

### **MACHINE LEARNING**

# Support Vector Machine For Classification

Part 1 – Linear SVM

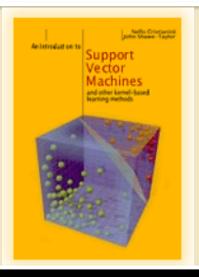


### Support Vector Machine (SVM)

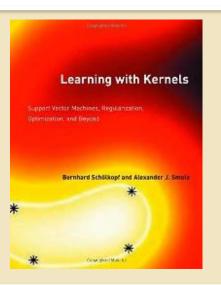
#### **Brief history**:

- SVM is traced back to the work by Vapnik and Chervonekis on statistical learning theory (Vapnik1979) and the notion of VC dimension.
- The current form of SVM was presented in (Boser, Guyon and Vapnik 1992) and Cortes and Vapnik (1995).

#### Textbooks:



A good survey of the theory behind SVM is given in *Support Vector Machines and other Kernel Based Learning methods* by Nello Cristianini and John-Shawe Taylor.



An easy introduction to SVM is given in *Learning with Kernels* by Bernhard Scholkopf and Alexander Smola.



### Support Vector Machine (SVM)

#### SVM was applied to numerous classification problems:

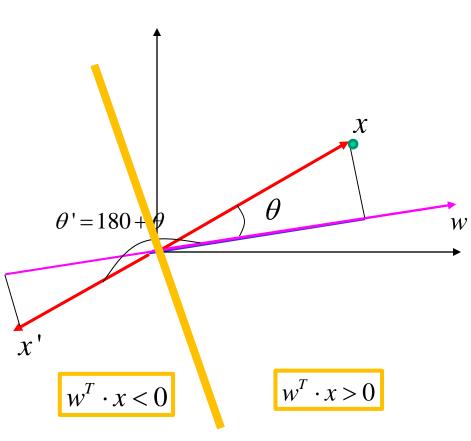
- Computer vision (face detection, object recognition, feature categorization)
- Bioinformatics (categorization of gene expression, of microarray data)
- WWW (categorization of websites)
- Production (control of quality, detection of defaults)
- Robotics (categorization of sensor readings)
- Finance (bankruptcy prediction)

#### The success of SVM is mainly due to:

- SVM depends on convex optimization.
- Its ease of use (lots of software available, good documentation).
- Excellent performance on variety of datasets.
- Good solvers making optimization (learning phase) very quick.
- Very fast at retrieval time does not hinder practical applications.



### Recap - Constructing a projection



#### Datapoint *x*

Projection vector w

The norm of the projection of x onto w is:

$$w^T \cdot x = ||w|| ||x|| \cos(\theta)$$

$$\cos(\theta) > 0 \implies w^T \cdot x > 0$$

$$w^T \cdot x' = ||w|| ||x'|| \cos(\theta')$$

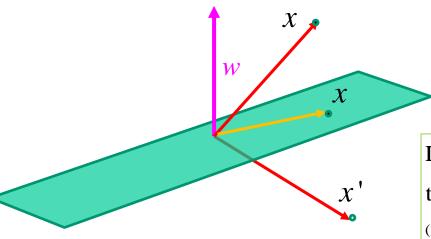
$$\cos(\theta') < 0 \implies w^T \cdot x' < 0$$



### Recap - Constructing a projection



Normal to plane w



Looking at the sign of  $(w^T x)$  allows to separate points on either side of the plane.

(we ignored the intercept and assumed the plane passed by the origin)

 $w^T \cdot x > 0 \implies$  point lies on the left handside of plane

 $w^T \cdot x' < 0 \implies$  point lies on the right handside of plane

$$w^T \cdot x = 0$$
?  $\Rightarrow$  point on the plane



fiers

(w, b)

x

f

Class label y={-1;1}

- denotes -1
- ° denotes +1

Separating Hyperplane

Separating hyperplane is defined by:

 $y = f(x; w, b) = \operatorname{sgn}(w^T x + b)$ 

w: the normal to the plane

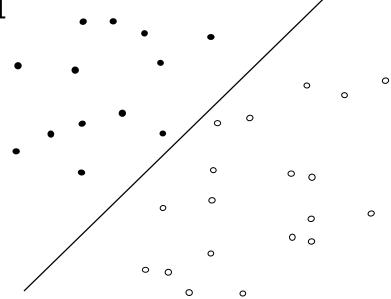
b: the intercept

6



Class label  $y=\{-1;1\}$ 

- denotes -1
- ° denotes +1

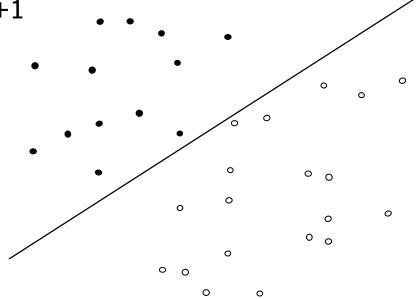


How would you classify this data?



Class label  $y=\{-1;1\}$ 

- denotes -1
- ° denotes +1

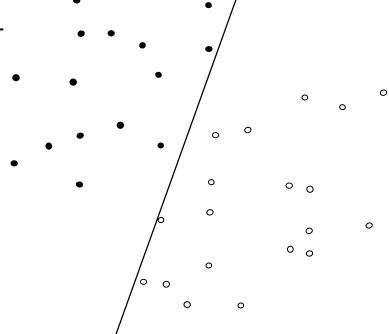


How would you classify this data?



Class label  $y=\{-1;1\}$ 

- denotes -1
- ° denotes +1



How would you classify this data?



Class label y={-1;1}

- denotes -1
- ° denotes +1

Any of these would be fine..

0 0

0

..but which is best?



# Classifier Margin

Class label y={-1;1}

- denotes -1
- ° denotes +1

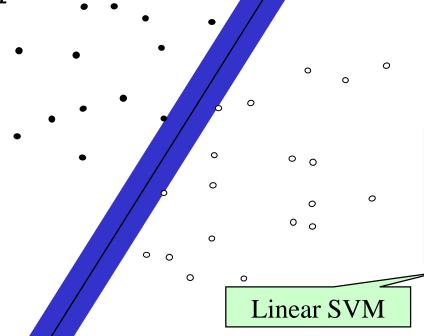
Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.



# Classifier Margin

Class label  $y=\{-1;1\}$ 

- denotes -1
- ° denotes +1



The maximum margin linear classifier is the linear classifier with the maximum margin.



# Classifier Margin

Class label  $y=\{-1;1\}$ 

- denotes -1
- ° denotes +1

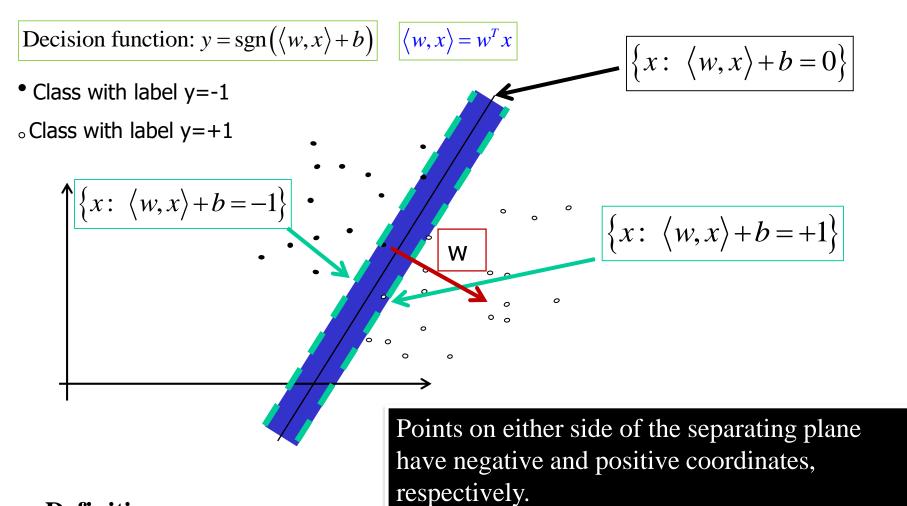
**Support Vectors** 

are those datapoints that are closest to the boundary. They define the margin. Need to determine a measure of the width of the margin, so as to maximize for this measure.

