

# Support Vector Machines - II

## Slide Set 14

Prof. Dan A. Simovici

UMB

- 1 Linear Classification
- 2 SVM - The Separable Case
- 3 SVM - The Non-Separable Case
- 4 Margins

# Problem Setting

- **the input space** is  $\mathcal{X} \subseteq \mathbb{R}^n$ ;
- **the output space** is  $\mathcal{Y} = \{-1, 1\}$ ;
- **concept sought**: a function  $f : \mathcal{X} \rightarrow \mathcal{Y} = \{-1, 1\}$ ;
- **sample**: a sequence  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)) \in (\mathcal{X} \times \mathcal{Y})^m$  extracted from a distribution  $\mathcal{D}$ .

## Problem Statement

We are exploring a hypothesis space  $H$  that consists of functions of the form  $h : \mathcal{X} \rightarrow \{-1, 1\}$  such that

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}'\mathbf{x} + b),$$

where

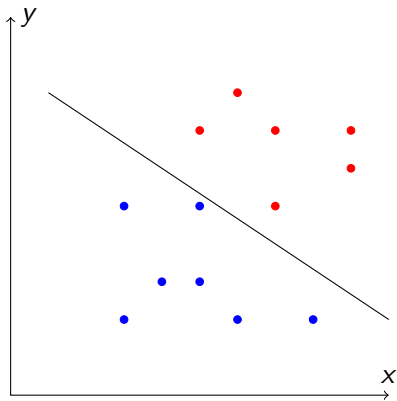
$$\text{sign}(a) = \begin{cases} 1 & \text{if } a \geq 0, \\ -1 & \text{if } a < 0. \end{cases}$$

such that the quantity

$$L(h) = P_{\mathbf{x} \sim \mathcal{D}}(h(\mathbf{x}) \neq f(\mathbf{x}))$$

is small. This is the **generalization error** of  $h$ .

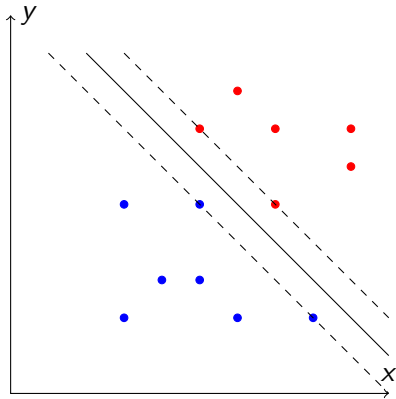
## A Fundamental Assumption: Linear Separability of $S$



If  $S$  is linearly separable there are, in general, infinitely many hyperplanes that can do the separation.

## Solution returned by SVMs

SVMs seek the hyperplane with the **maximum separation margin**.



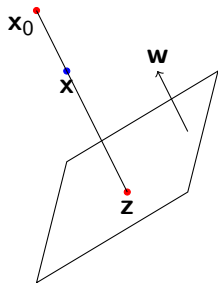
## The distance of a point $\mathbf{x}_0$ to a hyperplane $\mathbf{w}'\mathbf{x} + b = 0$

Equation of the line passing through  $\mathbf{x}_0$  and perpendicular on the hyperplane is

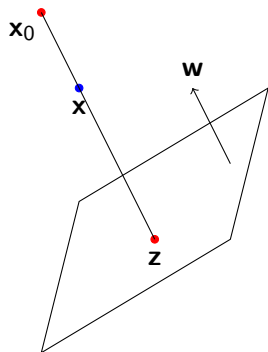
$$\mathbf{x} - \mathbf{x}_0 = t\mathbf{w};$$

Since  $\mathbf{z}$  is a point on this line that belongs to the hyperplane, to find the value of  $t$  that corresponds to  $\mathbf{z}$  we must have  $\mathbf{w}'(\mathbf{x}_0 + t\mathbf{w}) + b = 0$ , that is,

$$t = -\frac{\mathbf{w}'\mathbf{x}_0 + b}{\|\mathbf{w}\|^2}$$



The distance of a point  $\mathbf{x}_0$  to a hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$



