

Supervised learning

- Multi-class
- Support vector machines
- Neural networks
- Feature selection

Multi-class classification

Classification task with K classes $\{C_1, C_2, \dots, C_K\}$

Currently, there is no clear solution to the problem of multi-class classification

3 approaches:

- One vs All
- All vs All
- Hierarchical classification

One vs All

K classifiers: a class VS all other classes

$$\left\{ \begin{array}{l} C_1 \text{ VS } \{C_2, C_3, \dots, C_K\} \rightarrow \text{score}_1 \\ C_2 \text{ VS } \{C_1, C_3, \dots, C_K\} \rightarrow \text{score}_2 \\ \dots \\ C_i \text{ VS } \{C_1, \dots, C_{i-1}, C_{i+1}, \dots, C_K\} \rightarrow \text{score}_i \\ \dots \\ C_K \text{ VS } \{C_1, C_2, \dots, C_{K-1}\} \rightarrow \text{score}_K \end{array} \right.$$

$$C^* \leftarrow \operatorname{argmax}_i \{\text{score}_1, \text{score}_2, \dots, \text{score}_i, \dots, \text{score}_K\}$$

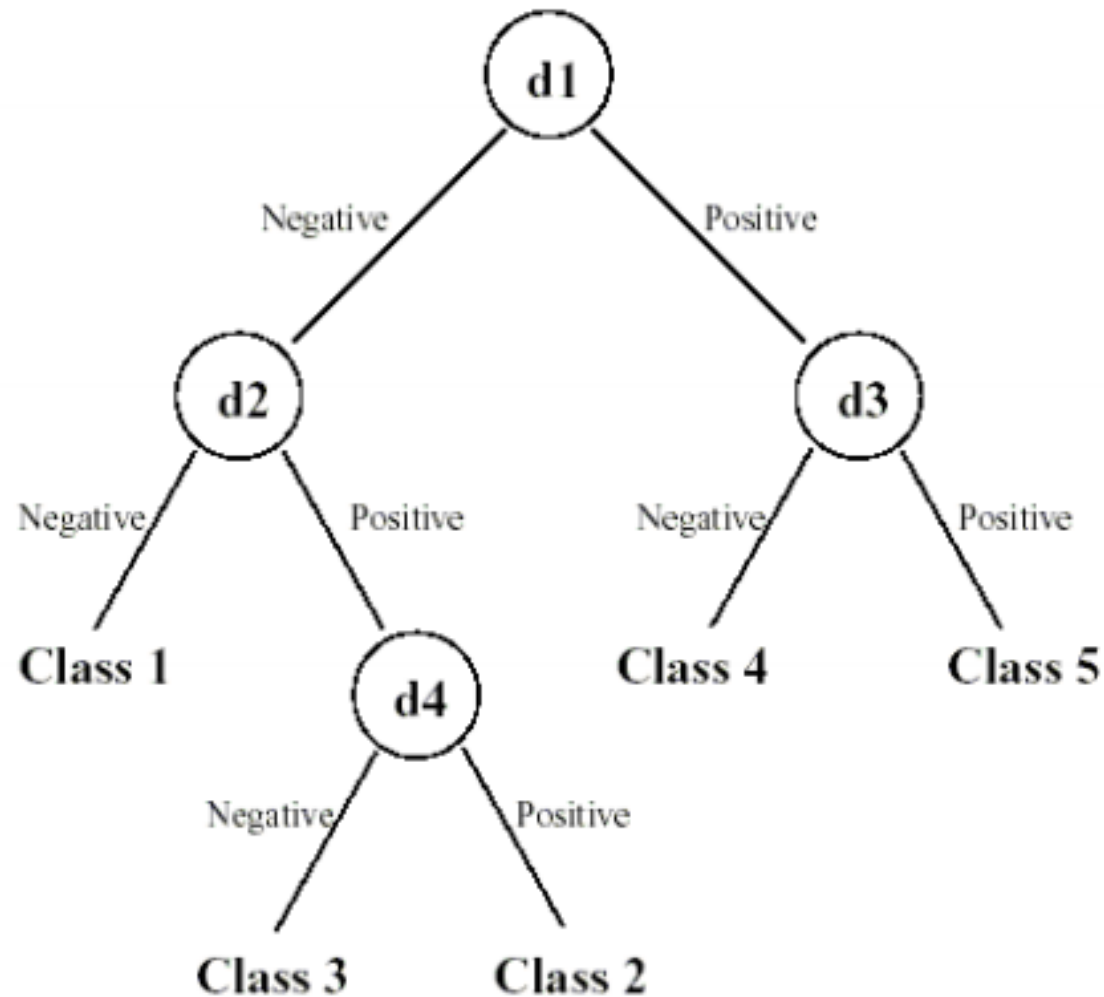
All VS All

$K(K-1)/2$ binary classifiers: C_i VS C_j for all $i > j$

	C_1	C_2	...	C_j	...	C_K	Total
C_1		0		1		0	6
C_2	1			1		0	5
...						0	
C_i	1	1				0	10
...							
C_K	0	0		0			2

$C^* \leftarrow \operatorname{argmax}_i \{ \text{Total}_1, \text{Total}_2, \dots, \text{Total}_i, \dots, \text{Total}_K \}$

Hierarchical classification



Performance of multi-class classifiers

Confusion matrix

	C_1	C_2	...	C_j	...	C_K
C_1	M_{11}	M_{12}		M_{1j}		M_{1K}
C_2	M_{21}	M_{22}				
...						
C_i	M_{j1}	M_{j2}		M_{jj}		M_{jK}
...						
C_K	M_{K1}	M_{K2}		M_{Kj}		M_{KK}

Cost Matrix

	C_1	C_2	...	C_j	...	C_K
C_1	α_{11}	α_{12}		α_{1j}		α_{1K}
C_2	α_{21}	α_{22}				
...						
C_i	α_{j1}	α_{j2}		α_{jj}		α_{jK}
...						
C_K	α_{K1}	α_{K2}		α_{Kj}		α_{KK}

No ROC space can be defined

Evaluation done by classification cost

Support Vector Machine (SVM)

Supervised Learning

Introduction

- Binary classification method
- Introduced by Vladimir Vapnik in 1995
- Competitive in many cases
- Relatively easy to use
- Many extention: Regression, density estimation, kernel PCA,...

