

# SY19 – Machine Learning

## Chapter 8: Support Vector Machines

Thierry Denœux

Université de technologie de Compiègne

<https://www.hds.utc.fr/~tdenoieux>  
email: tdenoeux@utc.fr

Automne 2024

The screenshot shows a presentation slide with the following elements:

- Header: Thierry Denœux, SY19 – Support Vector Machines
- Page number: A24, 1 / 93
- Content: Optimal Separating hyperplane
- Navigation: Standard presentation navigation icons.

### Overview

#### 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

#### 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

#### 3 Support Vector Regression

- Formalization
- Solution and interpretation

# Support Vector Machines

- In this chapter we describe new methods for classification and regression.
- Optimal separating hyperplanes are first introduced for binary classification in the case where two classes are linearly separable.
- Then we cover extensions to the nonseparable case, where the two classes overlap.
- These techniques are then generalized to the support vector machine (SVM), which produces nonlinear boundaries by constructing a linear boundary in a large, transformed version of the predictor space.
- Finally, we adapt the approach to regression.

The screenshot shows a presentation slide with the following elements:

- Header: Thierry Denœux, SY19 – Support Vector Machines
- Page number: A24, 2 / 93
- Content: Optimal Separating hyperplane | Formalization
- Navigation: Standard presentation navigation icons.

### Overview

#### 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

#### 2 Support Vector Machines

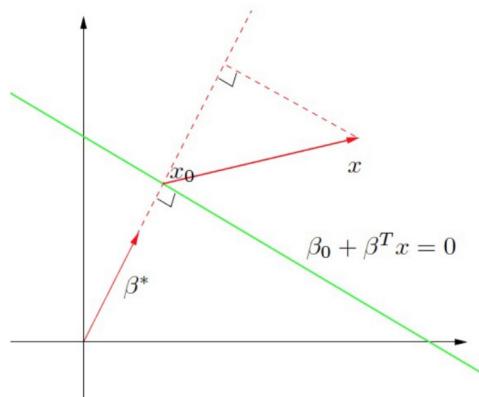
- The kernel trick
- Kernel functions
- Remarks, variants and extensions

#### 3 Support Vector Regression

- Formalization
- Solution and interpretation

## Hyperplane

In  $\mathbb{R}^p$ , a **hyperplane**  $H$  is a  $p - 1$  dimensional subspace defined by the equation  $g(x) = 0$  with  $g(x) = \beta_0 + \beta^T x$ . We have  $g(x) > 0$  on one side of  $H$  and  $g(x) < 0$  on the other side.



For any two points  $x_1$  and  $x_2$  lying in  $H$ , we have

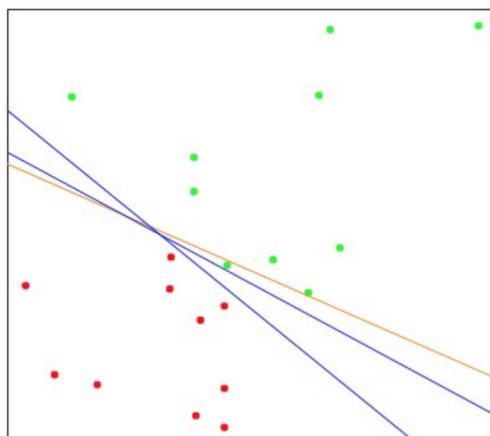
$$\beta_0 + \beta^T x_1 = 0$$

$$\beta_0 + \beta^T x_2 = 0$$

Consequently,  $\beta^T(x_1 - x_2) = 0$ , hence  $\beta^* = \beta / \|\beta\|$  is the unit vector normal to the surface of  $H$ .



## Linearly separable data



- Consider a two-class data set  $\{(x_i, y_i)\}_{i=1}^n$  with  $y_i \in \{-1, 1\}$ .
- It is said to be **linearly separable** if there exists a hyperplane  $H : g(x) = 0$  that separates the two classes, i.e., such that

$$g(x_i)y_i > 0, \quad \forall i$$

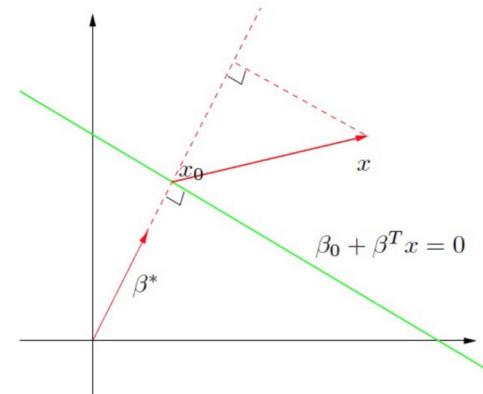


## Signed distance to a hyperplane

Let  $x_0 \in H$ . The **s**igned **d**istance of any point  $x$  to  $H$  is

$$d_s(x, H) = \beta^{*T}(x - x_0)$$

$$= \frac{\beta^T x - \beta^T x_0}{\|\beta\|}$$



As  $\beta_0 = -\beta^T x_0$ , we have

$$d_s(x, H) = \frac{\beta^T x + \beta_0}{\|\beta\|} = \frac{g(x)}{\|\beta\|}$$



## Optimal separating hyperplane

Let  $H : g(x) = 0$  be a separating hyperplane. The distance between  $H$  and a learning vector  $x_i$  is

$$d(x_i, H) = \frac{g(x_i)y_i}{\|\beta\|}$$

### Definition (Margin)

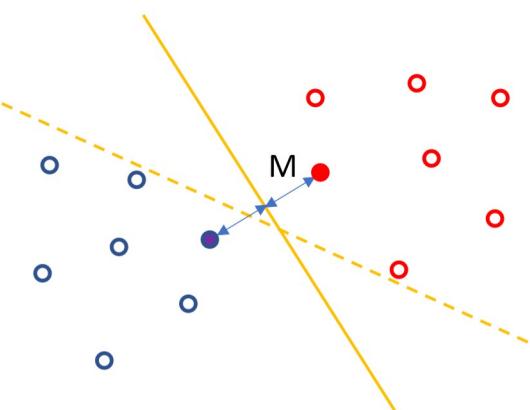
The **margin** of  $H$  is the smallest distance between  $H$  and a learning vector  $x_i$ :

$$M = \min_i d(x_i, H)$$

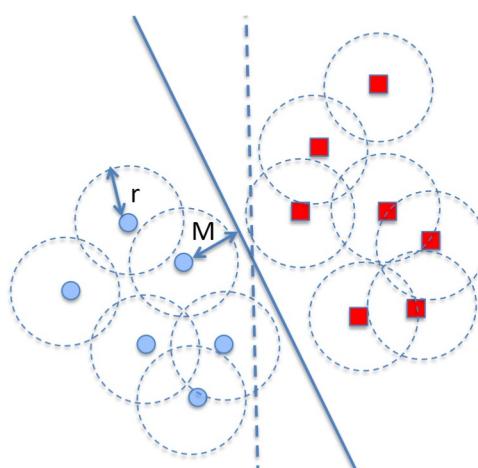
### Definition (Optimal separating hyperplane, support vectors)

The **optimal separating hyperplane** (OSH) is the separating hyperplane with the largest margin. The learning vectors  $x_i$  such that  $d(x_i, H) = M$  are called the **support vectors** (SVs) of  $H$ .

## Example 1

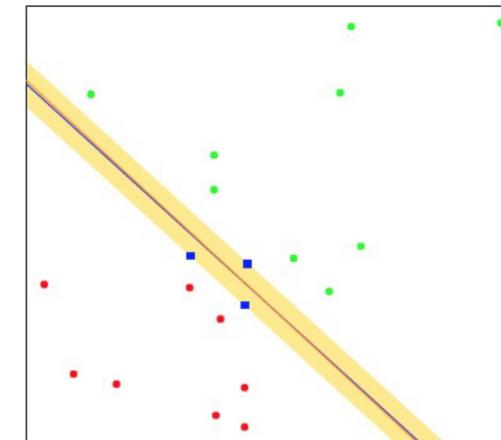


The OSH is more likely to separate future data



- Future data can be assumed to be “close” to past data.
- Assume they will lie within a distance  $r$  to a past data point.
- If  $M > r$ , the hyperplane will classify future data perfectly.