Thus,  $\mathbf{z} = \mathbf{x}_0 - \frac{\mathbf{w}'\mathbf{x}_0 + b}{\|\mathbf{w}\|^2} \mathbf{w}$ , hence the distance from  $\mathbf{x}_0$  to the hyperplane is

$$\|\mathbf{x}_0 - \mathbf{z}\| = \frac{|\mathbf{w}'\mathbf{x}_0 + b|}{\|\mathbf{w}\|}.$$



# Primal Optimization Problem

We seek a hyperplane in  $\mathbb{R}^n$  having the equation

$$\mathbf{w}'\mathbf{x} + b = 0,$$

where  $\mathbf{w} \in \mathbb{R}^n$  is a vector normal to the hyperplane and  $b \in \mathbb{R}$  is a scalar. A hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$  that does not pass through a point of  $S$  is in **canonical form** relative to a sample  $S$  if

$$\min_{(\mathbf{x},y) \in S} |\mathbf{w}'\mathbf{x} + b| = 1.$$

Note that we may always assume that the separating hyperplane are in canonical form relative by  $S$  by rescaling the coefficients of the equation that define the hyperplane (the components of  $\mathbf{w}$  and  $b$ ).

## Example

Consider the points:

$$A = \begin{pmatrix} 1 \\ 9 \end{pmatrix}, B = \begin{pmatrix} 4 \\ 2 \end{pmatrix}, C = \begin{pmatrix} 11 \\ 1 \end{pmatrix}, D = \begin{pmatrix} 10 \\ 6 \end{pmatrix}, E = \begin{pmatrix} 10 \\ 3 \end{pmatrix},$$

in  $\mathbb{R}^2$  and the hyperplane  $P$  (in this case, a line)

$$3x_1 + 10x_2 - 60 = 0.$$

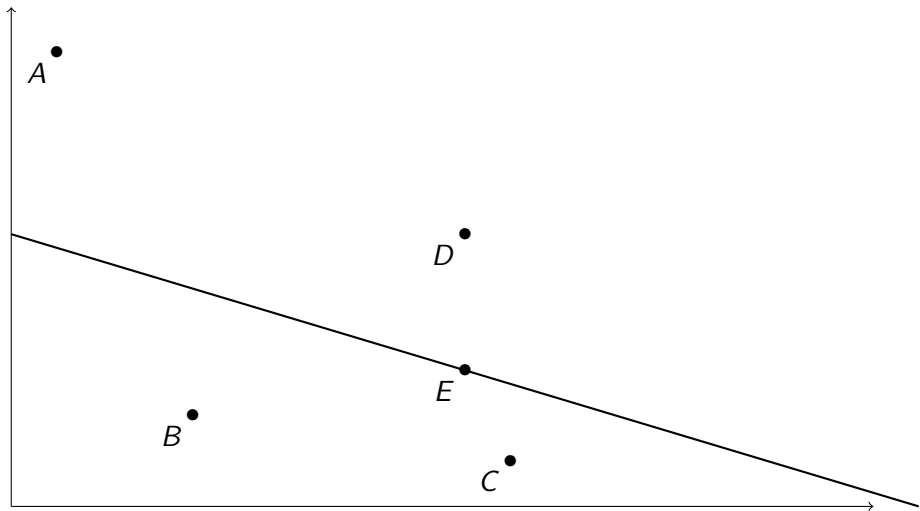
For this hyperplane we have  $\mathbf{w} = \begin{pmatrix} 3 \\ 10 \end{pmatrix}$  and  $b = -60$ . Also,

$$\|\mathbf{w}\| = \sqrt{109}.$$

Note that  $A, B, C, D$  do not belong to the hyperplane (e.g.