Classification problem

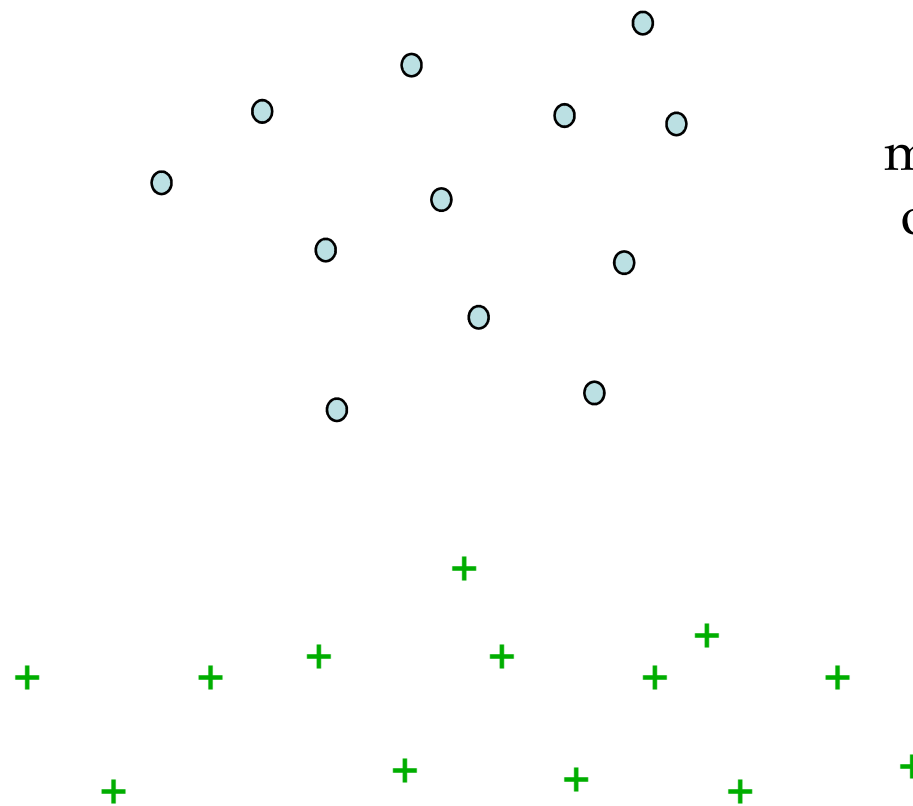
- Training examples: x_i , $i=1,\dots,l$
- Consider a simple case with two classes:

Define a class vector y

$$y_i = \begin{cases} 1 & \text{if } x_i \text{ in class } P \\ -1 & \text{if } x_i \text{ in class } N \end{cases}$$

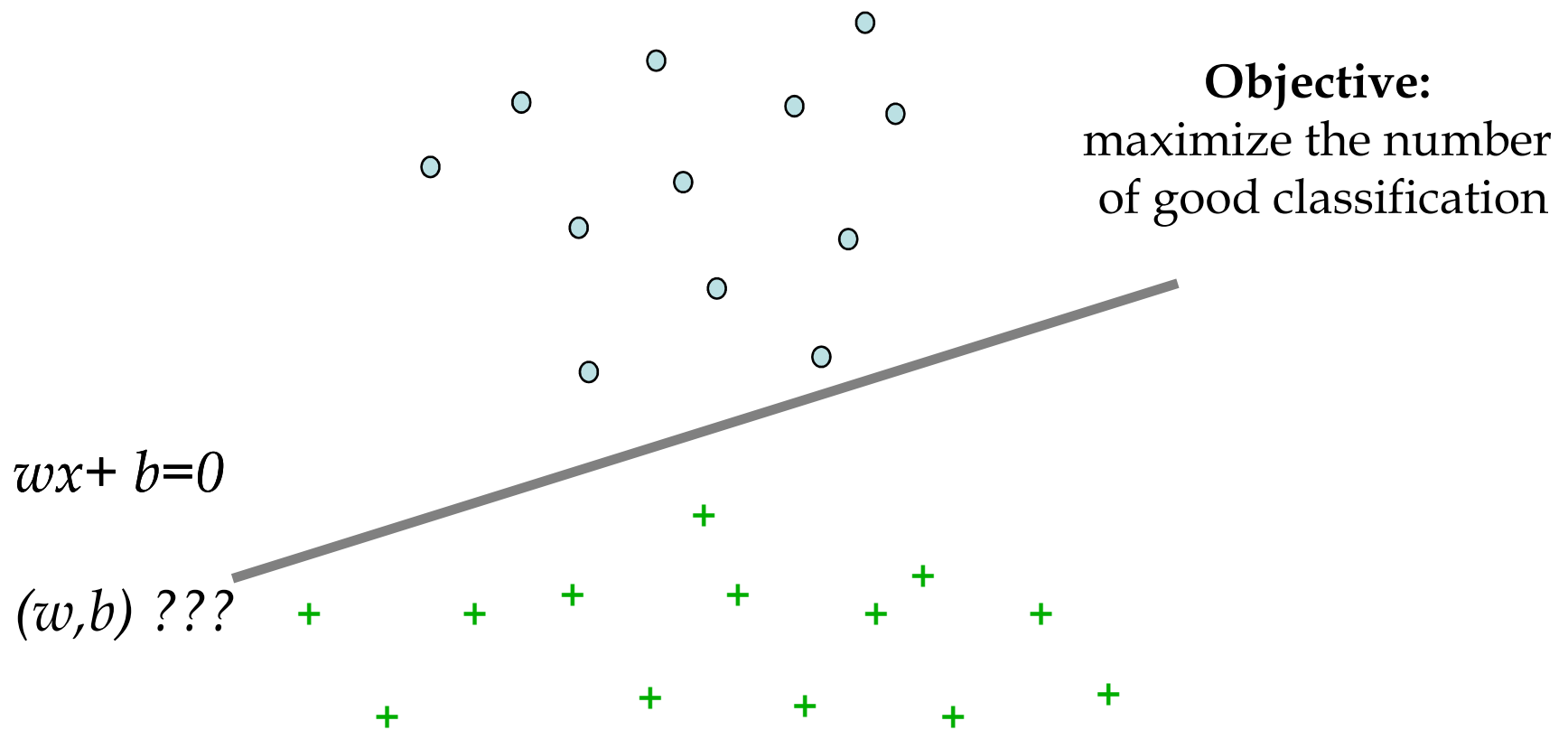
- Classifier: hyperplane that separates all examples