## Example 2



The shaded region delineates the maximum margin separating the two classes. There are 3 SVs, and the OSH is the blue line. The boundary found using logistic regression is the red line. In this case, it is very close to the OSH.

## Overview

- ➊ Optimal Separating hyperplane
  - Formalization
  - Solution in the separable case
  - Non-separable case
- ➋ Support Vector Machines
  - The kernel trick
  - Kernel functions
  - Remarks, variants and extensions
- ➌ Support Vector Regression
  - Formalization
  - Solution and interpretation

## How to find the OSH?

- The OSH can be found by solving the following optimization problem:

$$\max_{\beta, \beta_0} M$$

subject to  $\frac{y_i(\beta^T x_i + \beta_0)}{\|\beta\|} \geq M, \quad i = 1, \dots, n$

- If  $(\beta, \beta_0)$  is a solution, so is  $(\lambda\beta, \lambda\beta_0)$  for any  $\lambda$ . Hence, we can fix  $\|\beta\| = 1/M$  and reformulate the problem as

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$$

subject to  $y_i(\beta^T x_i + \beta_0) \geq 1, \quad i = 1, \dots, n$



## Reminder on constrained optimization

### Lagrange function

Consider the following minimization problem:

$$\min_{\beta} f(\beta) \tag{1}$$

subject to the constraints  $c_i(\beta) \geq 0, i = 1, \dots, n$ , where  $f$  and the  $c_i$ 's are differentiable functions.

### Definition (Lagrange function)

The Lagrange function is defined by

$$L(\beta, \alpha) = f(\beta) - \sum_{i=1}^n \alpha_i c_i(\beta)$$

where  $\alpha = (\alpha_1, \dots, \alpha_n)$  is the vector of Lagrange multipliers.

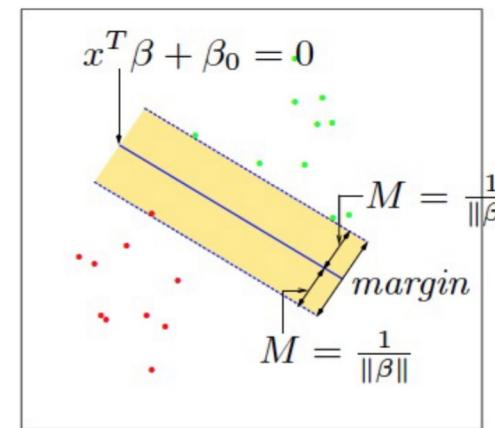
## Interpretation

- The constraints define an empty band or margin around the linear decision boundary of thickness  $1/\|\beta\|$ .

- The vectors  $x_i$  such that

$$\frac{y_i(\beta^T x_i + \beta_0)}{\|\beta\|} = \frac{1}{\|\beta\|}$$

i.e.,  $y_i(\beta^T x_i + \beta_0) = 1$ , are the SVs.



## Reminder on constrained optimization

### Karush-Kuhn-Tucker conditions

### Theorem (Karush-Kuhn-Tucker)

If function  $f$  has a minimum for some value  $\beta^*$  in the feasibility region, the following Karush-Kuhn-Tucker (KKT) conditions are verified for some vector  $\alpha^* = (\alpha_1^*, \dots, \alpha_n^*)$ :

$$\frac{\partial L}{\partial \beta}(\beta^*, \alpha^*) = 0 \tag{2a}$$

$$c_i(\beta^*) \geq 0, \quad i = 1, \dots, n \tag{2b}$$

$$\alpha_i^* c_i(\beta^*) = 0, \quad i = 1, \dots, n \tag{2c}$$

$$\alpha_i^* \geq 0, \quad i = 1, \dots, n \tag{2d}$$

Remark: if  $\alpha_i^* > 0$ , then  $c_i(\beta^*) = 0$ : constraint  $i$  is active.



## Reminder on constrained optimization

Wolfe dual

### Theorem (Wolfe dual)

Problem (1) is equivalent to the following problem (Wolfe dual):

$$\max_{\beta, \alpha} L(\beta, \alpha) \quad (3)$$

subject to

$$\frac{\partial L}{\partial \beta} = 0 \quad (4)$$

$$\alpha_i \geq 0 \quad i = 1, \dots, n \quad (5)$$

### Lagrangian of the dual problem

Substituting (7) and (8) in (6), we get

$$L_D(\alpha) = \underbrace{\frac{1}{2} \left( \sum_{i=1}^n \alpha_i y_i x_i \right)^T \left( \sum_{j=1}^n \alpha_j y_j x_j \right)}_{\sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j} - \underbrace{\sum_i \alpha_i y_i \left( \sum_{j=1}^n \alpha_j y_j x_j \right)^T x_i}_{\sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j} - \underbrace{\beta_0 \sum_i \alpha_i y_i}_{0} + \sum_i \alpha_i$$

which can be written as

$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

### Lagrange function

- Let us come back to the problem  $\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2$  subject to  $y_i(\beta^T x_i + \beta_0) \geq 1, i = 1, \dots, n$ .
- This is a **convex optimization problem** (quadratic criterion with linear inequality constraints), so the solution exists and it is unique.
- The Lagrange function is

$$L(\beta, \beta_0, \alpha) = \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^n \alpha_i [y_i(\beta^T x_i + \beta_0) - 1] \quad (6)$$

- Setting the derivatives to zero, we obtain:

$$\frac{\partial L}{\partial \beta} = \beta - \sum_{i=1}^n \alpha_i y_i x_i = 0 \Rightarrow \beta = \sum_{i=1}^n \alpha_i y_i x_i \quad (7)$$

$$\text{and } \frac{\partial L}{\partial \beta_0} = -\sum_{i=1}^n \alpha_i y_i = 0$$

### Solving the dual problem

- The solution is obtained by maximizing  $L_D(\alpha)$  subject to the constraints

$$\alpha_i \geq 0 \text{ and } \sum_{i=1}^n \alpha_i y_i = 0 \quad (9)$$

- This can be done using standard **quadratic programming** software. We will discuss a specialized optimization algorithm later.

## Interpreting the solution

### Support vectors

- The solution  $\alpha^*$  must satisfy the KKT conditions, which include (7), (8), (9) and

$$\alpha_i^*[y_i(\beta^{*T}x_i + \beta_0^*) - 1] = 0, \quad i = 1, \dots, n \quad (10)$$

- From these we can see that, if  $\alpha_i^* > 0$ , then  $y_i(\beta^{*T}x_i + \beta_0^*) = 1$ , i.e.,  $x_i$  is a SV.
- The SVs are the input vectors  $x_i$  such that  $\alpha_i^* > 0$ .

## SVM classifier

- The equation of the OSH is

$$g^*(x) = \beta^{*T}x + \beta_0^* = \sum_{i \in \mathcal{S}} \alpha_i^* y_i x_i^T x + \beta_0^* = 0$$

- The corresponding classifier is

$$D(x) = \text{sign } g^*(x)$$

- The classifier is based only on SVs, which are close to the boundary between classes.

## Interpreting the solution

### Computing $\beta^*$ and $\beta_0^*$

- From (7) we see that the solution vector  $\beta^*$  is defined in terms of a linear combination of the SVs:

$$\beta^* = \sum_{i=1}^n \alpha_i^* y_i x_i = \sum_{i \in \mathcal{S}} \alpha_i^* y_i x_i \quad (11)$$

with  $\mathcal{S} = \{i : \alpha_i^* > 0\}$ .

- The intercept  $\beta_0^*$  can be found from (10): for any SV  $x_i$ , we have

$$y_i(\beta^{*T}x_i + \beta_0^*) = 1$$

from which we can get  $\beta_0^*$ .

## Overview

### 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

### 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

### 3 Support Vector Regression

- Formalization
- Solution and interpretation

## Extension to non-separable data

- Until now, we have assumed that the data are linearly separable.
- This will generally not be the case with real data, so the technique derived so far is not really useful in practice.
- We need to propose an alternative formulation for the **non-separable case**.

