smooth and gives a poor representation of the underlying function which generated the data



## Summary Example (continue)



- As on the previous slide, but with  $\sigma = 10$
- The resulting network function is over-smoothed (the receptive fields of the RBF overlap too much) and again gives a poor representation of the underlying function which generated the data
- Conclusion: To get good generalisation, the width should be larger than the average distance between adjacent input vectors but smaller than the distance across all input space

# RBFS – USEFUL NOTES

## Width of the Gaussian (or Standard Deviation $\sigma$ )

**Small  $\sigma \rightarrow$  bumpy interpolations**

**Large  $\sigma \rightarrow$  extend the influence of each sample**

**Very large  $\sigma \rightarrow$  too global influence of each sample**

**To get a good generalisation, the width should be larger than the average distance between adjacent vectors, but smaller than the distance across all input space.**

# RBF NNS FOR CLASSIFICATION

## Cover's Theorem on the Separability of Patterns

**“A complex classification problem cast in *high dimensional* space nonlinearly is more likely to be *linearly separable* than in a *low dimensional* space” (Cover, 1965).** This includes two ingredients:

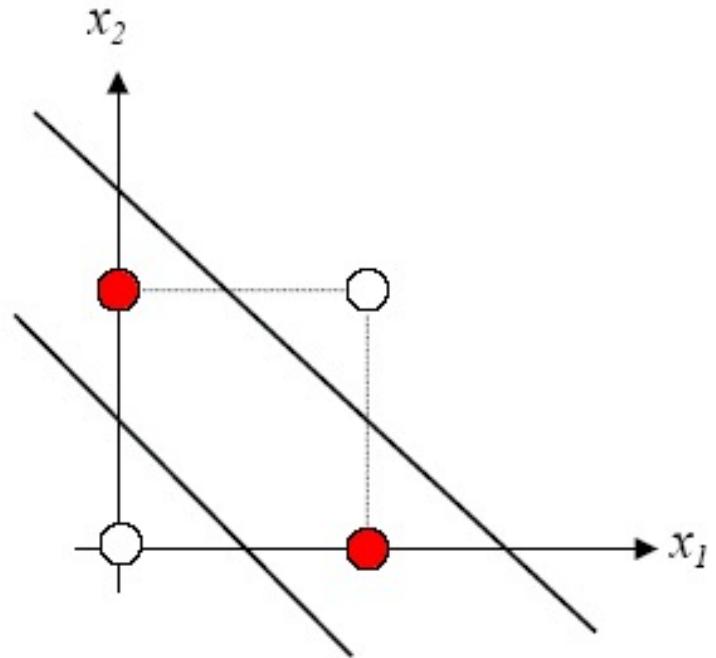
1. **Nonlinear formulation of hidden functions**
2. **High dimensionality of the hidden space compared to the input space (determined by the number of hidden neurons)**

**In some cases, however, the use of nonlinear mapping (i.e., point 1) may be sufficient to produce linear separability without having to increase the dimensionality of the hidden-neuron space (see example below)**

## THE XOR PROBLEM REVISITED

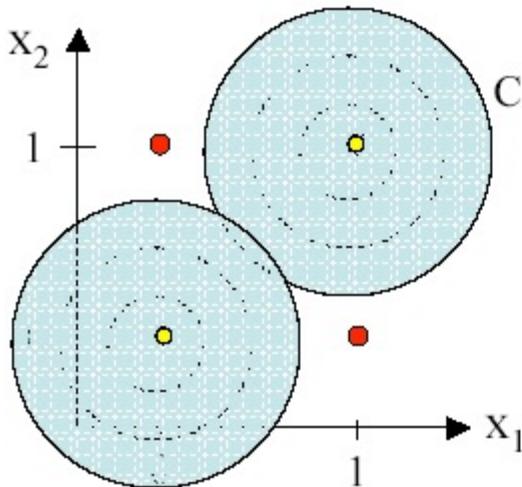
We are already familiar with the non-linearly separable XOR problem:

$p$	$x_1$	$x_2$	$t$
1	0	0	0
2	0	1	1
3	1	0	1
4	1	1	0



We know that Single Layer Perceptrons with step activation functions cannot generate the right outputs, because they are only able to form a single decision boundary. To deal with this problem we needed to change the activation function and introduce a non-linear hidden layer to form the Multi Layer Perceptron (MLP).