2-class problem

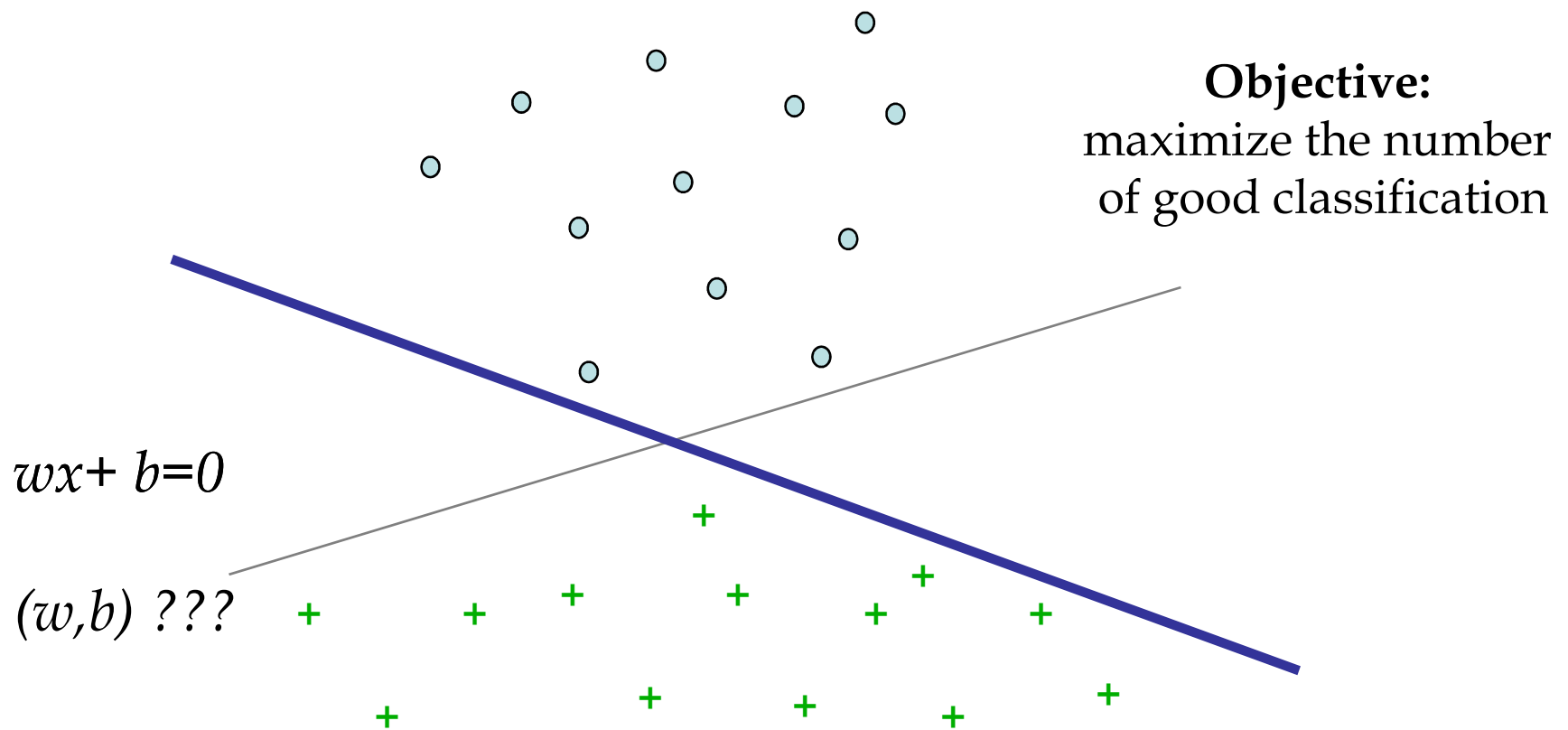


Objective:
maximize the number
of good classification

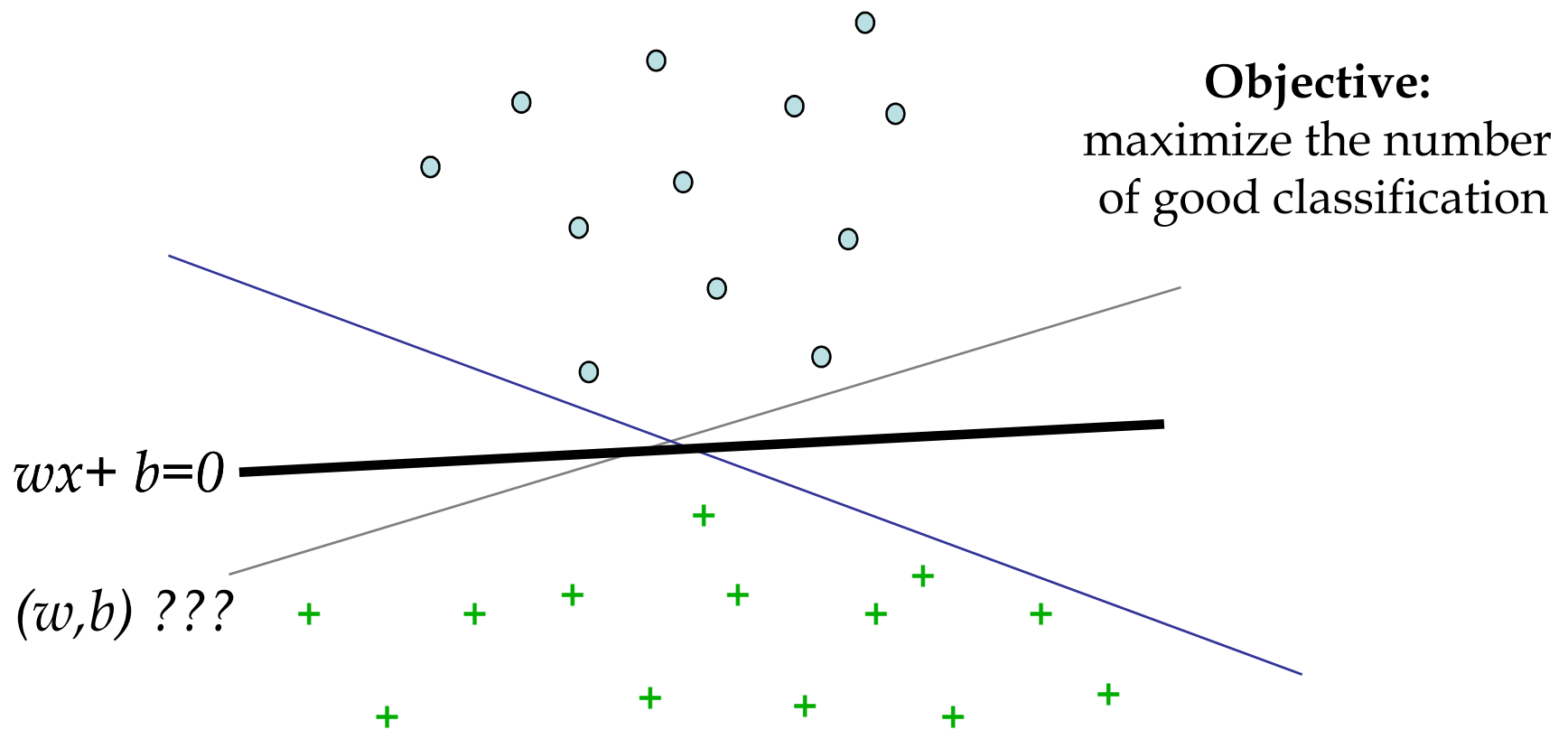
2-class problem



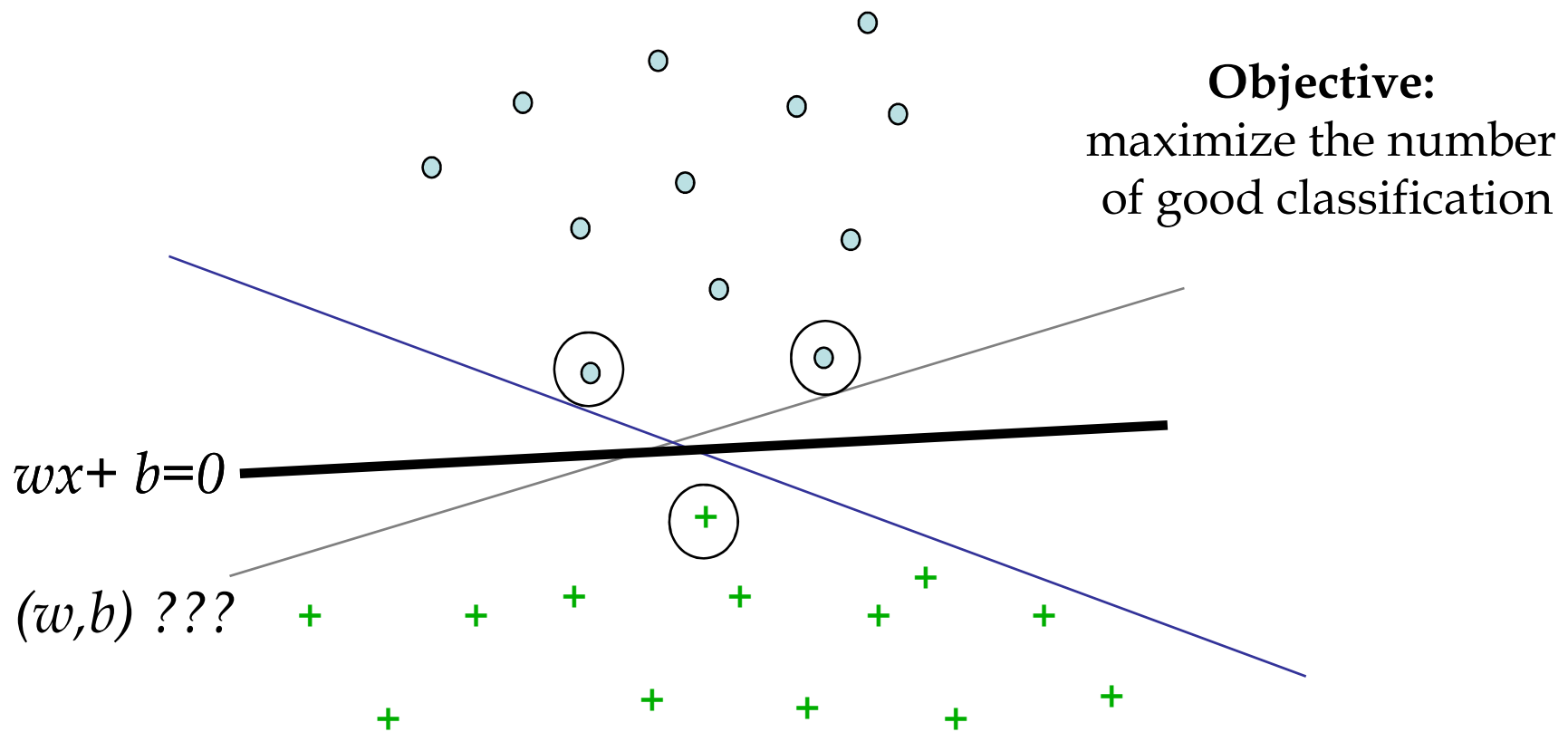
2-class problem



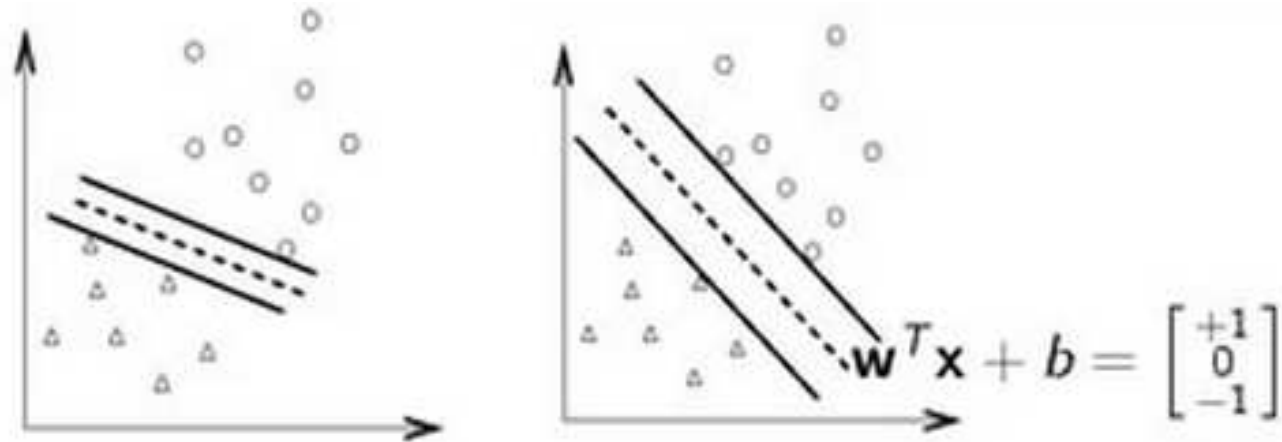
2-class problem