Thierry Denœux SY19 – Support Vector Machines A24 25 / 93

Optimal Separating hyperplane Non-separable case

utc

## Interpretation

- The filled points are on the wrong side of their margin by an amount  $M\xi_i$ .
- Points on the correct side have  $\xi_i = 0$ .
- Misclassified points have  $\xi_i > 1$ .

## Weakening the constraints

- Suppose that the classes overlap in predictor space.
- One way to deal with the overlap is to still maximize the margin  $M$ , but allow for some points to be on the wrong side of the margin.
- Define the **slack variables**  $\xi = (\xi_1, \xi_2, \dots, \xi_n)$  with  $\xi_i \geq 0$ . The constraints can be modified as

$$\frac{y_i(\beta^T x_i + \beta_0)}{\|\beta\|} \geq M(1 - \xi_i), \quad i = 1, \dots, n$$

- As before, fixing  $\|\beta\| = 1/M$ , this is equivalent to

$$y_i(\beta^T x_i + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

- The value  $\xi_i$  is the proportional amount by which vector  $x_i$  is on the wrong side of its margin. Misclassification occurs when  $\xi_i > 1$ .

Thierry Denœux SY19 – Support Vector Machines A24 26 / 93

Optimal Separating hyperplane Non-separable case

utc

## Optimization problem

The optimization problem is reformulated as:

$$\min_{\beta, \beta_0, \{\xi_i\}} \frac{1}{2} \|\beta\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \quad (12)$$

subject to

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

$$y_i(\beta^T x_i + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

where  $C$  is a hyperparameter.

## Lagrange function

- The Lagrange function is

$$\begin{aligned} L(\beta, \beta_0, \xi, \alpha, \mu) = & \frac{1}{2} \|\beta\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i \\ & - \sum_{i=1}^n \alpha_i [y_i (\beta^T x_i + \beta_0) - (1 - \xi_i)] - \sum_{i=1}^n \mu_i \xi_i \end{aligned}$$

- Setting the derivatives w.r.t.  $\beta$ ,  $\beta_0$  and  $\xi$  to zero, we get, as before,

$$\beta = \sum_{i=1}^n \alpha_i y_i x_i, \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad (13)$$

and

$$\frac{C}{n} - \alpha_i - \mu_i = 0 \quad (14)$$

## SMO algorithm

- The SMO algorithm is a grouped coordinate ascent procedure.
- Maximizing  $L_D(\alpha)$  one  $\alpha_i$  at a time does not work, because due to the constraint

$$\sum_{i=1}^n \alpha_i y_i = 0$$

variable  $\alpha_i$  is uniquely determined from the other  $\alpha_j$ 's through the equation

$$\alpha_i = -y_i \sum_{j \neq i} \alpha_j y_j$$

- Instead, the SMO algorithm maximizes  $L_D(\alpha)$  w.r.t. to each pair of variables  $(\alpha_i, \alpha_j)$  sequentially.

## Dual formulation

- By substituting (13), we obtain the Lagrangian dual objective function

$$\begin{aligned} L_D(\alpha) = & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \\ & + \underbrace{\sum_{i=1}^n \left( \frac{C}{n} - \alpha_i - \mu_i \right)}_0 \xi_i \quad (15) \end{aligned}$$

which has exactly the same form as in the previous problem.

- We maximize  $L_D$  subject to  $0 \leq \alpha_i \leq \frac{C}{n}$  and  $\sum_{i=1}^n \alpha_i y_i = 0$ .
- The sequential minimal optimization (SMO) algorithm gives an efficient way of solving this problem.

## SMO algorithm (continued)

Repeat until convergence {

- Select some pair  $\alpha_i$  and  $\alpha_j$  to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
- Optimize  $L_D(\alpha)$  with respect to  $\alpha_i$  and  $\alpha_j$ , while holding all the other  $\alpha_k$ 's ( $k \neq i, j$ ) fixed.

}

To test for convergence of this algorithm, we can check whether the KKT conditions are satisfied to within some tolerance (see next slide).

## Interpretation of the solution

- The solution verifies the KKT conditions (13)-(14) and

$$\alpha_i^*[y_i(\beta^{*T}x_i + \beta_0^*) - (1 - \xi_i^*)] = 0, \quad i = 1, \dots, n \quad (16)$$

$$\mu_i^*\xi_i^* = 0, \quad i = 1, \dots, n \quad (17)$$

- As before, the SVs are defined as the points such that  $\alpha_i^* > 0$ .
- From (14) and (17), the SVs such that  $\alpha_i^* < C/n$  verify  $\mu_i^* > 0$  and  $\xi_i^* = 0$ : they lie on the edge of the margin ("in-bound SVs"). The remainder ( $\xi_i^* > 0$ ) have  $\alpha_i^* = C/n$  and usually lie inside the margin ("margin errors").
- The SVs such that  $\xi_i^* > 1$  are misclassified.
- From (16) we can see that any of the in-bound SVs ( $\alpha_i^* > 0, \xi_i^* = 0$ ) can be used to solve for  $\beta_0^*$ , and we typically use an average of all the solutions for numerical stability.

Thierry Denœux

SY19 – Support Vector Machines

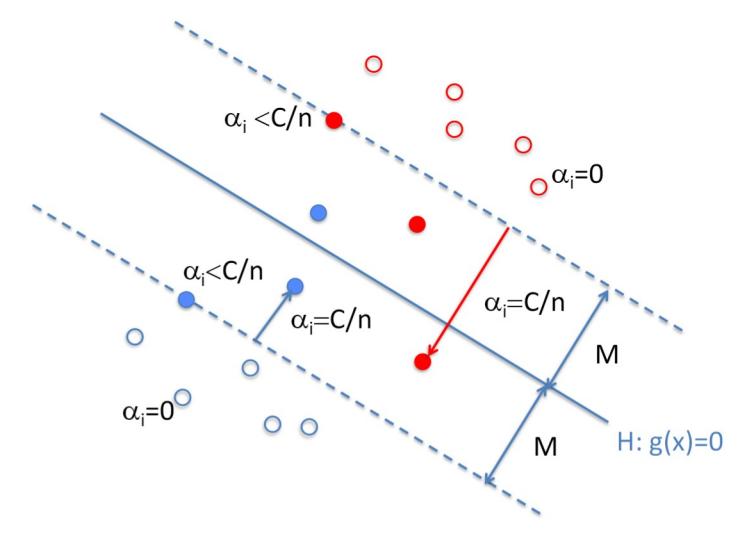


A24 33 / 93

## Tuning $C$

- The tuning parameter of this procedure is the cost parameter  $C$ .
- The optimal value for  $C$  can be estimated by cross-validation.
- From (12), the margin is smaller for larger  $C$ . Hence larger values of  $C$  focus attention more on points near the decision boundary, while smaller values involve data further away.

## Interpretation



Thierry Denœux

SY19 – Support Vector Machines

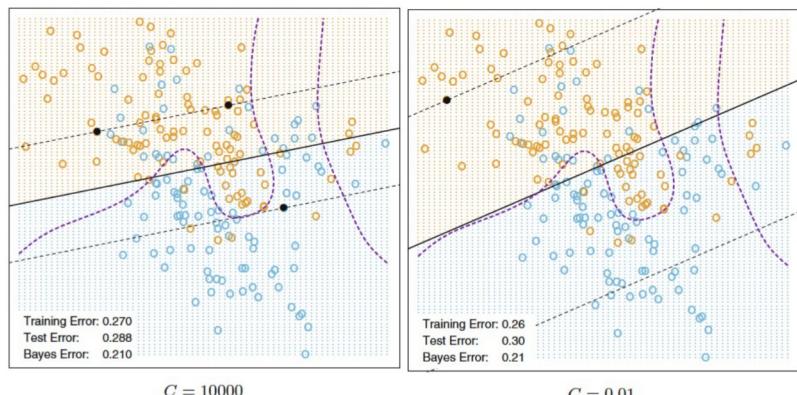
A24 34 / 93

## Bound on the LOO error

- The LOO cross-validation error can be bounded above by the proportion of SVs in the data.
- The reason is that leaving out an observation that is not a SV will not change the solution. Hence these observations, being classified correctly by the original boundary, will be classified correctly in the cross-validation process.
- However this bound tends to be too high, and is not generally useful for choosing  $C$ .