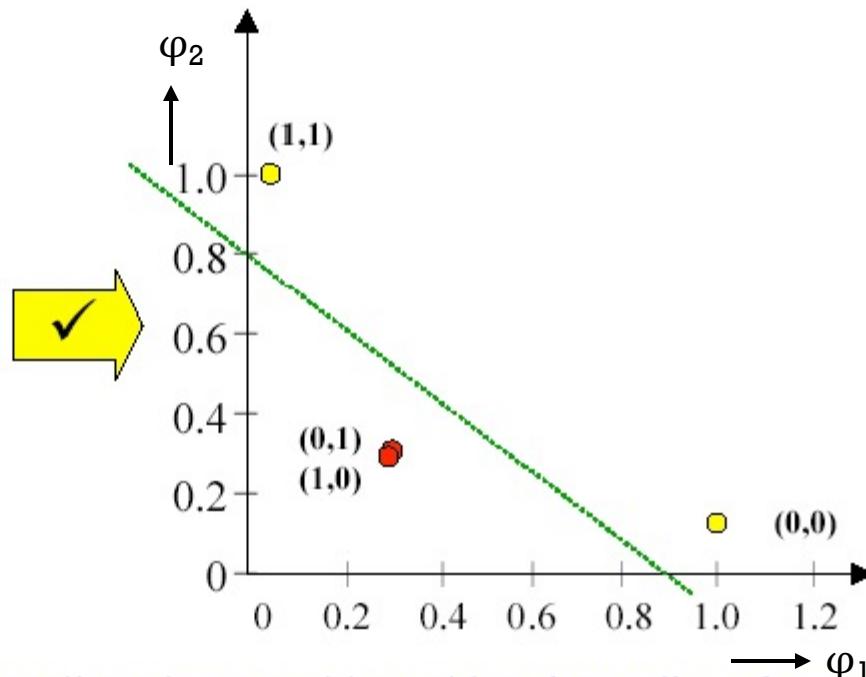
# RBF NNS ON XOR PROBLEM



Consider the nonlinear functions to map the input vector  $\mathbf{x}$  to the  $\varphi_1$ - $\varphi_2$  space

$$\mathbf{x} = [x_1 \ x_2] \quad \varphi_1(\mathbf{x}) = e^{-\|\mathbf{x}-\mathbf{t}_1\|^2} \quad \mathbf{t}_1 = [1 \ 1]^T$$
$$\varphi_2(\mathbf{x}) = e^{-\|\mathbf{x}-\mathbf{t}_2\|^2} \quad \mathbf{t}_2 = [0 \ 0]^T$$

Input $\mathbf{x}$	$\varphi_1(\mathbf{x})$	$\varphi_2(\mathbf{x})$
(1,1)	1	0.1353
(0,1)	0.3678	0.3678
(1,0)	0.3678	0.3678
(0,0)	0.1353	1



The nonlinear  $\varphi$  function transformed a nonlinearly separable problem into a linearly separable one !!!

## THE XOR PROBLEM OUTPUT WEIGHTS

In this case we just have one output  $y(\mathbf{x})$ , with one weight  $w_j$  to each hidden unit  $j$ , and one bias -  $\theta$ . This gives us the network's input-output relation for each input pattern  $\mathbf{x}$

$$y(\mathbf{x}) = w_1\varphi_1(\mathbf{x}) + w_2\varphi_2(\mathbf{x}) - \theta$$

Then, if we want the outputs  $y(\mathbf{x}^p)$  to equal the targets  $t^p$ , we get the four equations

$$1.0000w_1 + 0.1353w_2 - 1.0000\theta = 0$$

$$0.3678w_1 + 0.3678w_2 - 1.0000\theta = 1$$

$$0.3678w_1 + 0.3678w_2 - 1.0000\theta = 1$$

$$0.1353w_1 + 1.0000w_2 - 1.0000\theta = 0$$

Three are different, and we have three variables, so we can easily solve them to give

$$w_1 = w_2 = -2.5018 , \quad \theta = -2.8404$$

This completes our “training” of the RBF network for the XOR problem.