0 0



#### Computing the Distance to the Separating Hyperplane

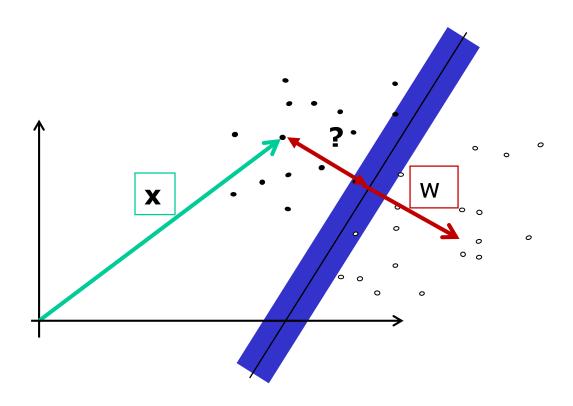


#### **Definition**:

The margins on either side of the hyperplane satisfy  $\langle w, x \rangle + b = \pm 1$ .



#### Computing the Distance to the Separating Hyperplane

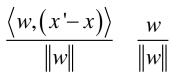


What is the distance from a point  $\mathbf{x}$  to the separating plane  $\langle \mathbf{w}, \mathbf{x} \rangle + \mathbf{b} = 0$ ?



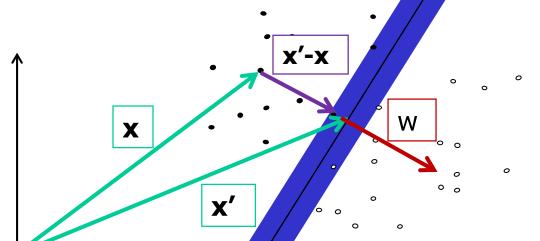
#### Computing the Distance to the Separating Hyperplane

Projection of (x'-x) onto w is then:



unitary vector

$$\langle w, x' - x \rangle = \langle w, x' \rangle - \langle w, x \rangle$$



We know that:

$$x'$$
 s.t.  $\langle w, x' \rangle + b = 0$   
 $\Rightarrow \langle w, x' \rangle = -b$ 

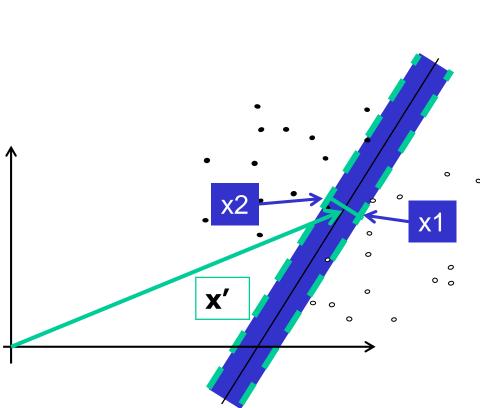
Projection of x - x' onto w is then:

$$\frac{-b - \langle w, x \rangle}{\|w\|} \frac{w}{\|w\|}$$
unitary vector

Distance to plane = 
$$\frac{\left| \left\langle w, x \right\rangle + b \right|}{\left\| w \right\|}$$



#### Computing the margin



Distance of each points on either side of the margin:

$$\|\mathbf{x}^{1} - \mathbf{x}'\| = \frac{\left|\left\langle w, x^{1} \right\rangle + b\right|}{\|w\|} = \frac{1}{\|w\|}$$

$$\|\mathbf{x}^{2} - \mathbf{x}'\| = \frac{\left|\left\langle w, x^{2} \right\rangle + b\right|}{\|w\|} = \frac{1}{\|w\|}$$



### Our objective function

Separating condition is measured by  $\frac{2}{\|w\|}$ .

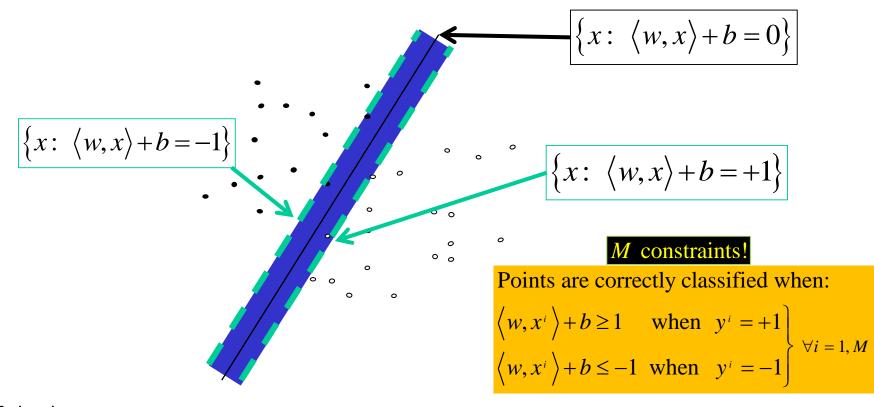
To maximize this condition is equivalent to minimizing  $\frac{\|w\|}{2}$ .

Better even is to minimize the convex form  $\frac{\|w\|^2}{2}$ .

How do we make sure the points sit on the correct side of the separating plane?



#### Determining the constraints



We have 3 situations:

$$0 < y^{i} \left( \left\langle w, x^{i} \right\rangle + b \right) < 1 \iff -1 < \left( \left\langle w, x^{i} \right\rangle + b \right) < 1 \quad \text{on the correct side but inside the margin}$$

$$1 < y^{i} \left( \left\langle w, x^{i} \right\rangle + b \right) \iff \left( \left\langle w, x^{i} \right\rangle + b \right) < -1 \quad \text{for } y^{i} = -1 \text{ or } \left( \left\langle w, x^{i} \right\rangle + b \right) > 1 \quad \text{for } y^{i} = +1 \text{ on the correct side } and \text{ outside the margin}$$

$$y^{i} \left( \left\langle w, x^{i} \right\rangle + b \right) < 0 \quad \Leftrightarrow \left( \left\langle w, x^{i} \right\rangle + b \right) > 0 \text{ for } y^{i} = -1 \text{ or } \left( \left\langle w, x^{i} \right\rangle + b \right) < 0 \quad \text{for } y^{i} = +1 \quad \text{on the wrong side!}$$



#### The complete problem

Finding the Optimal Separating Hyperplane turns out to be an optimization problem of the following form:

$$\min_{w,b} \frac{1}{2} \|w\|^{2}$$

$$\langle w, x^{i} \rangle + b \ge 1 \quad \text{when} \quad y^{i} = +1$$

$$\langle w, x^{i} \rangle + b \le -1 \quad \text{when} \quad y^{i} = -1$$

$$\Rightarrow y^{i} (\langle w, x^{i} \rangle + b) \ge 1, \text{ i=1,2,...,M.}$$

- *N*+1 parameters (N: dimension of data)
- *M* constraints (M: nm of datapoints)
- It is called the *primal problem*.



#### Solving the constrained optimization

$$\min_{w,b} \frac{1}{2} \|w\|^2$$
under constraints:  $y^i \left( \left\langle w, x^i \right\rangle + b \right) \ge 1$ , i=1,2,...,M.

This can be solved using the Lagrange method for inequality constraints:

$$L(w,b,\alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^{M} \alpha_i \left( y^i \left( \left\langle w, x^i \right\rangle + b \right) - 1 \right)$$
with  $\alpha_i \ge 0$ 

We have M Lagrange multipliers  $a_i$ , i = 1, ..., M (M, # of data points), one for each of the inequality constraints.

(Minimization of convex function under linear constraints through Lagrange gives the optimal solution, see complement of information on moodle).



#### Solving the constrained optimization

The solution of this problem is found when maximizing over  $\alpha$  and minimizing over w and b:

$$\max_{\alpha \ge 0} \left( \min_{w,b} L(w,b,\alpha) \right)$$

where

$$L(w,b,\alpha) \equiv \frac{1}{2} \|w\|^2 - \sum_{i=1}^{M} \alpha_i \left( y^i \left( \left\langle w, x^i \right\rangle + b \right) - 1 \right)$$



#### Solving the constrained optimization

The solution of this problem is found when maximizing over  $\alpha$  and minimizing over w and b:

$$\max_{\alpha \geq 0} \left( \min_{w,b} L(w,b,\alpha) \right)$$

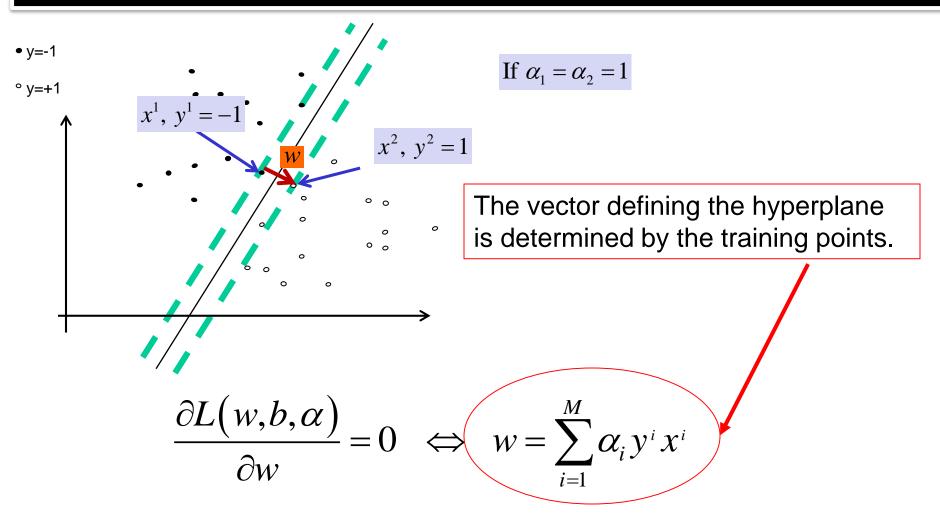
Requesting that the gradient of L vanishes with w.

$$\frac{\partial L(w,b,\alpha)}{\partial w} = 0 \iff w = \sum_{i=1}^{M} \alpha_i y^i x^i$$

(Take the partial derivatives on each coordinate of w)



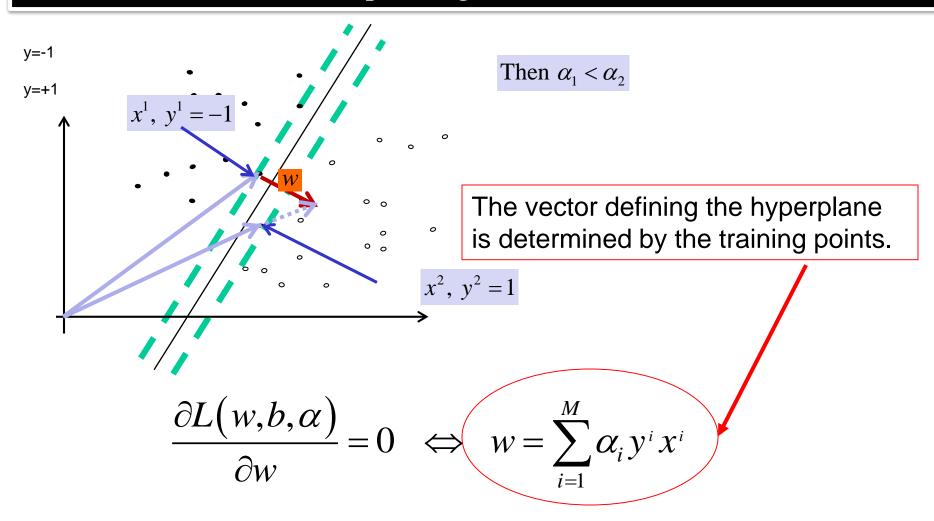
#### Interpreting the solution



Note that while w is unique (minimization of convex function), the alpha-s are not unique.



#### Interpreting the solution



Note that while w is unique (minimization of convex function), the alpha-s are not unique.



### Interpreting the solution

Requesting that the gradient of L vanishes with b.

$$\frac{\partial L(w,b,\alpha)}{\partial b} = 0 \iff \sum_{i=1}^{M} \alpha_i y^i = 0$$

Requires at minimum one datapoint in each class.



#### The dual optimization

Taking the definition of w and plugging it to the Lagrangian

$$L(\alpha) = \sum_{i=1}^{M} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} \alpha_{i} \alpha_{j} y_{i} y_{j} x_{i}^{T} x_{j} - \sum_{i=1}^{M} \alpha_{i} y_{i} b$$

$$= \sum_{i=1}^{M} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} \alpha_{i} \alpha_{j} y_{i} y_{j} x_{i}^{T} x_{j}$$

#### Dual optimization problem

$$\max_{\alpha} W(\alpha_{i}) = \sum_{i=1}^{M} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} \alpha_{i} \alpha_{j} y_{i} y_{j} x_{i}^{T} x_{j}$$

subject to: 
$$\alpha_i \ge 0$$
 and  $\sum_{i=1}^{M} \alpha_i y_i = 0$ 

This is usually solved through the: Sequential Minimal Optimization algorithm (SMO)



#### The Karush-Kuhn-Tucker (KKT) Conditions

The KKT conditions ensure that our primal and dual optimization problems have the same optimal solutions

Complete optimization problem:

$$\frac{\partial L(w,b,\alpha)}{\partial w} = 0 \iff w = \sum_{i=1}^{M} \alpha_i y^i x^i \qquad \text{(Dual feasibility)}$$

$$\frac{\partial L(w,b,\alpha)}{\partial b} = 0 \iff \sum_{i=1}^{M} \alpha_i y^i = 0$$
 (Dual feasibility)

$$\frac{\partial L(w,b,\alpha)}{\partial \alpha} \le 0 \iff y^i \left( \left\langle w, x^i \right\rangle + b \right) \ge 1 \quad \text{(Primal feasibility)}$$

Karush-Kuhn-Tucker conditions:

$$\alpha_{i}\left(y^{i}\left(\left\langle w,x^{i}\right\rangle +b\right)-1\right)=0\quad\forall\ i=1,..M$$
 (Complementarity conditions)  $\alpha_{i}\geq0,\qquad\forall\ i=1,..M$ 



#### Interpreting the conditions

All the pairs of data points  $(x^i, y^i)$  for which  $\alpha_i > 0$  are the support vectors.

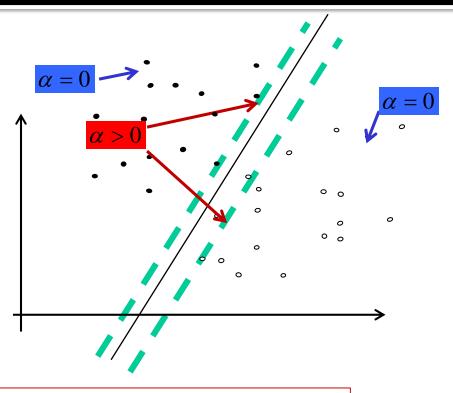
All the pairs of data points  $(x^i, y^i)$  for which  $\alpha_i = 0$  are "ignored".

#### Karush-Kuhn-Tucker conditions:

$$\alpha_{i} \left( y^{i} \left( \left\langle w, x^{i} \right\rangle + b \right) - 1 \right) = 0 \quad \forall i = 1, ..M$$
 (Complementarity conditions) 
$$\alpha_{i} \geq 0, \qquad \forall i = 1, ..M$$



#### Interpreting the conditions



Data points  $(x^i, y^i)$  for which  $\alpha_i > 0$  are the support vectors.

They participate in defining the hyperplane:  $w = \sum_{i=1}^{M} \alpha_i y^i x^i$ 

Data points  $(x^i, y^i)$  for which  $\alpha_i = 0$  are "ignored".



#### The decision function in SVM

The decision function is then expressed in terms of the support vectors:

$$f(x) = sgn(\langle w, x \rangle + b) \qquad \frac{\partial L(w, b, \alpha)}{\partial w} = 0 \iff w = \sum_{i=1}^{M} \alpha_i y^i x^i$$

$$= sgn\left(\sum_{i=1}^{M} \alpha_{i} y^{i} \left\langle x, x^{i} \right\rangle + b\right)$$

Use 
$$\left(y^{i}\left(\left\langle\sum_{i=1}^{M}\alpha_{i}y^{i}x^{i},x^{i}\right\rangle+b\right)-1\right)=0$$

to compute b.



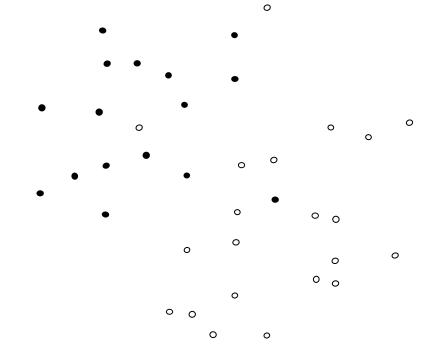
### Non-Separable Data Sets

What should we do?

Idea:

Introduce some slack on the constraints

- denotes +1
- denotes -1



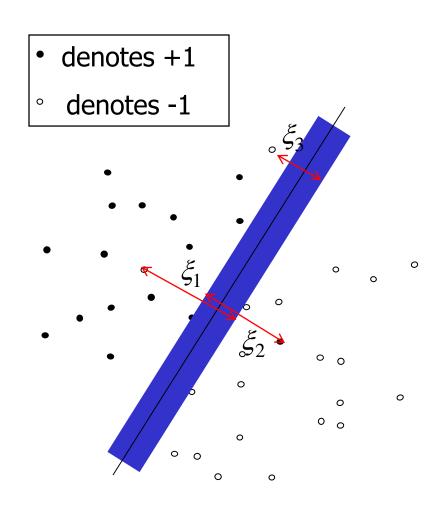


### Support Vector Machine for non-separable datasets

Introduce slack variables:  $\xi_i \ge 0$ 

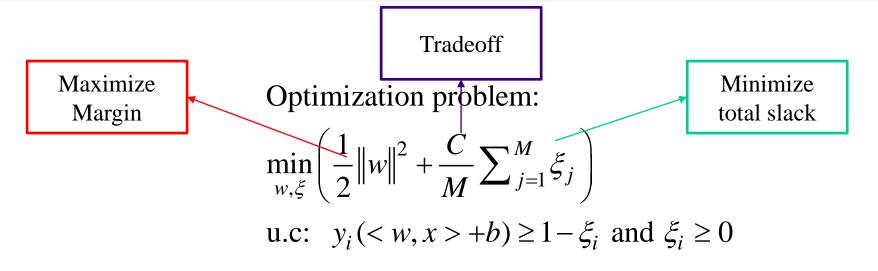
Relax the constraints:

$$y_i(< w, x > +b) \ge 1 - \xi_i \text{ and } \xi_i \ge 0$$





### Support Vector Machine for non-separable datasets



Three cases for  $\xi$ :

 $\xi_m = 0 \rightarrow$  correct classification and outside the margin  $0 < \xi_m < 1 \rightarrow$  correct classification inside margin  $\xi_m \ge 1 \rightarrow$  missclassification



### Support Vector Machine for non-separable datasets

The Dual Form is given by:

$$\max_{\alpha} L_{D}(\alpha) \equiv \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \langle x^{i}, x^{j} \rangle$$

Subject to these constraints:

$$0 \le \alpha_j \le \frac{C}{M} \quad \forall j = 1, ..., M \qquad \sum_{j=1}^{M} \alpha_j y_j = 0$$

$$\sum_{j=1}^{M} a_j y_j = 0$$

The hyperplane has the same solution:

$$w = \sum_{j=1}^{M} \alpha_j y_j x^j$$

Datapoints with  $\alpha_i > 0$ will be the support vectors.