Particular examples in classification



2-class problem



- A separating hyperplane: $\mathbf{w}^T \mathbf{x} + b = 0$

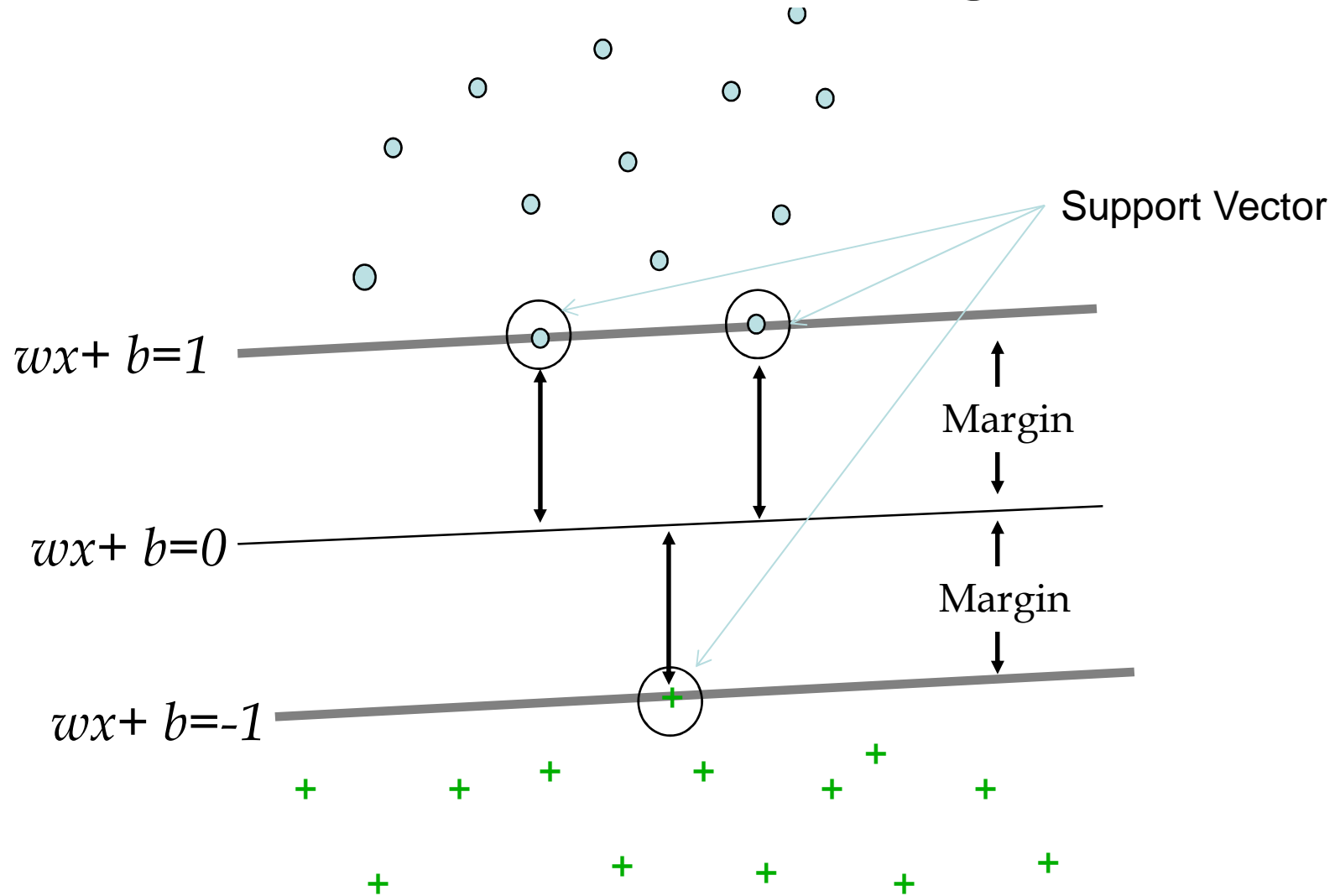
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) + b &> 0 & \text{if } y_i = 1 \\ (\mathbf{w}^T \mathbf{x}_i) + b &< 0 & \text{if } y_i = -1 \end{aligned}$$