# CREATING RBF NETWORKS FOR CLASSIFICATION

## 1. Select the locations of the centers $R_i$

- Choose them randomly from the training data or
- Use a clustering algorithm to find them, e.g. SOM

## 2. Fit the RBFs, i.e. set the width of the Gaussians

- Set their width using the p-nearest neighbour method
  - $\sigma_i$  of a centre  $i$  is the root mean squared distance from that centre to the p nearest centres
  - p is set by trial and error, typically p=2
  - Small p => narrower Gaussians (may be expected to emphasize the nature of the training data, great risk of overtraining)
  - Higher p => will broaden the Gaussians, may be expected to suppress the distinction between clusters
- Other heuristics:  $\sigma_i = \sqrt{d_i}$  here  $d_i$  is the distance to the nearest centre

## • 3. Gradient descent to train the linear layer

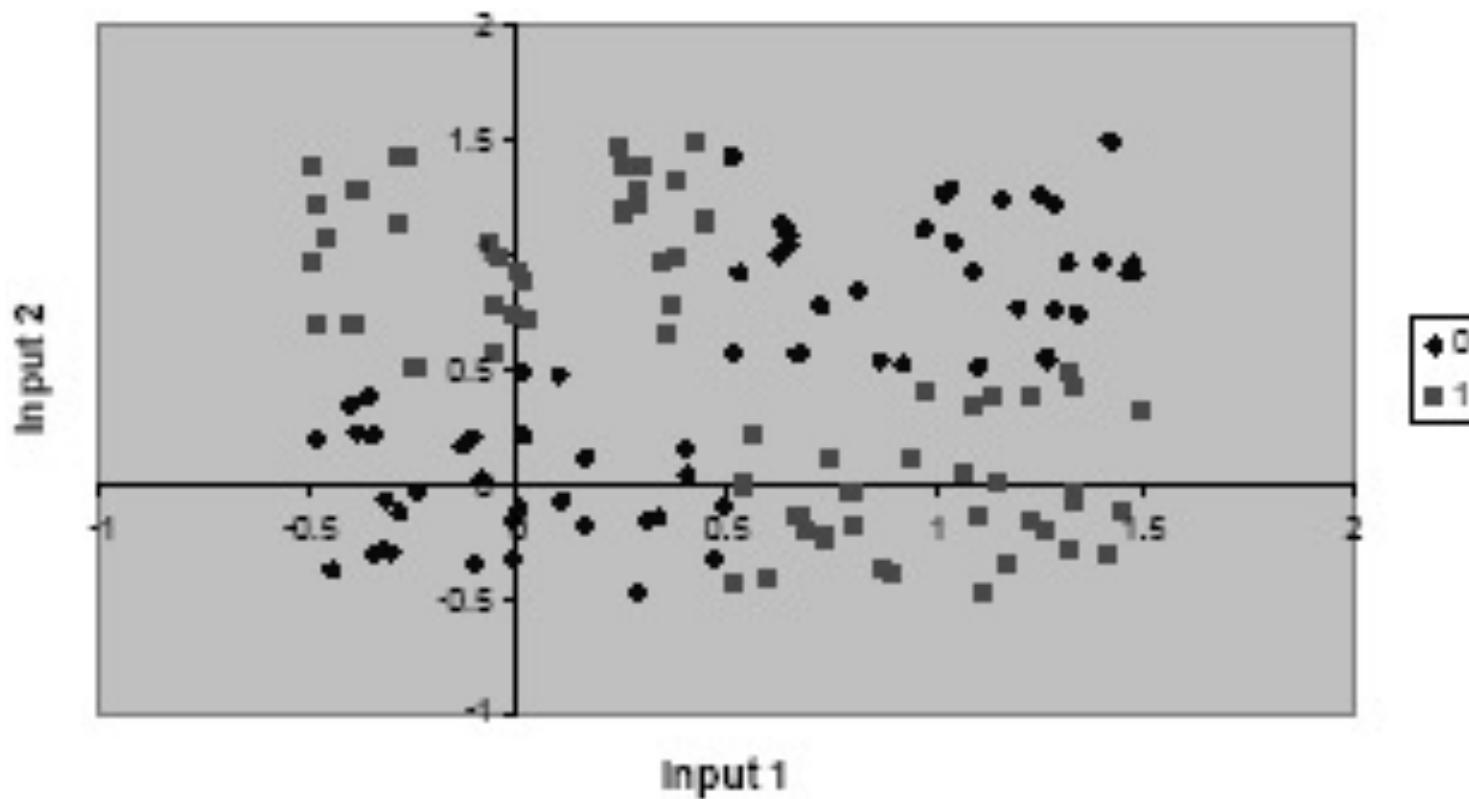
- Typically one output neuron for each class => i.e. binary encoding of the targets