$3 \cdot 1 + 10 \cdot 9 = 93 \neq 60$ ), except  $E$  for which we have

$3 \cdot 10 + 10 \cdot 3 - 60 = 0$ . We say that  $E$  is a **support point** for  $P$ .



$$A = \begin{pmatrix} 1 \\ 9 \end{pmatrix}, B = \begin{pmatrix} 4 \\ 2 \end{pmatrix}, C = \begin{pmatrix} 11 \\ 1 \end{pmatrix}, D = \begin{pmatrix} 10 \\ 6 \end{pmatrix}, D = \begin{pmatrix} 10 \\ 3 \end{pmatrix},$$

## Example (cont'd)

### Example

Distances from the points to the hyperplane are:

$$\begin{aligned}d(A, P) &= \frac{|3 + 90 - 60|}{\sqrt{109}} = \frac{33}{\sqrt{109}}, \\d(B, P) &= \frac{|12 + 20 - 60|}{\sqrt{109}} = \frac{28}{\sqrt{109}}, \\d(C, P) &= \frac{33 + 10 - 60}{\sqrt{109}} = \frac{17}{\sqrt{109}}, \\d(D, P) &= \frac{30 + 60 - 60}{\sqrt{109}} = \frac{30}{\sqrt{109}}, \\d(E, P) &= \frac{30 + 30 - 60}{\sqrt{109}} = 0.\end{aligned}$$

## Example

The closest point to  $P$  is  $C$  (except  $E$ ), which means that we can rescale the coefficients of the hyperplane by dividing them by 17. Thus, the equation of the hyperplane in canonical form becomes

$$\frac{3}{17}x_1 + \frac{10}{17}x_2 - \frac{60}{17} = 0.$$

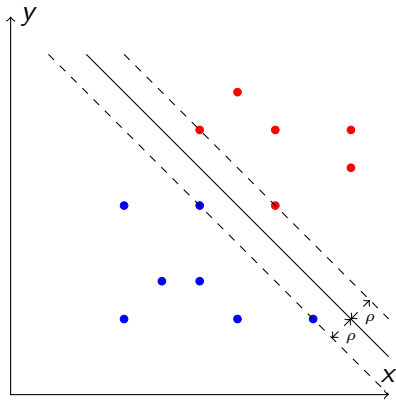
### Example

The minimum distance from one of the external points to the hyperplane is

$$d(C, P) = \frac{\left| \frac{3}{17}x_1 + \frac{10}{17}x_2 - \frac{60}{17} \right|}{\sqrt{109}} = \frac{\left| \frac{33}{17} + \frac{10}{17} - \frac{60}{17} \right|}{\sqrt{109}} = \frac{1}{\sqrt{109}}.$$

If the hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$  is in canonical form relative to the sample  $S$ , then the distance to the hyperplane to the closest points in  $S$  (the margin of the hyperplane) is the same, namely,

$$\rho = \min_{(\mathbf{x}, y) \in S} \frac{|\mathbf{w}'\mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}.$$



# Canonical Separating Hyperplane

For a canonical separating hyperplane we have

$$|\mathbf{w}'\mathbf{x} + b| \geq 1$$

for any point  $(\mathbf{x}, y)$  of the sample and

$$|\mathbf{w}'\mathbf{x} + b| = 1$$

for every support point. The point  $(\mathbf{x}_i, y_i)$  is classified correctly if  $y_i$  has the same sign as  $\mathbf{w}'\mathbf{x}_i + b$ , that is,  $y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1$ .

Maximizing the margin is equivalent to minimizing  $\|\mathbf{w}\|$  or, equivalently, to minimizing  $\frac{1}{2} \|\mathbf{w}\|^2$ . Thus, in the separable case the SVM problem is equivalent to the following convex optimization problem:

- minimize  $\frac{1}{2} \|\mathbf{w}\|^2$ ;
- subjected to  $y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1$  for  $1 \leq i \leq m$ .



Why  $\frac{1}{2} \|\mathbf{w}\|^2$ ?

Note that this objective function,

$$\frac{1}{2} \|\mathbf{w}\|^2 = \frac{1}{2}(w_1^2 + \dots + w_n^2)$$

is differentiable!