- Decision function $f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$, \mathbf{x} : test data

Many possible choices of \mathbf{w} and b



Maximal Margin



Maximal Margin

- Margin = distance between the support vector lines
- Linear $\rightarrow f(x) = w \cdot x + b = 0$
- Distance between a point and the hyperplane separator
$$d(x) = |w \cdot x + b| / ||w||$$
- Margin size: $2 / ||W||$
- Maximize the margin = minimize $||w||$

Maximal Margin

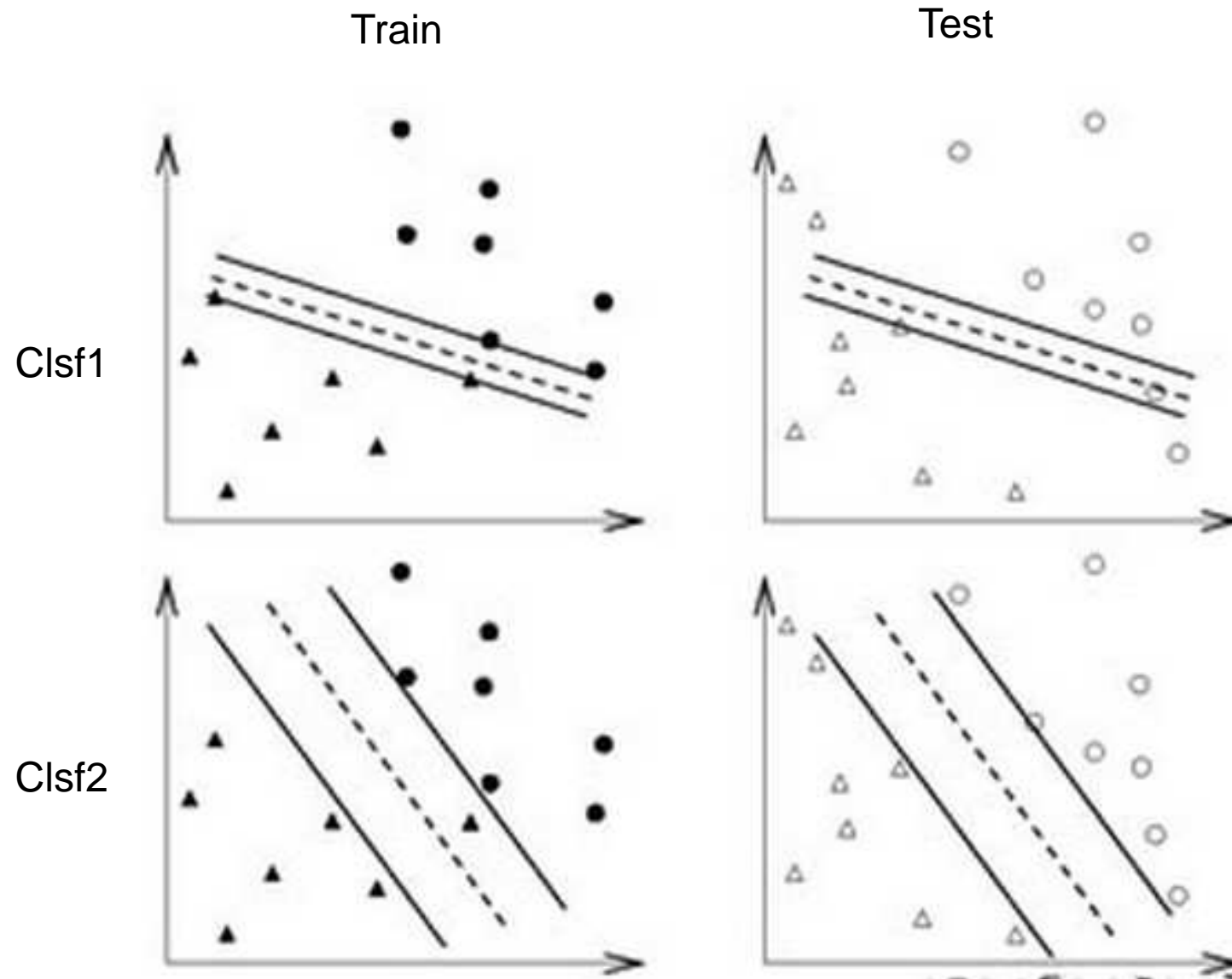
- Distance between $\mathbf{w}^T \mathbf{x} + b = 1$ and -1 :

$$2/\|\mathbf{w}\| = 2/\sqrt{\mathbf{w}^T \mathbf{w}}$$

- A **quadratic programming** problem
[Boser et al., 1992]

$$\begin{array}{ll} \min_{\mathbf{w}, b} & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \\ & i = 1, \dots, l. \end{array}$$