# CLASSIFICATION EXAMPLE – NOISY XOR PROBLEM

Noisy XOR truth table:

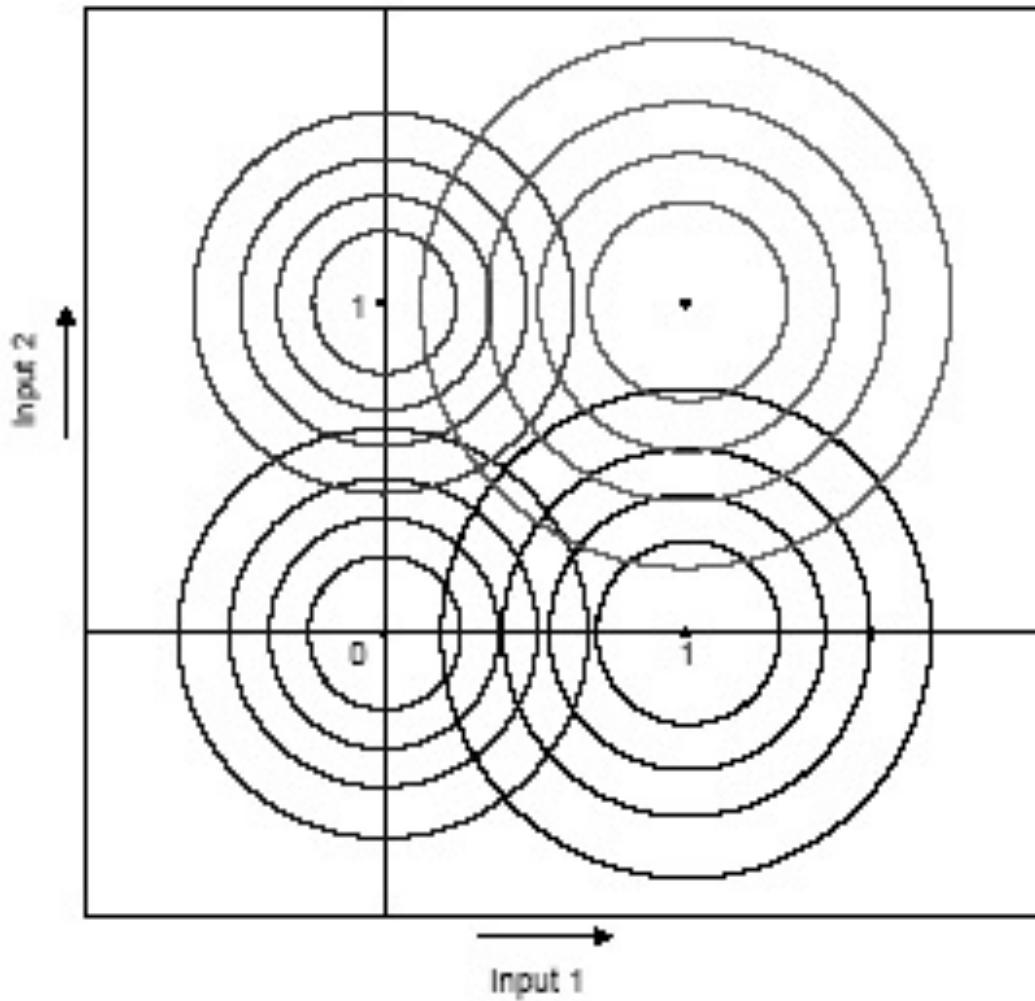
	Input 1	Input 2	Output
	<0.5	<0.5	0
	<0.5	>0.5	1
	>0.5	<0.5	1
	>0.5	>0.5	0

Correct Classifications



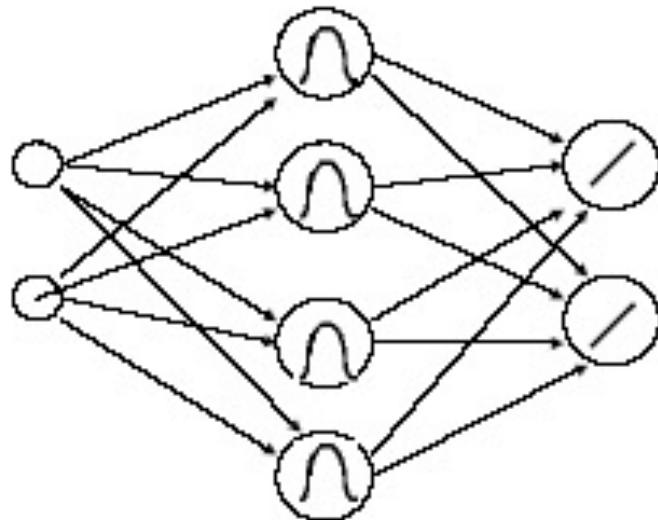
# NOISY XOR PROBLEM (CONTINUE)

1. Clustering the data using Kohonen SOM
2. Locating the centres of the RBF by setting them to the centres of the clusters
3. Determining their widths
4. The corresponding RBF net:  
\* 2 input and 4 RBF hidden neurons



## NOISY XOR PROBLEM (CONTINUE)

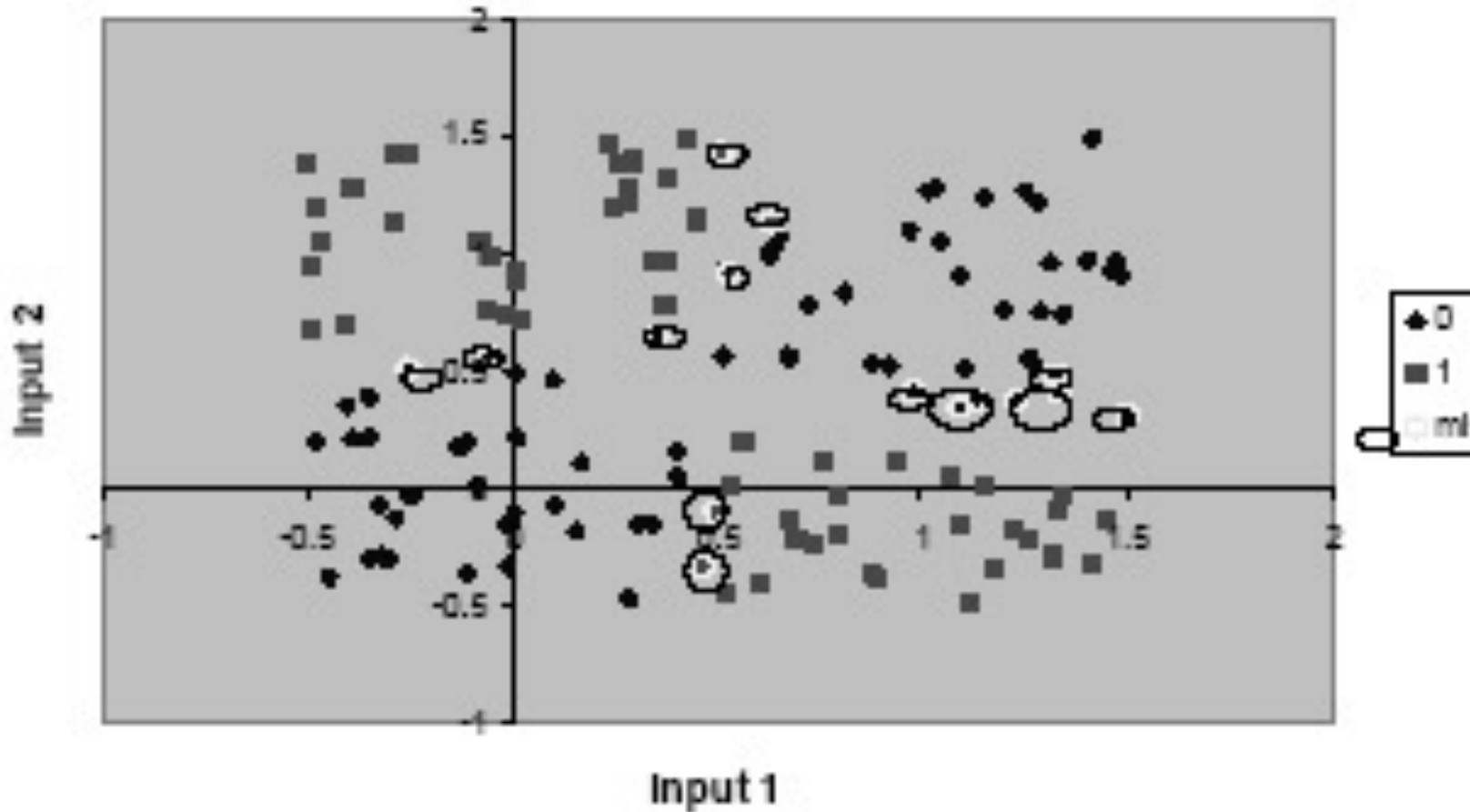
5. Create an output neuron for each class and use binary encoding for the targets (1 0 for class 1 and 0 1 for class 2)



6. Initialize the weights between hidden and output layer to small random values
7. Train the weights between the hidden and output layer using the gradient descent algorithm

## NOISY XOR PROBLEM (CONTINUE)

Classification by the RBF network (misclassifications are circled)



- Where are the misclassifications?
- Why RBF nets are not perfect fit for this data?

# RBF CODE

# MIDTERM

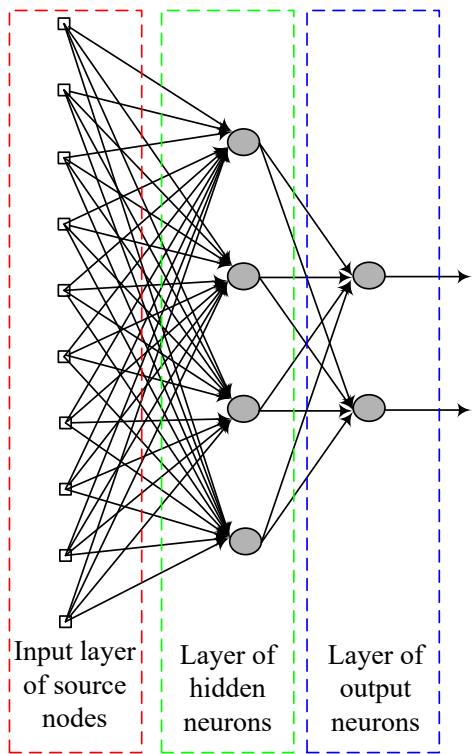
- 1. Course Material in Moodle (*Midterm Material!!!!*)**
- 2. Use textbooks just to help you understand the course material**
  - 1. Fundamentals of Neural Networks: Architectures, Algorithms, and Applications by Laurene Fausett, Prentice Hall, Englewood Cliffs, NJ, 1994.**
  - 2. Neural Networks and Learning Machines (3rd ed.) by Simon O. Haykin, Prentice Hall, Englewood Cliffs, NJ, 2008.**
    - Chapter 1: 1.1 to 1.6
    - Chapter 2: 2.2, 2.8, 2.9
    - Chapter 3: 3.5, 3.8
    - Chapter 4: 4.1 to 4.10, 4.11 to 4.14
    - Chapter 5: 5.1- 5.11

# OVERVIEW

- Radial Basis Function Networks

ANY  
QUESTIONS?

# NEXT LECTURE



- Midterm!!!
- BRNNs