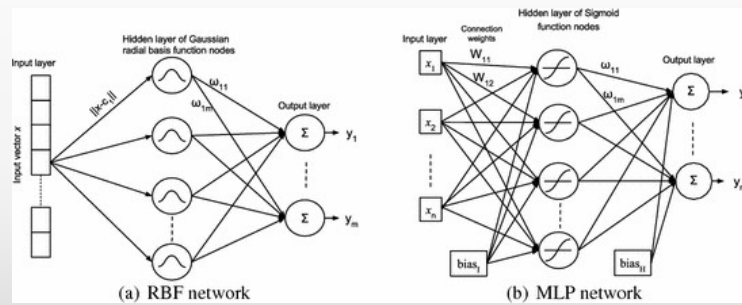# RBF Networks and  Support Vector Machine (SVM)

## RADIAL BASIS FUNCTION (RBF) NETWORKS

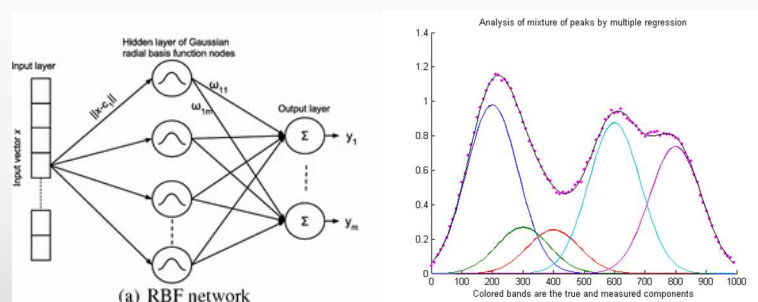# INTRODUCTION

- A radial basis function (RBF) network is an artificial neural network that uses RBFs as activation functions. The output of the network is a linear combination of radial basis functions of the inputs and neuron parameters.
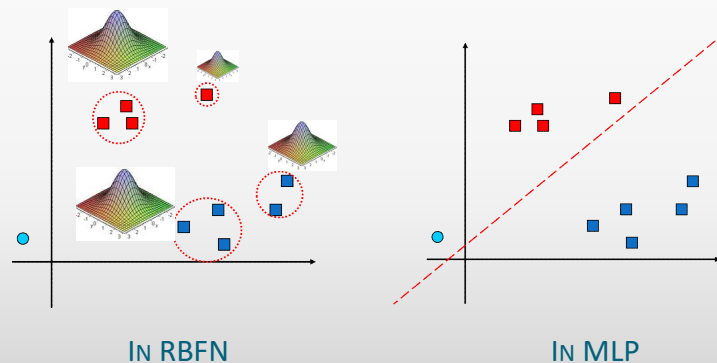


(a) RBF network     (b) MLP network

---

# INTRODUCTION

- On the other hand, an RBF neuron measures similarity or distance. RBF network can be seen as a completely different approach by <u>designing a neural network as a curve-fitting (approximation) problem in high-dimensional space.</u>



(a) RBF network

## INTRODUCTION

- Similar to a MLP, an RBF network is a kind of supervised neural network. However, it designs a neural network as curve-fitting problem. That is, approximate function with linear combination of Radial basis functions
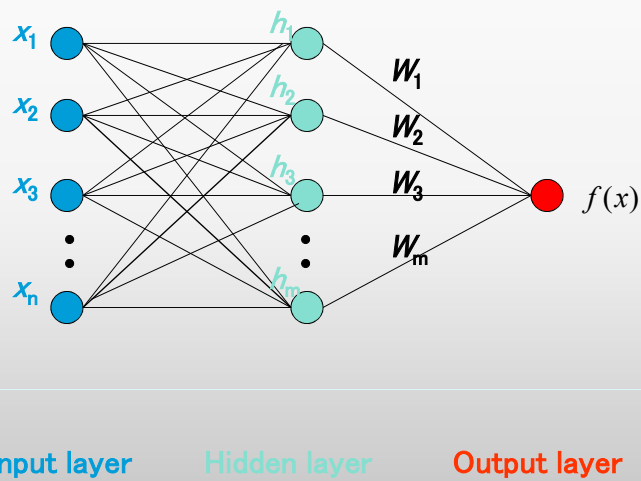


IN RBFN — IN MLP

---

## RADIAL BASIS FUNCTION NETWORK (SUMMARY)

- A kind of supervised neural networks
- Design of NN as curve-fitting problem. That is, approximate function with linear combination of Radial basis functions
- Learning
  - find surface in multidimensional space best fit to training data
- Generalization
  - Use of this multidimensional surface to interpolate the test data

## ARCHITECTURE

$x_1$   $h_1$   $W_1$

$x_2$   $h_2$   $W_2$

$x_3$   $h_3$   $W_3$   $f(x)$

$x_n$   $h_m$   $W_m$

**Input layer**    Hidden layer    **Output layer**

---

architecture

## RADIAL BASIS FUNCTION

$$f(x) = \sum_{j=1}^{m} w_j h_j(x)$$

$$h_j(x) = \exp\left(-\beta_j \parallel x - c_j \parallel^2\right)$$

Where $c_j$ is center of a region,

$\beta_j$ is a parameter describing the

width of the receptive field

## DESIGNING

- Require
  - Determine the number of radial basis neurons
  - Determine the radial basis function centers
  - Select the radial basis function width parameter
    - ✓ smaller width

      alerting in untrained test data
    - ✓ larger width

      network of smaller size & faster
      execution

---

**learning strategies**

## CAN ALSO BE TRAINED LIKE AN MLP (1/2)

$$E = \frac{1}{2}\sum_{i=1}^{N}(d_i(x) - f_i(x))^2 = \frac{1}{2}\sum_{i=1}^{N}e_i(x)^2 \qquad f(x) = \sum_{j=1}^{m}w_j h_j(x)$$

$$h_j(x) = \exp\left(-\beta_j \| x - c_j \|^2\right)$$

- Linear weights (output layer)

$$w_j \leftarrow w_j - \eta_1 \frac{\partial E}{\partial w_j}, \qquad j = 1,2,...,m$$

- Positions of centers (hidden layer)

$$c_j \leftarrow c_j - \eta_2 \frac{\partial E}{\partial c_j}, \qquad j = 1,2,...,m$$

- Spreads of centers (hidden layer)

$$\beta_j \leftarrow \beta_j - \eta_3 \frac{\partial E}{\partial \beta_j}, \qquad j = 1,2,...,m$$

**learning strategies**

## CAN ALSO BE TRAINED LIKE AN MLP (2/2)

$$E = \frac{1}{2}\sum_{i=1}^{N}(d_i(x) - f_i(x))^2 = \frac{1}{2}\sum_{i=1}^{N}e_i(x)^2 \qquad f(x) = \sum_{j=1}^{m}w_j h_j(x)$$

$$h_j(x) = \exp\left(-\beta_j \| x - c_j \|^2\right)$$

- Gradient of linear weights (output layer)

$$\frac{\partial E}{\partial w_j} = \frac{\partial E}{\partial e_i(x)}\frac{\partial e_i(x)}{\partial w_j} = -\sum_{i=1}^{N}e_i(x)h_j(x), \qquad j = 1,2,...,m$$

- Gradient of positions of centers (hidden layer)

$$\frac{\partial E}{\partial c_j} = \frac{\partial E}{\partial e_i(x)}\frac{\partial e_i(x)}{\partial h_j(x)}\frac{\partial h_j(x)}{\partial c_j} = -\sum_{i=1}^{N}e_i(x)w_j h_j(x)\beta_j \| x - c_j \|, \quad j = 1,2,...,m$$

- Gradient of spreads of centers (hidden layer)

$$\frac{\partial E}{\partial \beta_j} = \frac{\partial E}{\partial e_i(x)}\frac{\partial e_i(x)}{\partial h_j(x)}\frac{\partial h_j(x)}{\partial \beta_j} = \sum_{i=1}^{N}e_i(x)w_j h_j(x) \| x - c_j \|^2, \quad j = 1,2,...,m$$

## MLP VS RBFN

| | |
|---|---|
| Global hyperplane | Local receptive field |
| EBP (Error BP) | LMS (Least mean square) |
| Local minima | Serious local minima |
| Smaller number of hidden neurons | Larger number of hidden neurons |
| Shorter computation time | Longer computation time |
| Longer learning time | Shorter learning time |

## APPROXIMATION

- ○ MLP : Global network
  - • All inputs cause an output
- ○ RBF : Local network
  - • Only inputs near a receptive field produce an activation
  - • Can give "don't know" output

## SUPPORT VECTOR MACHINE (SVM)
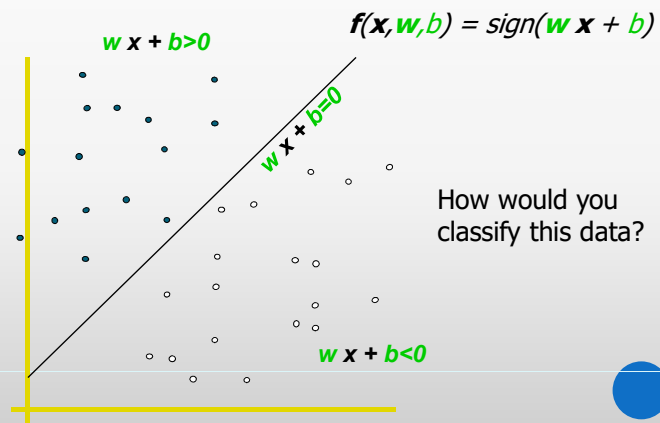
**SVM : a margin optimized classifier**

**SVMs developed by Vapnik in 1963, have gained wide acceptance because of their high generalization ability for a wide range of applications in 1990's.**

# Introduction
## -- SVM : a margin optimized classifier

○ Linear classifier
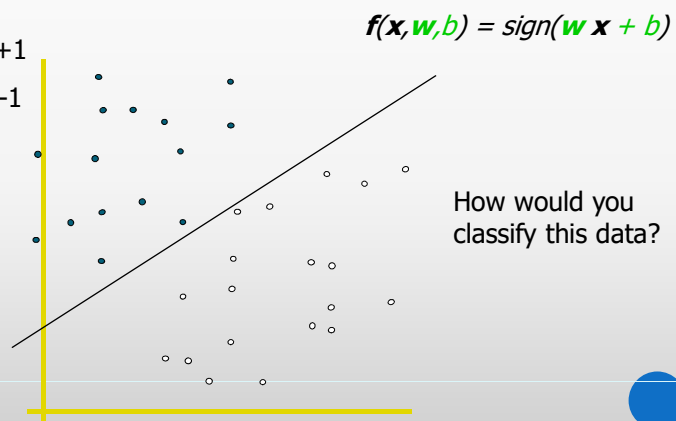
- denotes +1
○ denotes -1

$f(\textbf{x},\textbf{w},b) = sign(\textbf{w}\ \textbf{x} + b)$

**w x + b>0**

**w x + b=0**
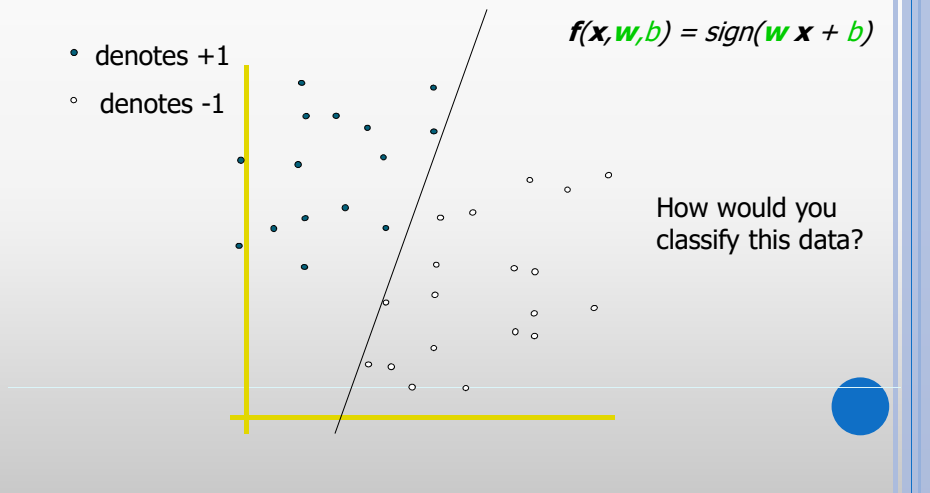
**w x + b<0**

How would you classify this data?

# Introduction
## -- SVM : a margin optimized classifier

○ Linear classifier

- denotes +1
○ denotes -1

$f(\textbf{x},\textbf{w},b) = sign(\textbf{w}\ \textbf{x} + b)$

How would you classify this data?
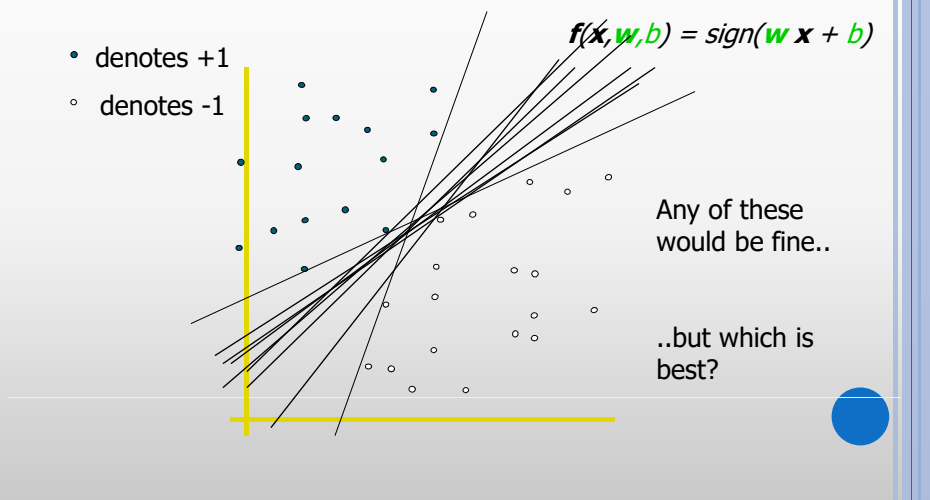
# Introduction

## -- SVM : a margin optimized classifier

- Linear classifier

$f(x,w,b) = sign(w\ x + b)$

- · denotes +1
- ○ denotes -1

How would you
classify this data?

---

# Introduction

## -- SVM : a margin optimized classifier

- Linear classifier

$f(x,w,b) = sign(w\ x + b)$

- · denotes +1
- ○ denotes -1

Any of these
would be fine..

..but which is
best?

# Introduction

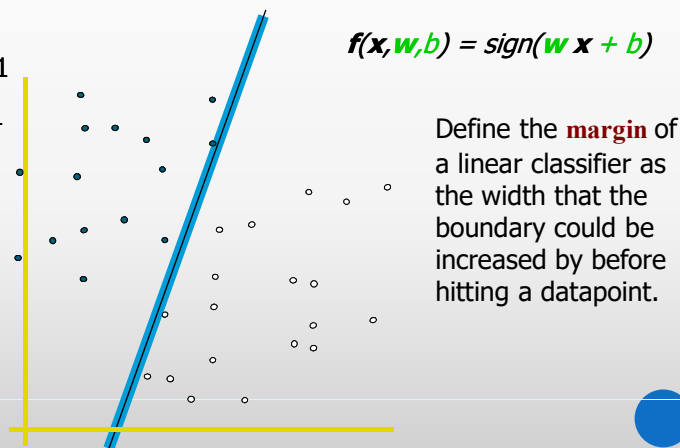## -- SVM : a margin optimized classifier

○ Linear classifier

$f(x,w,b) = sign(w\ x + b)$

• denotes +1

○ denotes -1

How would you classify this data?

**Misclassified to +1 class**

---

# Introduction

## -- SVM : a margin optimized classifier

○ Linear classifier

$f(x,w,b) = sign(w\ x + b)$

• denotes +1

○ denotes -1

Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Introduction

## -- SVM : a margin optimized classifier

○ Linear classifier

• denotes +1

○ denotes -1

1. Maximizing the margin is good according to intuition and PAC theory
2. Implies that only support vectors are important; other training examples are ignorable.
3. Empirically it works very very well.

**Support Vectors** are those datapoints that the margin pushes up against
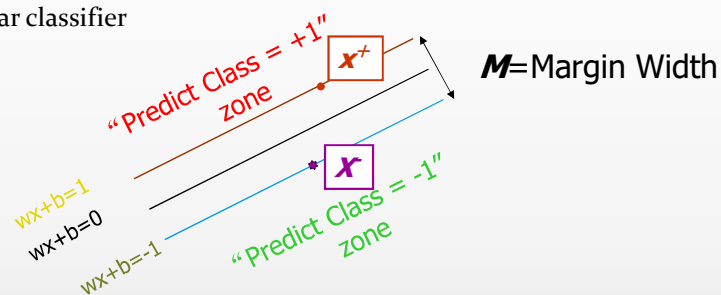
**classifier** is the linear classifier with the maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

---

# SVM Formulation
# -- Definition of margin

○ Linear classifier

"Predict Class = +1" zone

$x^+$

$M$=Margin Width

$x^-$

"Predict Class = -1" zone

wx+b=1

wx+b=0

wx+b=-1

What we know:

○ $w \cdot x^+ + b = +1$

○ $w \cdot x^- + b = -1$

○ $w \cdot (x^+ - x^-) = 2$

$$M = \frac{(x^+ - x^-) \cdot w}{|w|} = \frac{2}{|w|}$$

# SVM Formulation
## -- Maximizing the margin

○ Goal:  **1) Correctly classify all training data**

$$wx_i + b \geq 1 \qquad \textit{if } y_i = +1$$
$$wx_i + b \leq -1 \qquad \textit{if } y_i = -1$$
$$y_i(wx_i + b) \geq 1 \qquad \text{for all i}$$

**2) Maximize the Margin** $\quad M = \frac{2}{|w|}$

**same as minimize** $\quad \frac{1}{2} w^t w$

○ **We can formulate a Quadratic Optimization Problem and solve for w and b**

○

**Minimize** $\qquad \Phi(w) = \frac{1}{2} w^t w$

**subject to** $\qquad y_i(wx_i + b) \geq 1 \qquad \forall i$

---

# SVM Formulation
## -- Problem: overfitting in noisy case

- • denotes +1
- ○ denotes -1

○ **Hard Margin: So far we require all data points be classified correctly**

   **- No training error**

○ **What if the training set is noisy?**
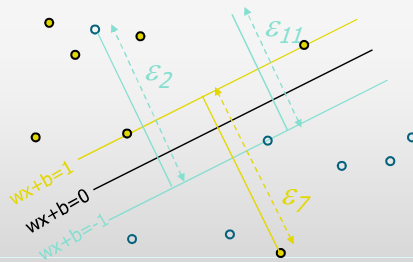
   **- Solution 1: use very powerful kernels**

**OVERFITTING!**

# SVM Formulation
# -- Introducing a soft margin

○ Non-separable classifier

***Slack variables* $\xi_i$ can be added to allow**
**misclassification of difficult or noisy examples.**



What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2}\mathbf{w}.\mathbf{w} + C\sum_{k=1}^{R}\varepsilon_k$$

# SVM Formulation
# -- Maximizing the soft margin

○ Hard Margin vs. Soft Margin

- **The old formulation:**

  Find $\mathbf{w}$ and $b$ such that
  $\Phi(\mathbf{w}) = $½ $\mathbf{w}^T\mathbf{w}$ is minimized and for all $\{(\mathbf{x_i}, y_i)\}$
  $y_i(\mathbf{w^Tx_i} + b) \geq 1$

- **The new formulation incorporating slack variables:**

  Find $\mathbf{w}$ and $b$ such that
  $\Phi(\mathbf{w}) = $½ $\mathbf{w}^T\mathbf{w} + C\Sigma\xi_i$   is minimized and for all $\{(\mathbf{x_i}, y_i)\}$
  $y_i(\mathbf{w^Tx_i} + b) \geq 1 - \xi_i$   and   $\xi_i \geq 0$ for all $i$

- **Parameter $C$ can be viewed as a way to control overfitting.**

# SVM Formulation
## -- Maximizing the soft margin (cont'd)

**The following Lagrangian should be considered**

$$L(w,b,\xi;\ \alpha,v) = \frac{1}{2}w^T w + C\sum_{K=1}^{N}\xi_k - \sum_{k=1}^{N}\alpha_k(y_k[w^T x_k + b] - 1 + \xi_k) + \sum_{k=1}^{N}v_k\xi_k$$

where $\alpha_k \geq 0,\ v_k \geq 0$ for $k = 1,...,N$

- The solution is given by the saddle point of the Lagrangian

$$\max_{\alpha,v}\ \min_{w,b,\xi}(L(w,b,\xi;\alpha,v))$$

$$\begin{cases} \frac{\partial L}{\partial w} = 0 \rightarrow w = \sum_{K=1}^{N}\alpha_k y_k x_k \\ \frac{\partial L}{\partial b} = 0 \rightarrow \sum_{k=1}^{N}\alpha_k y_k = 0 \\ \frac{\partial L}{\partial \xi} = 0 \rightarrow 0 \leq \alpha_k \leq C, k = 1,...,N \end{cases}$$

# SVM Formulation
## -- Maximizing the soft margin (dual form)

- Linear SVMs:  Overview

- **The classifier is a *separating hyperplane*.**
- **Quadratic optimization algorithms can identify which training points $x_i$ are support vectors with non-zero Lagrangian multipliers $\alpha_i$.**

**Find $\alpha_1...\alpha_N$ such that**
**$Q(\alpha) = \Sigma\alpha_i - \frac{1}{2}\Sigma\Sigma\alpha_i\alpha_j y_i y_j x_i^T x_j$ is maximized and**
**(1)  $\Sigma\alpha_i y_i = 0$**
**(2)  $0 \leq \alpha_i \leq C$ for all $\alpha_i$**

SVM
An margin optimized
linear classifier

**$f(x) = \Sigma\alpha_i y_i x_i^T x + b$**

# SVM Extension
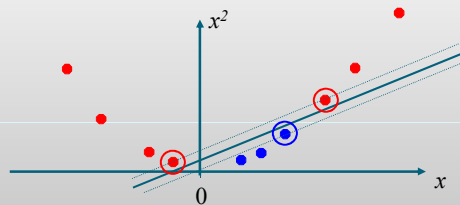## -- classification on different feature spaces

- Nonlinear classifier

- Datasets that are linearly separable with some noise work out great:

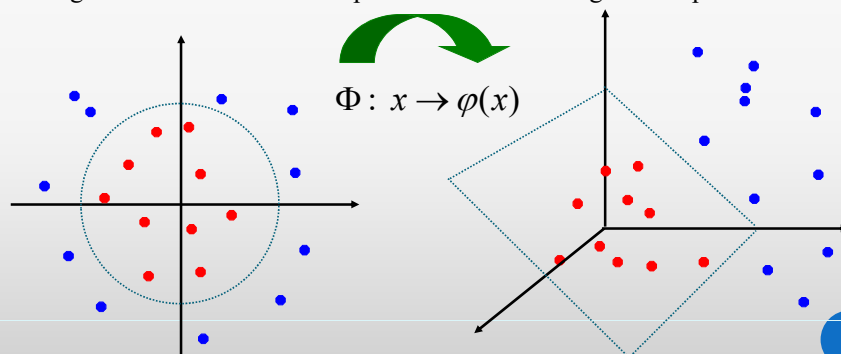- But what are we going to do if the dataset is just too hard?

- How about… mapping data to a higher-dimensional space:

# SVM Extension
## -- Mapping to higher-dimensional feature space

- Nonlinear classifier

- General idea:  the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

$$\Phi : x \rightarrow \varphi(x)$$

15

# SVM Extension: nonlinear classifier
# -- Kernel tricky

◆ **Find** $\alpha_1 \dots \alpha_N$ **such that**

$$\max_{\alpha_1 \dots \alpha_N} Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

subject to  (1) $\sum_{i=1}^{N} \alpha_i y_i = 0$ ,  (2) $0 \leq \alpha_i \leq C$  for  $\forall \alpha_i$

◆ The kernel function

$$K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$$

◆ *Decision value* is

$$v = \sum_{i=1}^{N} \alpha_i y_i K(x, x_i) + b$$

◆ *Classifier* is  $y(x) = \text{sign}(v)$

# SVM classifier
# -- RBF network with optimized margin