Problem: not optimal

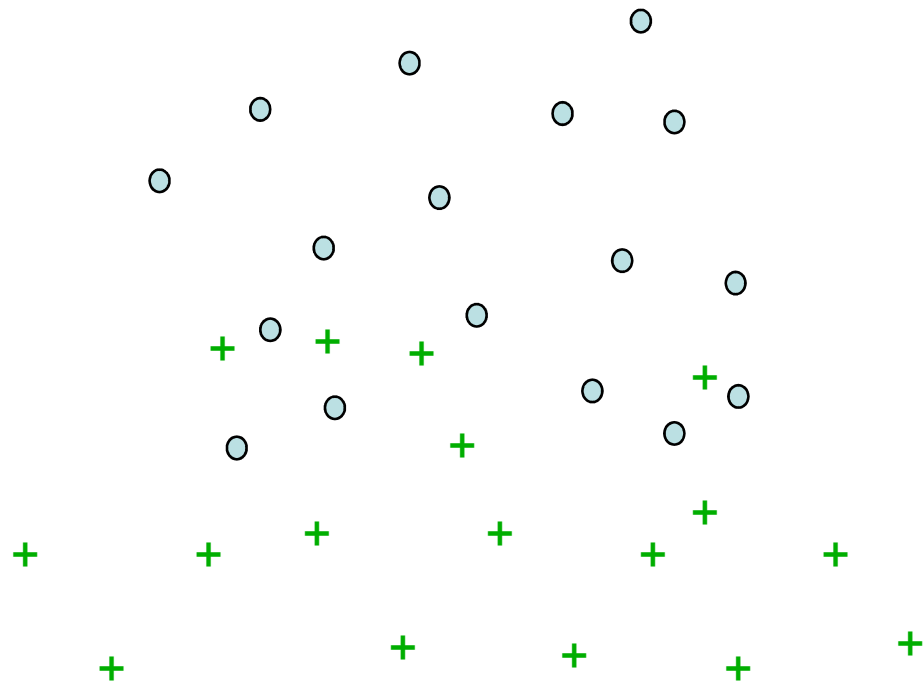


Slacking Variables

- It may be not possible or not optimal to separate perfectly the two classes in training set.
- We add slacking variables to make the model tolerant to the errors.
- C defines the tradeoff between margin and constraint.

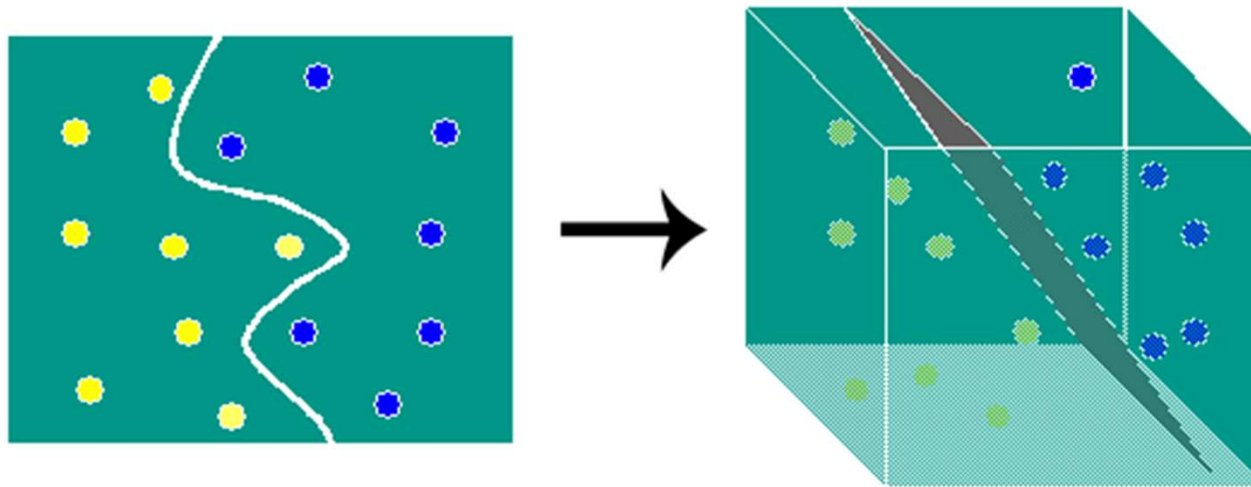
$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, \dots, l. \end{aligned}$$

Problem: not linearly separable



Kernel Trick

- Solving a non linear problem ?
- We change the dimensionality in using a kernel :
- The data are projected in a new higher dimension space
- A linear classifier is consructed in this kernel space



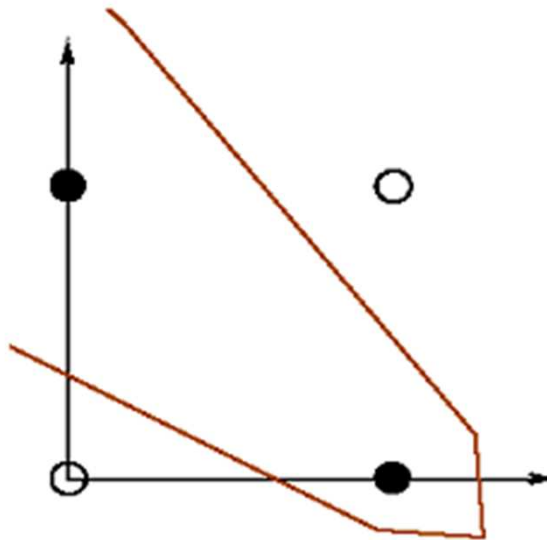
Kernel Trick: example

Example of the XOR problem:

The XOR problem is not linearly separable.

Let see in a 2-D space:

Points : $(0,0)$; $(0,1)$; $(1,0)$; $(1,1)$



Kernel Trick: example

We take a polynomial kernel : $(x, y) \rightarrow (x, y, x \cdot y)$

We project the example from a 2-D space to a 3-D space

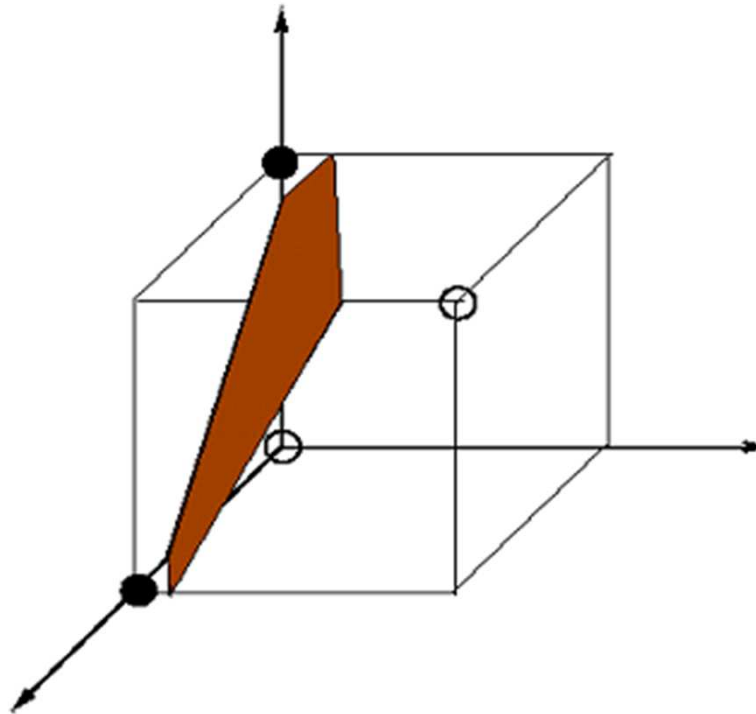
The class are linearly separable in this new space:

$(0,0) \rightarrow (0,0,0)$

$(0,1) \rightarrow (0,1,0)$

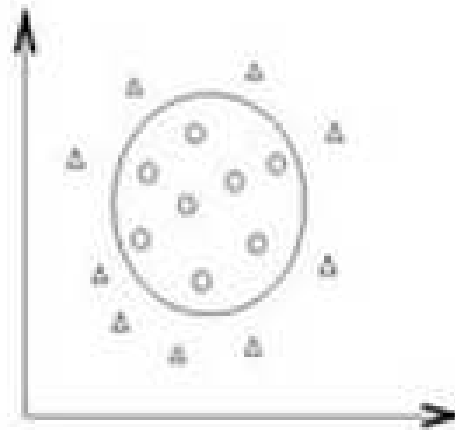
$(1,0) \rightarrow (1,0,0)$

$(1,1) \rightarrow (1,1,1)$



Kernel Trick

- An example:



- Allow training errors
- Higher dimensional (maybe infinite) feature space

$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots).$$

SVM Formulation

- Standard SVM [Cortes and Vapnik, 1995]

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, \dots, l. \end{aligned}$$

- Example: $\mathbf{x} \in R^3, \phi(\mathbf{x}) \in R^{10}$

$$\begin{aligned} \phi(\mathbf{x}) = & (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, x_1^2, \\ & x_2^2, x_3^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3) \end{aligned}$$

Finding the decision function

- \mathbf{w} : maybe **infinite** variables
- The **dual** problem

$$\begin{array}{ll}\min_{\alpha} & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} & 0 \leq \alpha_i \leq C, i = 1, \dots, l \\ & \mathbf{y}^T \alpha = 0,\end{array}$$

where $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- A **finite** problem: #variables = #training data

Decision function

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- Decision function

$$\begin{aligned} & \mathbf{w}^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^l \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \end{aligned}$$

- Only $\phi(\mathbf{x}_i)$ of $\alpha_i > 0$ used \Rightarrow **support vectors**

Kernel function

- The problem and the solution depend only on the dot product $\Phi(x) \cdot \Phi(x')$
- In practice we do not choose Φ but K
 $K: X \times X \rightarrow R$ such that $k(x, x') = \Phi(x) \cdot \Phi(x')$
- K is called kernel function
- Pb: How to choose the kernel function

Mercer Condition

- K is a kernel $\Leftrightarrow (k(x_i, x_j))_{i,j}$ is a definite positive matrix
- K is symmetric
- In this case the function Φ tq: $k(x, x') = \Phi(x) \cdot \Phi(x')$ is defined
- Probleme:
 - These conditions are difficult to check
 - Φ stays unknown
 - Do not help for the choice of the kernel

Kernel example

General Kernel:

- Linéaire

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x} \cdot \mathbf{x}'$$

- Polynomial

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}')^d \text{ ou } (c + \mathbf{x} \cdot \mathbf{x}')^d$$

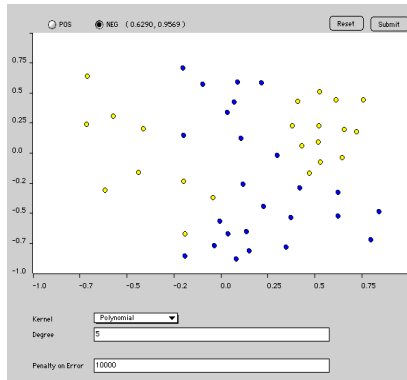
- Gaussien

$$k(\mathbf{x}, \mathbf{x}') = e^{-\|\mathbf{x} - \mathbf{x}'\|^2 / \sigma}$$

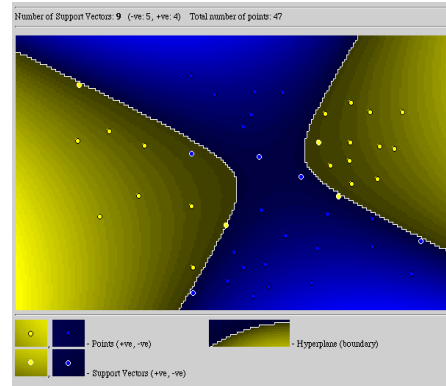
- Laplacien

$$k(\mathbf{x}, \mathbf{x}') = e^{-\|\mathbf{x} - \mathbf{x}'\|_1 / \sigma}$$

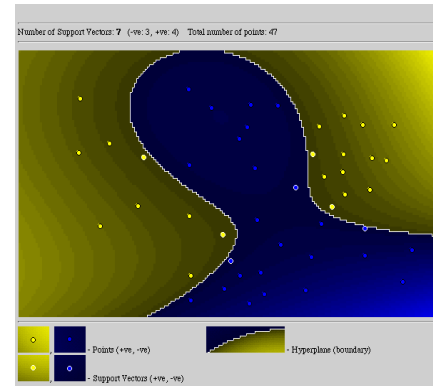
Kernel example



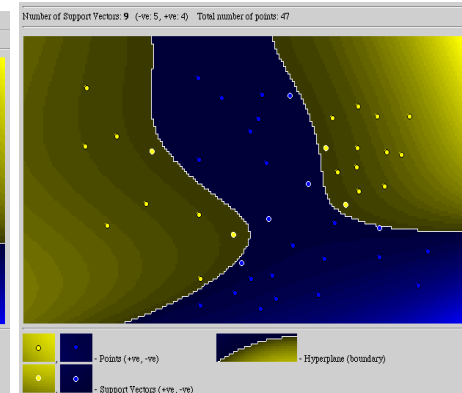
- 47 examples
(22 +, 25 -)



(5-, 4+)

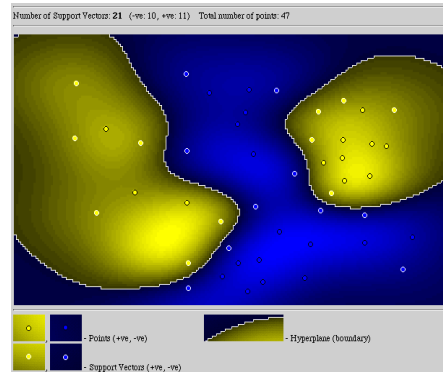


(3-, 4+)

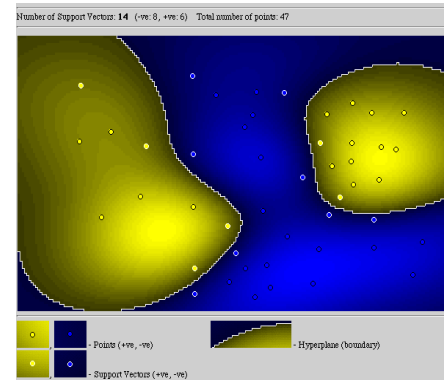


(5-, 4+)

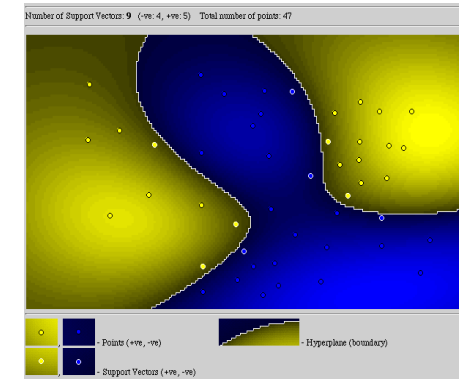
Polynomial function, degree=2, 5, 8 and $C = 10000$



(10-, 11+)