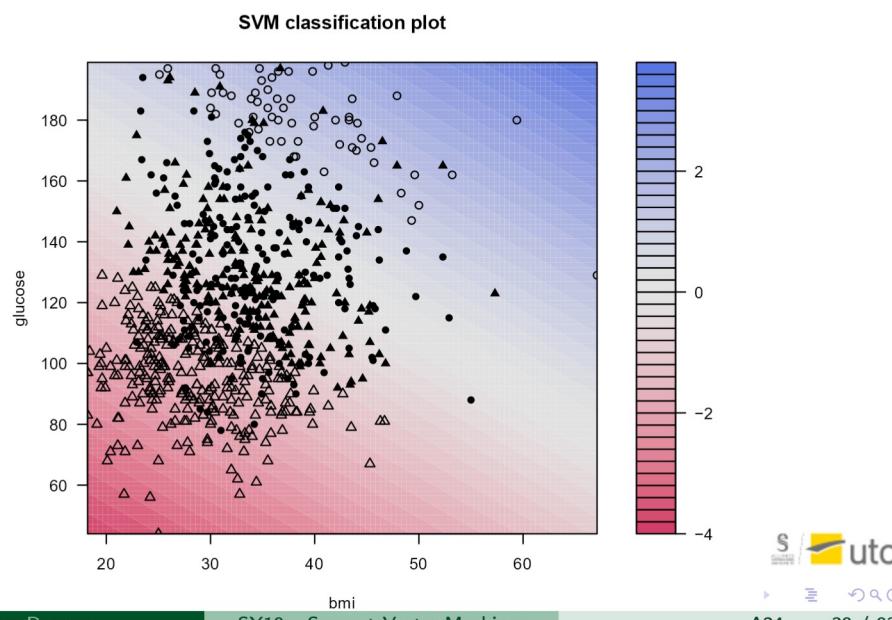
## Example



The SVs ( $\alpha_i^* > 0$ ) are all the points on the wrong side of their margin. The black solid dots are in-bound SVs ( $\alpha_i^* < C/n$ ). In the left (resp., right) panel 62% (resp., 85%) of the observations are SVs.



## Result

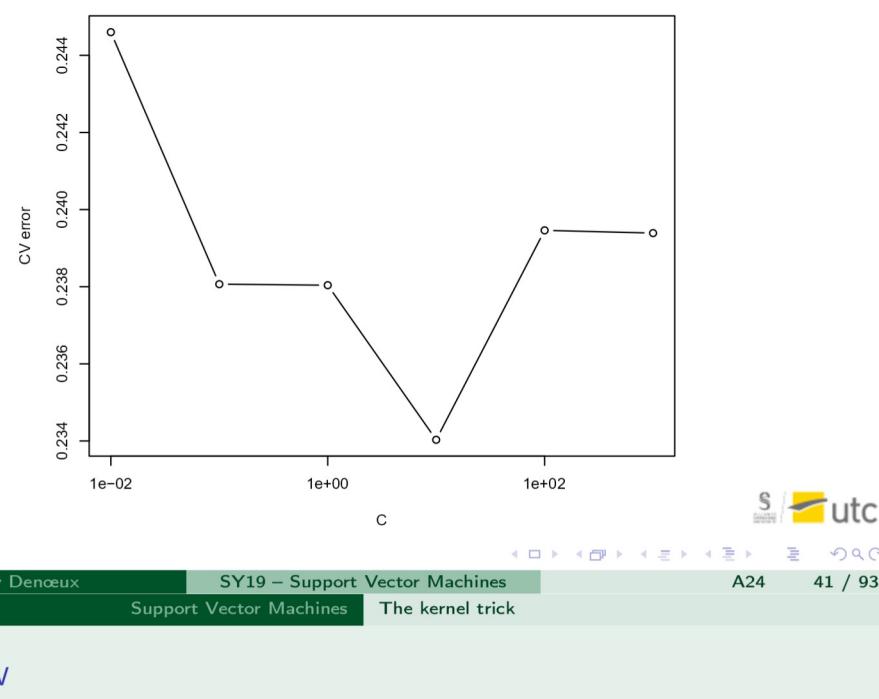


## Application in R

```
library("kernlab")
ii<-which((pima$glucose>0) & (pima$bmi>0))
svmfit<-ksvm(as.factor(class)~ glucose+bmi,data=pima[ii,],
              type="C-svc",kernel="vanilladot",C=10)
plot(svmfit,data=pima[ii,],grid=100)
```



## Cross-validation result



## Overview

### 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

### 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

### 3 Support Vector Regression

- Formalization
- Solution and interpretation



## Extension to non-linear classification

- The support vector classifier described so far finds **linear boundaries** in the predictor space.
- As with other linear methods, we could make the procedure more flexible by **enlarging the predictor space** using basis expansions such as, e.g., polynomials or splines.
- Linear boundaries in the enlarged space generally achieve better training-class separation, and translate to nonlinear boundaries in the original space.



## Extension to non-linear classification (continued)

- Once the basis functions  $\Phi_j(x)$ ,  $j = 1, \dots, J$  are selected, the procedure is the same as before:

  - We fit the SV classifier using predictors

$$\Phi(x_i) = (\Phi_1(x_i), \Phi_2(x_i), \dots, \Phi_J(x_i)), \quad i = 1, \dots, n$$

and produce the (nonlinear) discriminant function

$$g^*(x) = \Phi(x)^T \beta^* + \beta_0.$$

  - The classifier is  $D^*(x) = \text{sign}(g^*(x))$  as before.
  - In SVM, the mapping  $x \rightarrow \Phi(x)$  will be defined implicitly, and  $J$  will be potentially very large (even infinite!).



## Dot-products in the transformed input space

- Assume that the input vector  $x$  is replaced by  $\Phi(x)$  for some transformation  $\Phi : \mathbb{R}^p \rightarrow \mathcal{H}$ .
- The objective function will become

$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \Phi(x_i), \Phi(x_j) \rangle \quad (19)$$

and the optimal discriminant function will be

$$g^*(x) = \sum_{i \in \mathcal{S}} \alpha_i^* y_i \langle \Phi(x_i), \Phi(x) \rangle + \beta_0^* = 0$$

where  $\langle \cdot, \cdot \rangle$  denotes the dot-product in  $\mathcal{H}$ .

- All we need is a method to compute dot-products in  $\mathcal{H}$ .



## The OSH depends only on dot-products

A key feature of the OSH is that it depends only on the dot products between input vectors:

- The solution is found by maximizing

$$L_D(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \quad (18)$$

subject to  $0 \leq \alpha_i \leq C/n$  and  $\sum_{i=1}^n \alpha_i y_i = 0$ .

- The optimal discriminant function is

$$g^*(x) = \sum_{i \in \mathcal{S}} \alpha_i^* y_i x_i^T x + \beta_0^* = 0$$

where  $\beta_0^*$  also depends only on the dot products  $x_i^T x_j$ .



## The “kernel trick”

- If there exists a mapping  $\mathcal{K} : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}_+$  such that

$$\mathcal{K}(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

then the transformation  $\Phi$  will be defined implicitly.

- Such a mapping is called a kernel function, and defining a kernel function is the “kernel trick”.



## Example

- Assume  $p = 2$  and  $\mathcal{K}(x, x') = (x^T x')^2$ .
- We have

$$\begin{aligned}\mathcal{K}(x, x') &= (x_1 x'_1 + x_2 x'_2)^2 \\ &= x_1^2(x'_1)^2 + 2x_1 x_2 x'_1 x'_2 + x_2^2(x'_2)^2 \\ &= \Phi(x)^T \Phi(x')\end{aligned}$$

with

$$\Phi : x \longrightarrow \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \end{pmatrix}$$

- Function  $\Phi$  is defined implicitly by the kernel function  $\mathcal{K}$ .



## Mercer condition

- In practice, a kernel function must verify the following **Mercer condition**:

$$\forall f : \mathbb{R}^p \rightarrow \mathbb{R} \text{ s.t. } \int f(x)^2 dx < \infty, \int \mathcal{K}(x, x') f(x) f(x') dx dx' \geq 0$$

- If this condition is not verified, the Wolfe dual problem may not have a solution.
- In practice, the method may still work most of the time with a kernel function that does not meet this condition.



## Overview

- Optimal Separating hyperplane
  - Formalization
  - Solution in the separable case
  - Non-separable case

### 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

### 3 Support Vector Regression

- Formalization
- Solution and interpretation



## Popular kernel functions

- Three popular choices for  $\mathcal{K}$  in the SVM literature are

$$\begin{aligned}\mathcal{K}(x, x') &= (a + b \cdot x^T x')^d, \quad d > 0 && \text{(polynomial kernel)} \\ \mathcal{K}(x, x') &= \exp[-\sigma \|x - x'\|^2], \quad \sigma > 0 && \text{(RBF or Gaussian kernel)} \\ \mathcal{K}(x, x') &= \tanh(a + b \cdot x^T x') && \text{(MLP kernel).}\end{aligned}$$

- The polynomial and Gaussian verify the Mercer condition, but the MLP kernel does not.
- With the MLP kernel, the discriminant function is

$$g(x) = \sum_{i \in S} \alpha_i^* y_i \tanh(a + b \cdot x_i^T x) + \beta_0^*$$

It is the transfer function of a neural network with  $n_S = \text{card}(S)$  hidden units (see chapter on neural networks).



## Influence of $C$

- The role of parameter  $C$  is clearer in an enlarged predictor space, since perfect separation is often achievable there. (The dimension of  $\mathcal{H}$  may be very large and even infinite.)
- A small value of  $C$  will encourage a small value of  $\|\beta\|$ , which in turn causes  $g(x)$  and hence the boundary to be smoother.
- Both  $C$  and the kernel parameters are usually tuned by cross-validation.

Thierry Denœux SY19 – Support Vector Machines A24 53 / 93

Support Vector Machines Kernel functions

## Application in R

```
x<-matrix(rnorm(200*2),ncol=2)
y<as.factor(c(rep(-1,150),rep(1,50)))
x[1:100,]<- x[1:100,]+2
x[101:150,]<- x[101:150,]-2

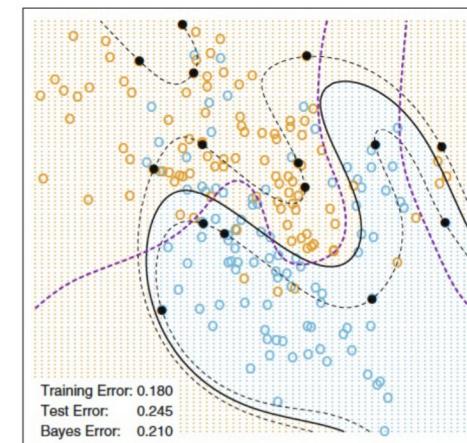
svmfit<-ksvm(x,y,type="C-svc",kernel="rbfdot",
               kpar=list(sigma=1),C=1)
plot(svmfit,data=x,grid=100)
```

Thierry Denœux SY19 – Support Vector Machines A24 55 / 93

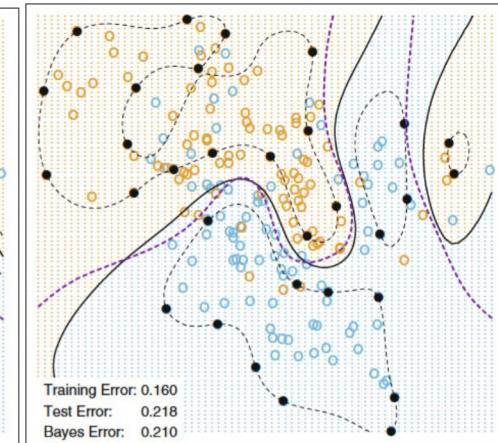
Support Vector Machines Kernel functions

## Example

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space

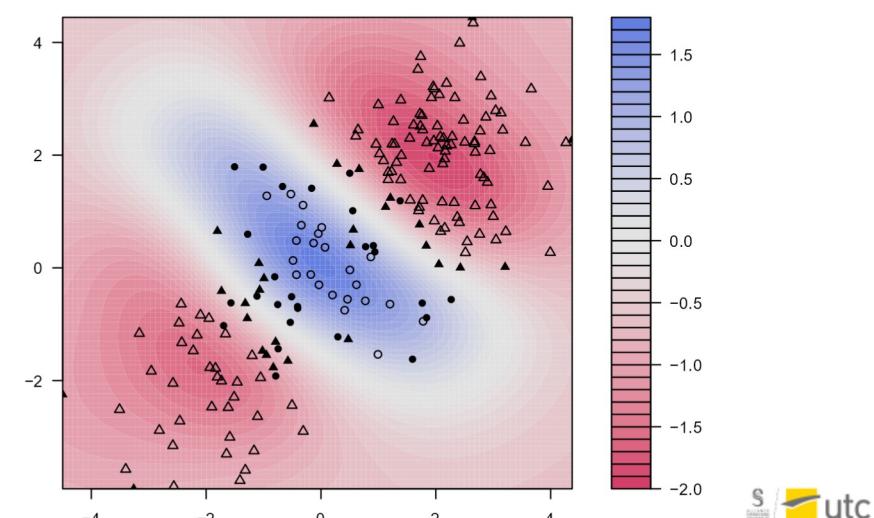


Thierry Denœux SY19 – Support Vector Machines A24 54 / 93

Support Vector Machines Kernel functions

## Result

SVM classification plot



Thierry Denœux SY19 – Support Vector Machines A24 56 / 93

Support Vector Machines Kernel functions

## Estimation of posterior probabilities

- The SVM classifier gives us a decision function, but it does not provide estimates of conditional class probabilities  
 $P(x) = \mathbb{P}(Y = +1 | X = x)$ .
  - One approach to estimate these probabilities is to use [logistic regression](#), with the output  $g(x)$  of the SVM classifier as the predictor.  
 We then have
- $$\hat{P}(x) = \frac{1}{1 + \exp[-(a + b \cdot g(x))]}$$
- To avoid overfitting, it is preferable to estimate the additional parameters  $a$  and  $b$  from a [validation dataset](#) or using [cross-validation](#).

Thierry Denœux SY19 – Support Vector Machines A24 57 / 93

Support Vector Machines Remarks, variants and extensions

S / utc

## SVM as a penalization method

Unconstrained formulation

- Let  $g(x_i) = \langle \beta, \Phi(x_i) \rangle + \beta_0$ . The problem

$$\min_{\beta, \beta_0, \{\xi_i\}} \frac{1}{2} \|\beta\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i$$

subject to  $\xi_i \geq 0$  and  $y_i g(x_i) \geq 1 - \xi_i$ ,  $i = 1, \dots, n$  is equivalent to the unconstrained optimization problem:

$$\min_{\beta, \beta_0} \sum_{i=1}^n \underbrace{[1 - y_i g(x_i)]_+}_{\text{"hinge" loss}} + \underbrace{\frac{\lambda}{2} \|\beta\|^2}_{\text{penalty}}$$

where  $[ \cdot ]_+$  denotes the positive part, with  $\lambda = n/C$ .

- Proof: see next slide.

Thierry Denœux SY19 – Support Vector Machines A24 59 / 93

S / utc

## Overview

- Optimal Separating hyperplane
  - Formalization
  - Solution in the separable case
  - Non-separable case
- Support Vector Machines
  - The kernel trick
  - Kernel functions
  - Remarks, variants and extensions
- Support Vector Regression
  - Formalization
  - Solution and interpretation

Thierry Denœux SY19 – Support Vector Machines A24 58 / 93

Support Vector Machines Remarks, variants and extensions

S / utc

## SVM as a penalization method

Hinge loss

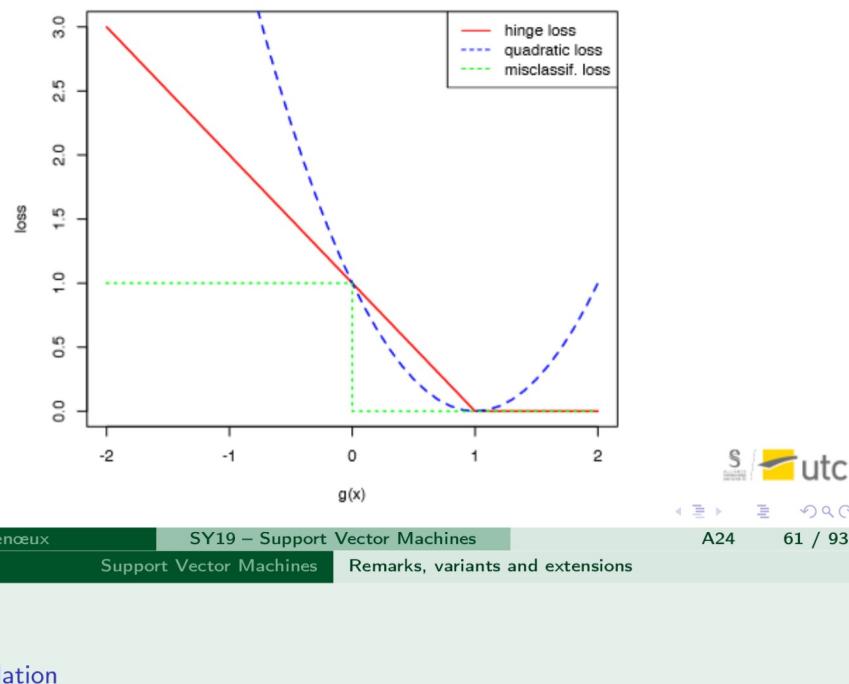
- Proof: from  $\xi_i \geq 0$  and  $\xi_i \geq 1 - y_i g(x_i)$ , we get equivalently  $\xi_i \geq \max(0, 1 - y_i g(x_i)) = [1 - y_i g(x_i)]_+$ . Since we minimize  $\xi_i$ , we set  $\xi_i = [1 - y_i g(x_i)]_+$ .
- The hinge loss can be compared to other loss functions such as:
  - Misclassification:  $I(\text{sign}(g(x)) \neq y)$
  - Squared error loss:  $(y - g(x))^2$
 (see next slide)

Thierry Denœux SY19 – Support Vector Machines A24 60 / 93

S / utc

## SVM as a penalization method

Hinge loss (case  $y = +1$ )



Thierry Denœux

SY19 – Support Vector Machines



## $\nu$ -SVM

Problem formulation

- New formulation:

$$\min_{\beta, \beta_0, \{\xi_i\}, \rho} \frac{1}{2} \|\beta\|^2 - \nu \rho + \frac{1}{n} \sum_{i=1}^n \xi_i$$

subject to

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

$$y_i(\langle \beta, \Phi(x_i) \rangle + \beta_0) \geq \rho - \xi_i, \quad i = 1, \dots, n$$

$$\rho \geq 0$$

- Interpretation: if  $\xi = 0$ , the two classes are separated by a margin  $M = \rho/\|\beta\|$ .

## $\nu$ -SVM

- In the formulation considered until now, the coefficient  $C$  is a rather unintuitive parameter, and we have no a priori way to select it.
- A variant replaces  $C$  by a parameter  $\nu$ , which turns out to control the numbers of margin errors (points with  $\xi_i > 0$ ) and SVs.

Thierry Denœux

SY19 – Support Vector Machines



Thierry Denœux

SY19 – Support Vector Machines



## $\nu$ -SVM

Dual problem

- The derivation of the  $\nu$ -SVM is similar to that of the  $C$ -SVM.

- Dual problem:

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathcal{K}(x_i, x_j)$$

subject to:

$$0 \leq \alpha_i \leq \frac{1}{n}, \quad i = 1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y_i = 0, \quad \sum_{i=1}^n \alpha_i \geq \nu$$

- Remark: there is one more constraint than in the previous formulation.

$\nu$ -SVM

Properties

Proposition (Interpretation of  $\nu$ )

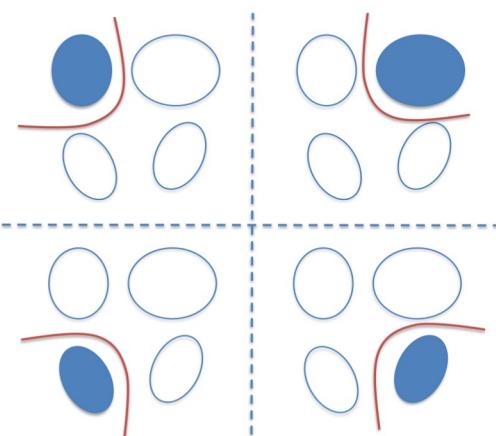
Suppose we run a  $\nu$ -SVM with a kernel  $\mathcal{K}$  on some data with the result that  $\rho > 0$ . Then:

- ①  $\nu$  is an upper bound on the fraction of margin errors (SVs with  $\xi_i > 0$ ).
- ②  $\nu$  is a lower bound on the fraction of SVs.
- ③ Suppose the data  $(x_1, y_1), \dots, (x_n, y_n)$  were generated iid from a continuous distribution. Under very general technical conditions, with probability 1,  $\nu$  equals asymptotically both the fraction of SVs and the fraction of margin errors.

Proposition (Connection between  $\nu$ -SVM and C-SVM)

If  $\nu$ -SVM classification leads to  $\rho > 0$ , then C-SVM classification, with  $C$  set a priori to  $1/\rho$ , leads to the same decision function.

## One-against-all approach



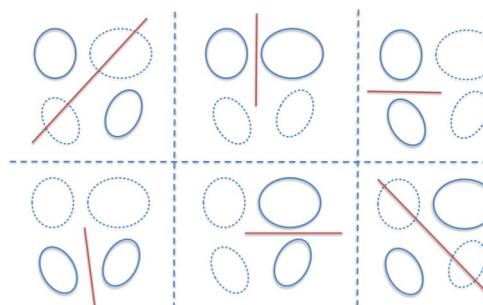
- $c$  binary classifiers are trained to discriminate between each class and the  $c - 1$  others.
- Let  $g_k^*$  be the discriminant function that classifies class  $k$  vs. all other classes. We have

$$D(x) = \arg \max_k g_k^*(x)$$

## From binary to multi-class classification

- So far, the discussion has been restricted to binary classification.
- To apply the SVM technique to multi-class classification, we usually decompose the multi-class problem into several binary problems.
- How?

## One-against-one approach



- $c(c - 1)/2$  binary classifiers are trained to discriminate between each pair of classes.
- Each classifier votes for a class.
- The majority class is selected.

## One-against-all or one-against-one?

- The one-against-one approach implies solving a larger number of binary classification problems, but each one of them is simpler and involves fewer data points.
- The one-against-one approach often outperforms one-against-all and is preferred in most cases, except when the number of classes is very large.
- Function `ksvm` in R package `kernlab` uses the one-against-one approach for multi-class classification.

Thierry Denœux SY19 – Support Vector Machines A24 69 / 93

Support Vector Regression

Overview

### 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

### 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

### 3 Support Vector Regression

- Formalization
- Solution and interpretation

## Example

```
letter <- read.table("letter-recognition.data",header=FALSE)
n<-nrow(letter)
napp=10000
train<-sample(1:n,napp)
letter.test<-letter[-train,]
letter.train<-letter[train,]

fit<- ksvm(as.factor(V1)~.,kernel="rbfdot",C=1,data=letter.train)

# Number of votes for each class
pred<-predict(fit,newdata=letter.test,type = "votes")

# Error rate calculation
pred<-predict(fit,newdata=letter.test,type = "response")
mean(letter.test$V1 != pred)
```

0.0807

Thierry Denœux SY19 – Support Vector Machines A24 70 / 93

Support Vector Regression

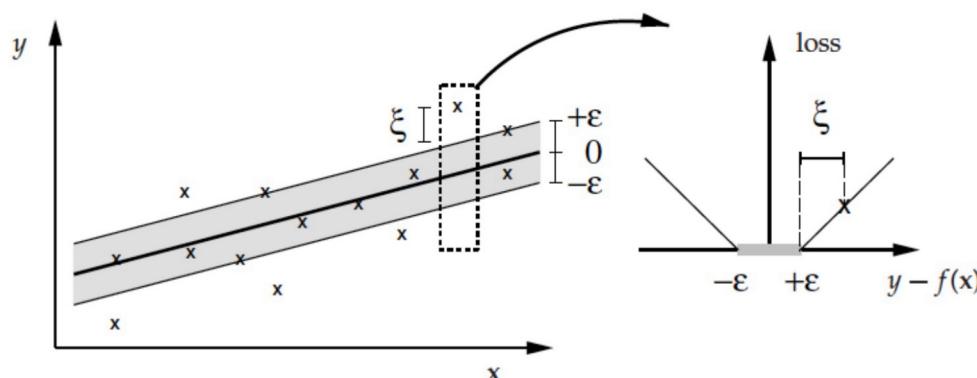
From classification to regression

- SVMs were first developed for classification.
- As described previously, they represent the decision boundary in terms of a typically small subset of all training examples – the support vectors.
- To transpose the SV algorithm to regression, we need to find a way of retaining this feature. This can be achieved using the  $\epsilon$ -insensitive loss function

$$|f(x) - y|_{\epsilon} = \begin{cases} 0 & \text{if } |f(x) - y| \leq \epsilon \\ |f(x) - y| - \epsilon & \text{otherwise} \end{cases}$$

$$= [|f(x) - y| - \epsilon]_+$$

Thierry Denœux SY19 – Support Vector Machines A24 72 / 93

$\epsilon$ -insensitive loss function

The  $\epsilon$ -insensitive loss function does not penalize errors below some  $\epsilon$ , chosen a priori. The  $\epsilon$ -insensitive zone is sometimes referred to as the  $\epsilon$ -tube.

## Overview

## 1 Optimal Separating hyperplane

- Formalization
- Solution in the separable case
- Non-separable case

## 2 Support Vector Machines

- The kernel trick
- Kernel functions
- Remarks, variants and extensions

## 3 Support Vector Regression

- Formalization
- Solution and interpretation

## Basic approach

- The regression algorithm is then developed in close analogy to the case of classification.
- Again, we estimate linear functions, use a  $\|\beta\|^2$  regularizer, and rewrite everything in terms of dot products to generalize to the nonlinear case.

## Problem formulation

- We search for the linear function  $f(x) = \beta^T x + \beta_0$  minimizing the following criterion:

$$\underbrace{\frac{1}{2} \|\beta\|^2}_{\text{regularization}} + \underbrace{\frac{C}{n} \sum_{i=1}^n |f(x_i) - y_i|_\epsilon}_{\text{loss function}}$$

- $C$  is a **hyperparameter**, which balances training error and model complexity.
- To solve this problem, we transform it into an equivalent constrained optimization problem.

## Reformulation as a constrained optimization problem

- We have

$$\begin{aligned}|f(x_i) - y_i|_\epsilon &= [|f(x_i) - y_i| - \epsilon]_+ \\&= \begin{cases} [f(x_i) - y_i - \epsilon]_+ = \xi_i^- & \text{if } f(x_i) \geq y_i \\ [y_i - f(x_i) - \epsilon]_+ = \xi_i^+ & \text{if } f(x_i) < y_i\end{cases}\end{aligned}$$

- Furthermore,  $\xi_i^- = 0$  if  $f(x_i) < y_i$ , and  $\xi_i^+ = 0$  if  $f(x_i) \geq y_i$ . We can thus write

$$|f(x_i) - y_i|_\epsilon = \xi_i^+ + \xi_i^-$$

- The quantities  $\xi_i^+$  and  $\xi_i^-$  are called “slack variables” (they will become slack variables in the constrained optimization formulation of the problem)



## Primal objective function

Using the slack variables, the previous problem can be reformulated as a quadratic optimization problem:

$$\min_{\beta, \beta_0, \xi^-, \xi^+} \frac{1}{2} \|\beta\|^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^- + \xi_i^+)$$

subject to:

$$\xi_i^+ \geq y_i - \beta^T x_i - \beta_0 - \epsilon$$

$$\xi_i^+ \geq 0$$

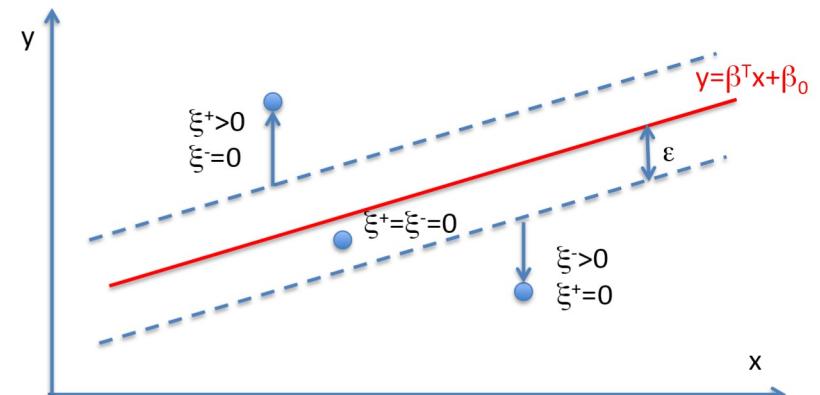
$$\xi_i^- \geq \beta^T x_i + \beta_0 - y_i - \epsilon$$

$$\xi_i^- \geq 0$$

for  $i = 1, \dots, n$ .



## Representation of the slack variables



## Lagrange function

The Lagrange function is

$$\begin{aligned}L(\beta, \beta_0, \xi^-, \xi^+, \alpha^-, \alpha^+, \eta^-, \eta^+) &= \frac{1}{2} \|\beta\|^2 + \frac{C}{n} \sum_{i=1}^n (\xi_i^- + \xi_i^+) \\&\quad - \sum_{i=1}^n (\eta_i^- \xi_i^- + \eta_i^+ \xi_i^+) - \sum_{i=1}^n \alpha_i^- (\epsilon + \xi_i^- + y_i - \beta^T x_i - \beta_0) \\&\quad - \sum_{i=1}^n \alpha_i^+ (\epsilon + \xi_i^+ - y_i + \beta^T x_i + \beta_0),\end{aligned}$$

where  $\alpha_i^-$ ,  $\alpha_i^+$ ,  $\eta_i^-$ ,  $\eta_i^+$  are Lagrange multipliers.



## Derivatives of the Lagrange function

- Setting the derivatives to zero, we obtain:

$$\frac{\partial L}{\partial \beta} = \beta - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) x_i = 0$$

$$\frac{\partial L}{\partial \beta_0} = \sum_{i=1}^n (\alpha_i^- - \alpha_i^+) = 0$$

$$\frac{\partial L}{\partial \xi_i^-} = \frac{C}{n} - \alpha_i^- - \eta_i^- = 0, \quad i = 1, \dots, n$$

$$\frac{\partial L}{\partial \xi_i^+} = \frac{C}{n} - \alpha_i^+ - \eta_i^+ = 0, \quad i = 1, \dots, n$$

- We use these relations to simplify the expression of the Lagrange function (see next slide).



## Dual problem

$$\begin{aligned} L_D(\alpha_i^-, \alpha_i^+) &= -\frac{1}{2} \sum_{i,j} (\alpha_i^+ - \alpha_i^-)(\alpha_j^+ - \alpha_j^-) x_i^T x_j \\ &\quad - \epsilon \sum_{i=1}^n (\alpha_i^+ + \alpha_i^-) + \sum_{i=1}^n y_i (\alpha_i^+ - \alpha_i^-) \end{aligned}$$

to be maximized subject to

$$\sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) = 0$$

$$0 \leq \alpha_i^- \leq \frac{C}{n}, \quad i = 1, \dots, n$$

$$0 \leq \alpha_i^+ \leq \frac{C}{n}, \quad i = 1, \dots, n$$



## Simplification of the Lagrange function

$$\begin{aligned} L &= \frac{1}{2} \left( \sum_i (\alpha_i^+ - \alpha_i^-) x_i \right)^T \left( \sum_j (\alpha_j^+ - \alpha_j^-) x_j \right) \\ &\quad + \sum_{i=1}^n \underbrace{\xi_i^- \left( \frac{C}{n} - \eta_i^- - \alpha_i^- \right)}_0 + \sum_{i=1}^n \underbrace{\xi_i^+ \left( \frac{C}{n} - \eta_i^+ - \alpha_i^+ \right)}_0 \\ &\quad - \epsilon \sum_i (\alpha_i^+ + \alpha_i^-) - \beta_0 \underbrace{\sum_i (\alpha_i^+ - \alpha_i^-)}_0 + \sum_i y_i (\alpha_i^+ - \alpha_i^-) \\ &\quad - \sum_{i=1}^n (\alpha_i^+ - \alpha_i^-) \left( \sum_j (\alpha_j^+ - \alpha_j^-) x_j \right)^T x_i \end{aligned}$$

## Overview

- 1 Optimal Separating hyperplane
  - Formalization
  - Solution in the separable case
  - Non-separable case
- 2 Support Vector Machines
  - The kernel trick
  - Kernel functions
  - Remarks, variants and extensions
- 3 Support Vector Regression
  - Formalization
  - Solution and interpretation



## Support vectors

- As in the case of SVMs, the dual problem can be solved using any quadratic programming solver.
- Let  $\alpha_i^{+*}, \alpha_i^{-*}, i = 1, \dots, n$  be the solution.
- The learning vectors  $x_i$  such that  $\alpha_i^{+*} > 0$  or  $\alpha_i^{-*} > 0$  are called the **support vectors**. They lie outside the tube (or at the border).
- Let  $\mathcal{S}$  be the set of support vectors. We have

$$\beta^* = \sum_{i \in \mathcal{S}} (\alpha_i^{+*} - \alpha_i^{-*}) x_i$$

and

$$f^*(x) = \sum_{i \in \mathcal{S}} (\alpha_i^{+*} - \alpha_i^{-*}) x_i^T x + \beta_0^*$$



## Karush-Kuhn-Tucker conditions

- The solution  $\alpha_i^{+*}, \alpha_i^{-*}, i = 1, \dots, n$  must satisfy the KKT conditions

$$\alpha_i^{+*}(\epsilon + \xi_i^{+*} + y_i - \beta^{*T} x_i - \beta_0^*) = 0 \quad (20a)$$

$$\alpha_i^{-*}(\epsilon + \xi_i^{-*} - y_i + \beta^{*T} x_i + \beta_0^*) = 0 \quad (20b)$$

$$\left(\frac{C}{n} - \alpha_i^{+*}\right)\xi_i^{+*} = 0, \quad \left(\frac{C}{n} - \alpha_i^{-*}\right)\xi_i^{-*} = 0 \quad (20c)$$

- Consequences:

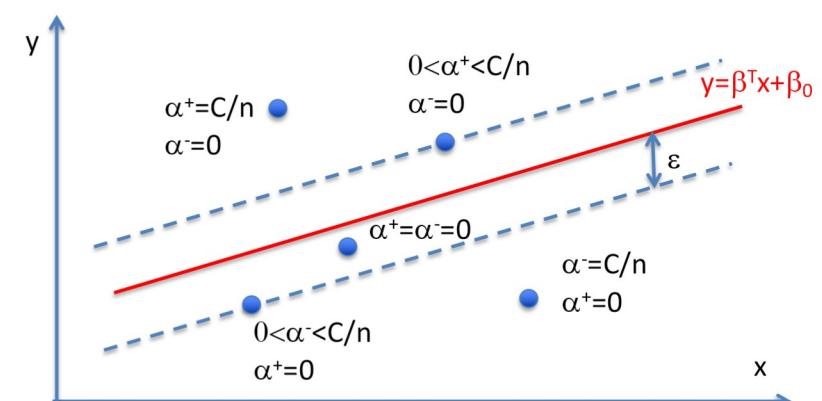
- When  $\alpha_i^{+*} \in (0, C/n)$  or  $\alpha_i^{-*} \in (0, C/n)$ , we have  $\xi_i^{+*} = \xi_i^{-*} = 0$ . The corresponding SVs lie at the border of the tube.
- The other SVs verify  $\alpha_i^{+*} = C/n$  or  $\alpha_i^{-*} = C/n$ ; they lie outside the tube (i.e.,  $\xi_i^{+*} > 0$  or  $\xi_i^{-*} > 0$ ). (See next slide)

## Sparsity of the SV expansion

- We thus have a **sparse expansion** of  $\beta$  in terms of  $x_i$  (we do not need all  $x_i$  to compute  $\beta^*$ ).
- The points inside the tube (i.e., which are not support vectors) do not contribute to the solution: we could remove any one of them, and still obtain the same solution.



## Interpretation of $\alpha_i^+$ and $\alpha_i^-$



## Calculation of $\beta_0$

- $\beta_0^*$  can be calculated from (20a) or (20b) for SVs at the border of the tube as

$$\beta_0^* = \begin{cases} \epsilon + y_i - \beta^{*T} x_i & \text{for } \alpha_i^{-*} \in (0, C/n) \\ y_i - \beta^{*T} x_i - \epsilon & \text{for } \alpha_i^{+*} \in (0, C/n) \end{cases}$$

- Theoretically, it suffices to use any Lagrange multiplier in  $(0, C/n)$ .
- If given the choice between several such multipliers in  $(0, C/n)$ , it is safer to use one that is not too close to 0 or  $C/n$ .

The slide title is 'Nonlinear extension'. The top navigation bar includes the logo 'S / utc', the course name 'SY19 – Support Vector Machines', the section 'Support Vector Regression', and the page number 'A24 89 / 93'. The main content area contains a bulleted list and a mathematical equation.

- As in the classification case, the complete algorithm can be described in terms of **dot products** between the data.
- This makes it possible to formulate a nonlinear extension using **kernels**, replacing dot products  $x_i^T x_j$  in  $\mathcal{X}$  with dot products

$$\langle \Phi(x_i), \Phi(x_j) \rangle = \mathcal{K}(x_i, x_j)$$

in  $\mathcal{H}$ .

- Additional kernel parameters may be determined by cross-validation.

## Parameter tuning

The solution depends on two parameters,  $\epsilon$  and  $C$ . These play different roles:

- Parameter  $\epsilon$  in the loss function specifies the desired accuracy of the approximation. If we scale our response, then we might consider using preset values for  $\epsilon$ .
- The quantity  $C$  is a more traditional regularization parameter. It can be estimated, for example, by cross-validation.

The slide title is 'Application in R'. The top navigation bar includes the logo 'S / utc', the course name 'SY19 – Support Vector Machines', the section 'Support Vector Regression', and the page number 'A24 90 / 93'. The main content area contains R code for fitting a support vector regression model.

```
library('kernlab')
library('MASS')
mcycle.data<-data.frame(mcycle)
mcycle.data$accel<-scale(mcycle.data$accel)
t<- seq(min(mcycle.data$times),max(mcycle.data$times),0.5)
testdat<-data.frame(times=t)

svmfit<-ksvm(accel~,data=mcycle.data,scaled=TRUE,type="eps-svr",
              kernel="rbfdot",C=100,epsilon=0.1,kpar=list(sigma=1))

yhat<-predict(svmfit,newdata=testdat)
plot(mcycle.data$times,mcycle.data$accel)
lines(t,yhat)
```

The slide title is 'Nonlinear extension'. The top navigation bar includes the logo 'S / utc', the course name 'SY19 – Support Vector Machines', the section 'Support Vector Regression', and the page number 'A24 91 / 93'. The bottom navigation bar shows standard presentation controls.

The slide title is 'Application in R'. The top navigation bar includes the logo 'S / utc', the course name 'SY19 – Support Vector Machines', the section 'Support Vector Regression', and the page number 'A24 92 / 93'. The bottom navigation bar shows standard presentation controls.

## Result