We have  $\nabla \left( \frac{1}{2} \|\mathbf{w}\|^2 \right) = \mathbf{w}$  and that

$$H_{\frac{1}{2}\|\mathbf{w}\|^2} = \mathbf{I}_n,$$

which shows that  $\frac{1}{2} \|\mathbf{w}\|^2$  is a convex function of  $\mathbf{w}$ .

# Support Vectors

The Lagrangean of the optimization problem

- minimize  $\frac{1}{2} \|\mathbf{w}\|^2$ ;
- subjected to  $y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1$  for  $1 \leq i \leq m$ .

is

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m a_i (y_i(\mathbf{w}'\mathbf{x}_i + b) - 1).$$

# The Karush-Kuhn-Tucker Optimality Conditions

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{i=1}^m a_i y_i \mathbf{x}_i = 0,$$

$$\nabla_b L = - \sum_{i=1}^m a_i y_i = 0,$$

$$a_i (y_i (\mathbf{w}' \mathbf{x}_i + b) - 1) = 0 \text{ for all } i$$

imply

$$\mathbf{w} = \sum_{i=1}^m a_i y_i \mathbf{x}_i = 0,$$

$$\sum_{i=1}^m a_i y_i = 0,$$

$$a_i = 0 \text{ or } y_i (\mathbf{w}' \mathbf{x}_i + b) = 1 \text{ for } 1 \leq i \leq m.$$

## Consequences of the KKT Conditions

- the weight vector is a linear combination of the training vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ , where  $\mathbf{x}_i$  appears in this combination only if  $a_i \neq 0$  (support vectors);
- since  $a_i = 0$  or  $y_i(\mathbf{w}'\mathbf{x}_i + b) = 1$  for all  $i$ , if  $a_i \neq 0$ , then  $y_i(\mathbf{w}'\mathbf{x}_i + b) = 1$  for the support vectors; thus, all these vectors lie on the marginal hyperplanes  $\mathbf{w}'\mathbf{x} + b = 1$  or  $\mathbf{w}'\mathbf{x} + b = -1$ ;
- if non-support vector are removed the solution remains the same;
- while the solution of the problem  $\mathbf{w}$  remains the same different choices may be possible for the support vectors.

Recall that the optimization problem for SVMs was

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to } y_i(\mathbf{w}'\mathbf{x} + b) \geq 1 \text{ for } 1 \leq i \leq m \end{aligned}$$

Equivalently, the constraints are

$$1 - y_i(\mathbf{w}'\mathbf{x} + b) \leq 0$$

for  $1 \leq i \leq m$ .

The Lagrangean is

$$\begin{aligned} L(\mathbf{w}, b, \mathbf{a}) &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i(1 - y_i(\mathbf{w}'\mathbf{x}_i + b)) \\ &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i - \sum_{i=1}^m a_i y_i \mathbf{w}'\mathbf{x}_i - b \sum_{i=1}^m a_i y_i. \end{aligned}$$

# The Dual Problem

maximize  $L(\mathbf{w}, b, \mathbf{a})$

The KKT conditions are

$$\begin{aligned}(\nabla_{\mathbf{w}} L) &= \mathbf{w} - \sum_{i=1}^m a_i y_i \mathbf{x}_i = \mathbf{0}, \\(\nabla_b L) &= -\sum_{i=1}^m a_i y_i = 0, \\a_i(1 - y_i(\mathbf{w}'\mathbf{x}_i + b)) &= 0,\end{aligned}$$

which are equivalent to

$$\begin{aligned}\mathbf{w} &= \sum_{i=1}^m a_i y_i \mathbf{x}_i, \\ \sum_{i=1}^m a_i y_i &= 0, \\ a_i = 0 &\text{ or } y_i(\mathbf{w}'\mathbf{x}_i + b) = 1,\end{aligned}$$

respectively.

# Implications

- the weight vector  $\mathbf{w}$  is a linear combination of the training vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ ;
- a vector  $\mathbf{x}_i$  appears in  $\mathbf{w}$  if and only if  $a_i \neq 0$  (such vectors are called **support vectors**);
- if  $a_i \neq 0$ , then  $y_i(\mathbf{w}'\mathbf{x}_i + b) = \pm 1$ .

Note that support vectors define the maximum margin hyperplane, or the SVM solution.

# Transforming the Lagrangean

Since

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i - \sum_{i=1}^m a_i y_i \mathbf{w}' \mathbf{x}_i - b \sum_{i=1}^m a_i y_i,$$

$\mathbf{w} = \sum_{j=1}^m a_j y_j \mathbf{x}_j$  (note that we changed the summation index from  $i$  to  $j$ ), and  $\sum_{i=1}^m a_i y_i = 0$ , we have

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i - \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}_j' \mathbf{x}_i.$$



## Further Transformation of the Lagrangean

Note that

$$\begin{aligned}\|\mathbf{w}\|^2 &= \mathbf{w}'\mathbf{w} = \left( \sum_{j=1}^m a_j y_j \mathbf{x}'_j \right) \left( \sum_{i=1}^m a_i y_i \mathbf{x}_i \right), \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_j \mathbf{x}_i.\end{aligned}$$

Therefore,

$$\begin{aligned}L(\mathbf{w}, b, \mathbf{a}) &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m a_i - \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_j \mathbf{x}_i \\ &= \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_j \mathbf{x}_i.\end{aligned}$$

## The Dual Optimization Problem for Separable Sets

$$\begin{aligned} & \text{maximize } \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_i \mathbf{x}_j \\ & \text{subject to } a_i \geq 0 \text{ for } 1 \leq i \leq m \text{ and } \sum_{i=1}^m a_i y_i = 0. \end{aligned}$$

Note that the objective function depends on  $a_1, \dots, a_m$ .

- in this case the strong duality holds; therefore, the primal and the dual problems are equivalent;
- the solution  $\mathbf{a}$  of the dual problem can be used directly to determine the hypothesis returned by the SVM as

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}'\mathbf{x} + b) = \text{sign}\left(\sum_{i=1}^m a_i y_i (\mathbf{x}'_i \mathbf{x}) + b\right);$$

- since support vectors lie on the marginal hyperplanes, for every support vector  $\mathbf{x}_i$  we have  $\mathbf{w}'\mathbf{x}_i + b = y_i$ , so

$$b = y_i - \sum_{j=1}^m a_j y_j (\mathbf{x}'_j \mathbf{x}_i).$$

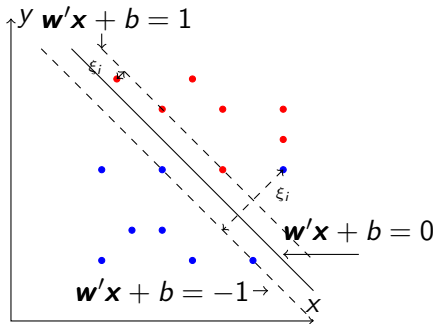
# Slack Variables

If data is not separable the conditions  $y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1$  cannot all hold (for  $1 \leq i \leq m$ ). Instead, we impose a relaxed version, namely

$$y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 - \xi_i,$$

where  $\xi_i$  are new variables known as **slack variables**.

A slack variable  $\xi_i$  measures the distance by which  $\mathbf{x}_i$  violates the desired inequality  $y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1$ .



A vector  $\mathbf{x}_i$  is an outlier if  $\mathbf{x}_i$  is not positioned correctly on the side of the appropriate hyperplane.

- a vector  $\mathbf{x}_i$  with  $0 < y_i(\mathbf{w}'\mathbf{x}_i + b) < 1$  is still an outlier even if it is correctly classified by the hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$  (see the red point);
- if we omit the outliers the data is correctly separated by the hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$  with a **soft margin**  $\rho = \frac{1}{\|\mathbf{w}\|}$ ;
- we wish to limit the amount of slack due to outliers ( $\sum_{i=1}^m \xi_i$ ), but we also seek a hyperplane with a large margin (even though this may lead to more outliers).

## Optimization for Non-Separable Data

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i^p \\ & \text{subject to } y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for } 1 \leq i \leq m. \end{aligned}$$

The parameter  $C$  is determined in the process of cross-validation.

This is a convex optimization problem with affine constraints.

# Support Vectors

As in the separable case:

- constraints are affine and thus, qualified;
- the objective function and the affine constraints are convex and differentiable;
- thus, the KKT conditions apply.



# Variables

- $a_i \geq 0$  for  $1 \leq i \leq m$  are variables associated with  $m$  constraints;
- $b_i \geq 0$  for  $1 \leq i \leq m$  are variables associated with the non-negativity constraints of the slack variables.

The Lagrangean is defined as:

$$L(\mathbf{w}, b, \xi_1, \dots, \xi_m, \mathbf{a}, \mathbf{b}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m a_i [y_i(\mathbf{w}'\mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^n b_i \xi_i.$$

The KKT conditions are:

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \mathbf{w} - \sum_{i=1}^m a_i y_i \mathbf{x}_i = 0 &\Rightarrow & \mathbf{w} = \sum_{i=1}^m a_i y_i \mathbf{x}_i \\ \nabla_b L &= -\sum_{i=1}^m a_i y_i = 0 &\Rightarrow & \sum_{i=1}^m a_i y_i = 0 \\ \nabla_{\xi_i} L &= C - a_i - b_i = 0 &\Rightarrow & a_i + b_i = C \end{aligned}$$

and

$$\begin{aligned} a_i [y_i(\mathbf{w}'\mathbf{x}_i + b) - 1 + \xi_i] &= 0 \text{ for } 1 \leq i \leq m \Rightarrow a_i = 0 \text{ or } \\ y_i(\mathbf{w}'\mathbf{x}_i + b) &= 1 - \xi_i, \\ b_i \xi_i &= 0 \Rightarrow b_i = 0 \text{ or } \xi_i = 0. \end{aligned}$$

# Consequences of the KKT Conditions

- $\mathbf{w}$  is a linear combination of the training vectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ , where  $\mathbf{x}_i$  appears in the combination only if  $a_i \neq 0$ ;
- if  $a_i \neq 0$ , then  $y_i(\mathbf{w}'\mathbf{x}_i + b) = 1 - \xi_i$ ;
- if  $\xi_i = 0$ , then  $y_i(\mathbf{w}'\mathbf{x}_i + b) = 1$  and  $\mathbf{x}_i$  lies on marginal hyperplane as in the separable case; otherwise,  $\mathbf{x}_i$  is an **outlier**;
- if  $\mathbf{x}_i$  is an outlier,  $b_i = 0$  and  $a_i = C$  or  $\mathbf{x}_i$  is located on the marginal hyperplane.
- $\mathbf{w}$  is unique; the support vectors are not.

# The Dual Optimization Problem

The Lagrangean can be rewritten by substituting  $\mathbf{w}$ :

$$\begin{aligned} L &= \frac{1}{2} \left\| \sum_{i=1}^m a_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_i \mathbf{x}_j \\ &\quad - \sum_{i=1}^m a_i y_i b + \sum_{i=1}^m a_i \\ &= \sum_{i=1}^m a_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m a_i a_j y_i y_j \mathbf{x}'_i \mathbf{x}_j, \end{aligned}$$

- the Lagrangean **has exactly the same form as in the separable case**;
- we need  $a_i \geq 0$  and, in addition  $b_i \geq 0$ , which is equivalent to  $a_i \leq C$  (because  $a_i + b_i = C$ );

The dual optimization problem for the non-separable case becomes:

$$\begin{aligned}
 & \text{maximize for } \mathbf{a} \quad \sum_{i=1}^m a_i - \frac{1}{2} a_i a_j y_i y_j \mathbf{x}'_i \mathbf{x}_j \\
 & \text{subject to } 0 \leq a_i \leq C \text{ and } \sum_{i=1}^m a_i y_i = 0 \\
 & \text{for } 1 \leq i \leq m.
 \end{aligned}$$

# Consequences

- the objective function is concave and differentiable;
- the solution can be used to determine the hypothesis

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}'\mathbf{x} + b);$$

- for any support vector  $b_i$  we have  $b = y_i - \sum_{j=1}^m a_j y_j \mathbf{x}'_i \mathbf{x}_j$ .
- the hypothesis returned depends only on the inner products between the vectors and not directly on the vectors themselves.

## Definition

The **geometric margin** relative to a linear classifier  $h(\mathbf{x}) = \mathbf{w}'\mathbf{x} + b$  is its distance to the hyperplane  $\mathbf{w}'\mathbf{x} + b = 0$ :

$$\rho(\mathbf{x}) = \frac{y(\mathbf{w}'\mathbf{x} + b)}{\|\mathbf{w}\|}.$$

The **margin for a linear classifier**  $h$  for a sample  $S = (\mathbf{x}_1, \dots, \mathbf{x}_m)$  is

$$\rho = \min_{1 \leq i \leq m} \frac{y_i(\mathbf{w}'\mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

## Theorem

Let  $S$  be a sample included in a sphere of radius  $r$ ,  $S \subseteq \{\mathbf{x} \mid \|\mathbf{x}\| \leq r\}$ .  
The VC dimension of the set of canonical hyperplanes of the form

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}'\mathbf{x}), \min_{\mathbf{x} \in S} |\mathbf{w}'\mathbf{x}| = 1 \text{ and } \|\mathbf{w}\| \leq \Lambda,$$

verifies  $d \leq r^2 \Lambda^2$ .



# Proof

Suppose that  $\{\mathbf{x}_1, \dots, \mathbf{x}_d\}$  is a set that can be fully shattered. Then, for all  $\mathbf{y} = (y_1, \dots, y_d) \in \{-1, 1\}^d$  there exists  $\mathbf{w}$  such that  $1 \leq y_i(\mathbf{w}'\mathbf{x}_i)$  for  $1 \leq i \leq d$ .

Summing up these inequalities yields:

$$d \leq \mathbf{w}' \sum_{i=1}^d y_i \mathbf{x}_i \leq \|\mathbf{w}\| \cdot \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\| \leq \Lambda \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|.$$

## Proof (cont'd)

Since  $y_1, \dots, y_d$  are independent, if  $i \neq j$ ,  $E(y_i y_j) = E(y_i)E(y_j) = 0$ ; also,  $E(y_i y_i) = 1$ .

Since  $d \leq \Lambda \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|$  holds for all  $\mathbf{y} \in \{-1, 1\}^d$ , it holds over expectations and we have

$$\begin{aligned} d &\leq \Lambda E_{\mathbf{y}} \left( \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\| \right) \leq \Lambda \left( E_{\mathbf{y}} \left( \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|^2 \right) \right)^{1/2} \\ &= \Lambda \left( \sum_{i=1}^m \sum_{j=1}^m E_{\mathbf{y}}(y_i y_j) (\mathbf{x}'_i \mathbf{x}_j) \right)^{1/2} \\ &= \Lambda \left( \sum_{i=1}^d \mathbf{x}'_i \mathbf{x}_i \right)^{1/2} \leq \Lambda (dr^2)^{1/2} = \Lambda r \sqrt{d}. \end{aligned}$$

Thus,

$$d \leq \Lambda^2 r^2$$

- recall that when the data is linearly separable the margin  $\rho$  is given by:

$$\rho = \min_{(\mathbf{x}, y) \in S} \frac{|\mathbf{w}'\mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|};$$

- if we restrict the sample  $S$  such that the resulting  $\mathbf{w}$  is such that  $\|\mathbf{w}\| = \frac{1}{\rho} = \Lambda$ , it follows that

$$d \leq \frac{r^2}{\rho^2}.$$