(8-, 6+)



(4-, 5+)

Gaussian function $\sigma = 2, 5, 10$ and $C = 10000$

Selecting Kernels

- For beginners, use Gaussian (RBF) first
- Linear Kernel: special case of RBF
- Polynomial: numerical difficulties

More parameters than RBF

- Researchers design many kernels specialized to different domains
- We can combine different kernel to construct a more complex kernel

Parameters selection

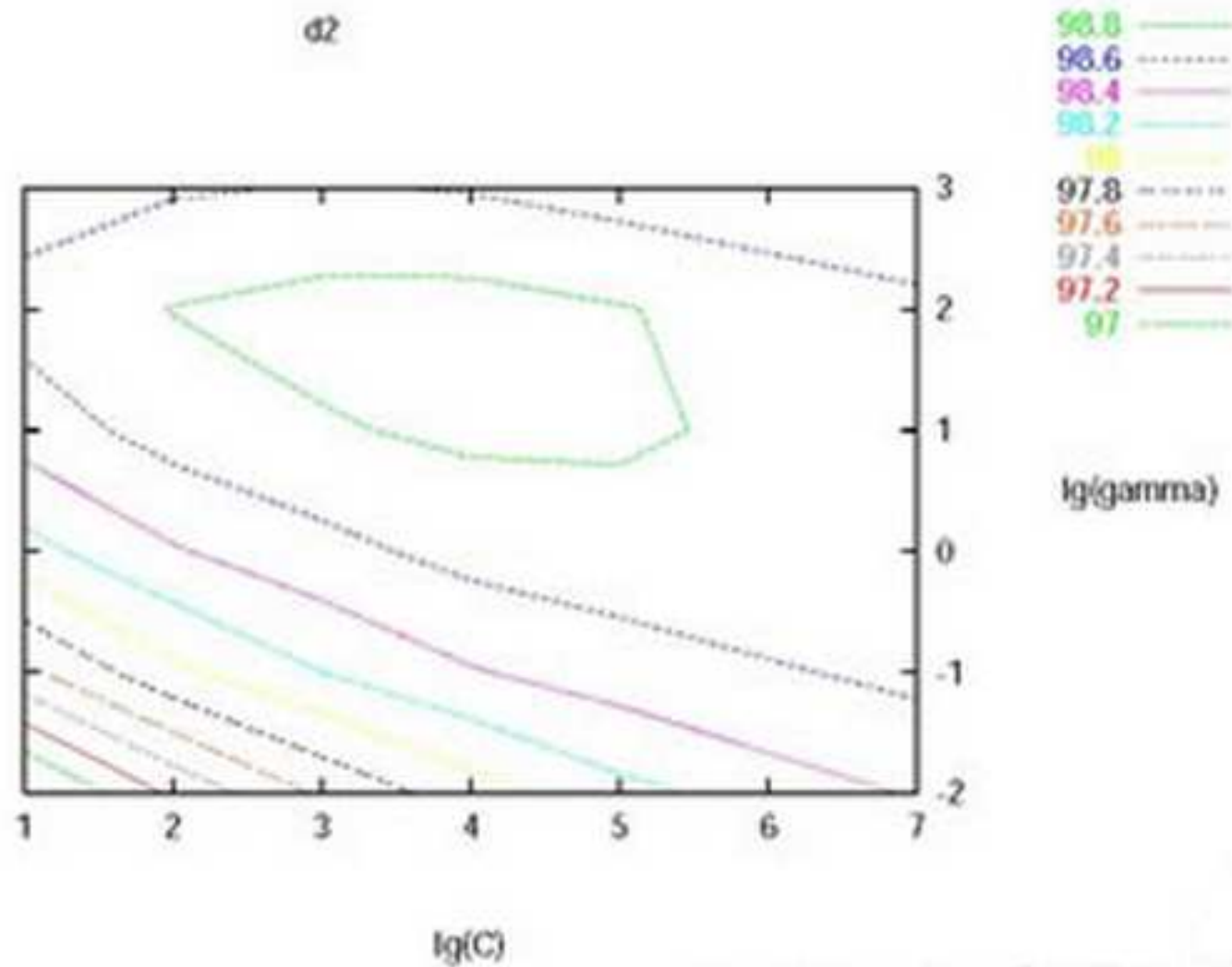
- Is important
- Now parameters are
C, kernel parameters
- Example:

$$\gamma \text{ of } e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}$$

$$a, b, d \text{ of } (\mathbf{x}_i^T \mathbf{x}_j / a + b)^d$$

- How to select them?
So performance better?

Parameters selection



Parameters selection

- In practice
Available data \Rightarrow training and validation
- Train the training
- Test the validation
- k -fold cross validation:
Data randomly separated to k groups
Each time $k - 1$ as training and one as testing
- Using CV on training + validation
- Predict testing with the best parameters from CV

Data scaling

- Problem: Features in greater numeric range may dominate the construction of the separator
- Scale the features to the range [0,1]

$$x_i^j = \frac{x_i^j - x_i^{\min}}{x_i^{\max} - x_i^{\min}}$$

- Scaling generally helps but not always
- Same scaling factor have to be used on training and testind set

A simple procedure

- Conduct simple **scaling** on the data
- Consider **RBF** kernel $K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$
- Use cross-validation to find the **best parameter** C and γ
- Use the best C and γ to **train the whole** training set
- Test

For beginners only, you can do a lot more

Computing time

- n = Number of training example
- d = Number of features
- $dn^2 \leq \text{complexity} \leq dn^3$
- Kernel matrix size= n^2
- \Rightarrow SVM can deal with problem up to 100.000 examples

SVM with R

- Function `svm()` du package `e1071`

`svm(x, y = NULL, scale = TRUE, type = NULL, kernel = "radial", degree = 3, coef0 = 0, cost = 1, nu = 0.5, class.weights = NULL, cachesize = 40, tolerance = 0.001, epsilon = 0.1, shrinking = TRUE, cross = 0, probability = FALSE, fitted = TRUE, ..., subset, na.action = na.omit)`

SVM with R

x: a data matrix.

y: a response vector with one label for each row/component of x

Scale: A logical vector indicating the variables to be scaled. If scale is of length 1, the value is recycled as many times as needed. Per default, data are scaled internally (both x and y variables) to zero mean and unit variance.

Type: svm can be used as a classification machine, as a regression machine. (C-classification,eps-regression)

SVM with R

Kernel: the kernel used in training and predicting.

- linear: $u'v$
- polynomial: $(\text{gamma} * u'v + \text{coef0})^{\text{degree}}$
- radial basis: $\exp(-\text{gamma} * |u-v|^2)$
- sigmoid: $\tanh(\text{gamma} * u'v + \text{coef0})$

Degree, gamma, coef0: parameter needed for kernel

Cost: cost of constraints violation (default: 1)—it is the 'C'-constant of the regularization term in the Lagrange formulation.

SVM with R

Class.weights: a named vector of weights for the different classes, used for asymmetric class sizes.

Cross: if an integer value $k > 0$ is specified, a k -fold cross validation on the training data is performed to assess the quality of the model

Probability: logical indicating whether the model should allow for probability predictions.

Demo: <http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml>