

Sparse kernel machinesMaximum margin classifiers

Consider  $K=2$  classification using a linear model:  $y(\vec{x}) = \vec{w}^T \vec{g}(\vec{x}) + b$

Training dataset:  $\{\vec{x}_1, \dots, \vec{x}_N\}$   
 $\{t_1, \dots, t_N\}$   $t_n \in \{-1, 1\}$

Thus new data points are classified by the sign of  $y(\vec{x})$ .

Consider a linearly separable dataset:  
 $\exists$  one or more set  $\{\vec{w}, b\}$  s.t.

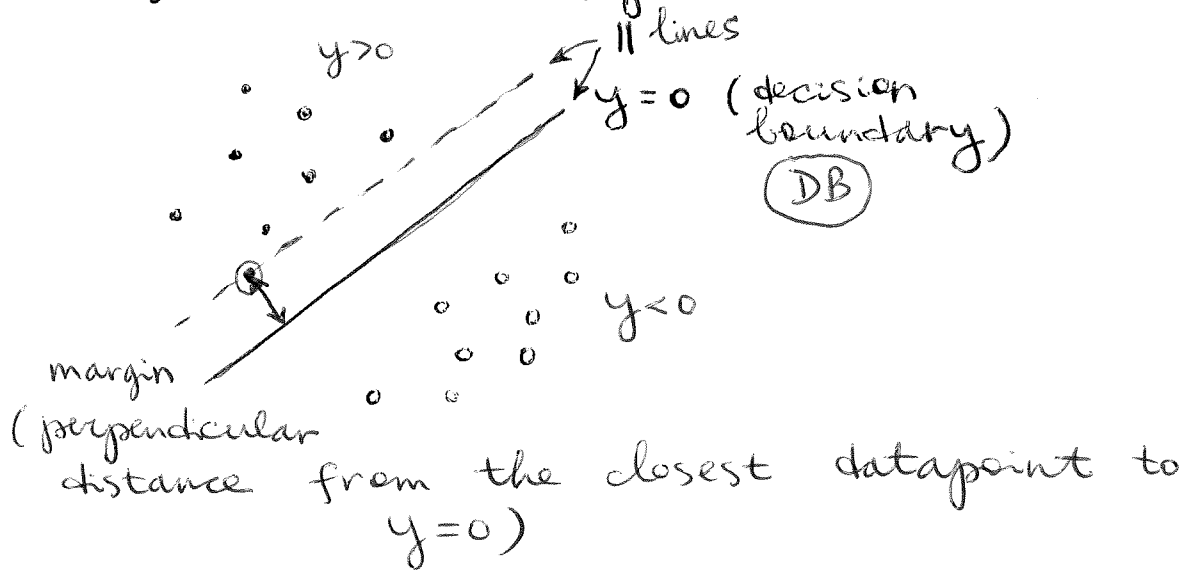
$$\begin{cases} y(\vec{x}_n) > 0 & \text{for all points with } t_n = 1 \\ y(\vec{x}_n) < 0 & \text{for all points with } t_n = -1 \end{cases}$$

$\Leftrightarrow$  always have  $t_n y(\vec{x}_n) > 0$   
 $n=1, \dots, N$

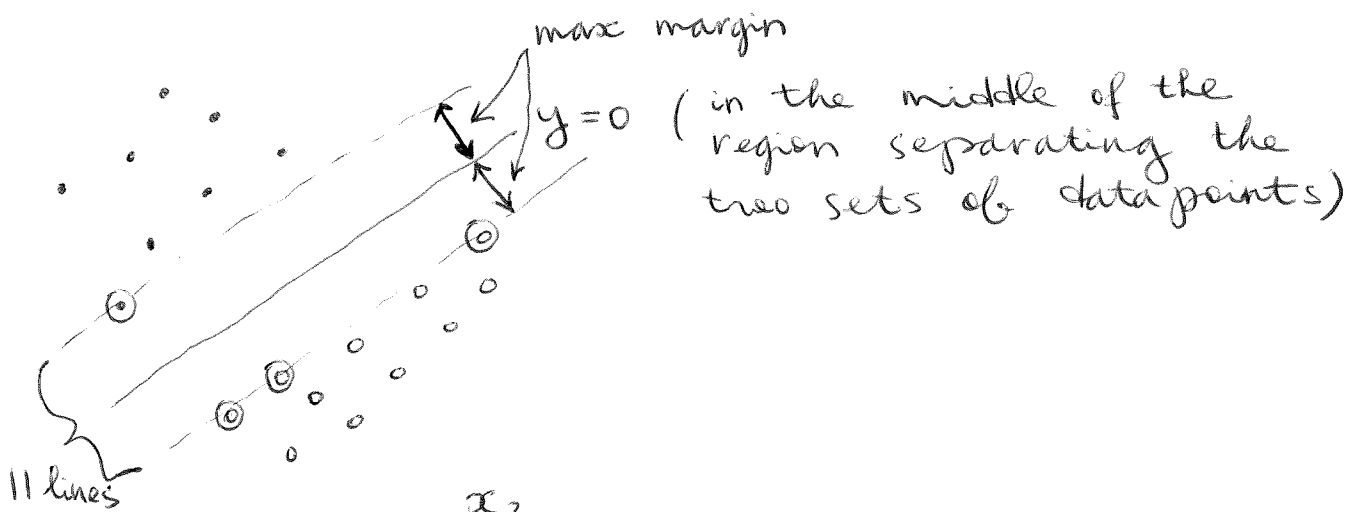
There may be many  $\{\vec{w}, b\}$ , as was clear with e.g. the perceptron algorithm.

Which one is the best?

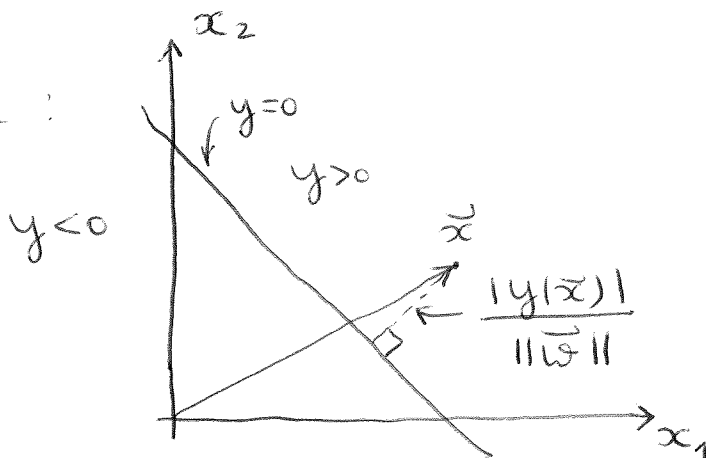
Consider an example:



Idea: choose  $\{\vec{w}, b\}$  which maximize the margin.



Recall:



Thus the distance from  $\vec{x}_n$  to  $y=0$  is given by:

$$\frac{t_n y(\vec{x}_n)}{\|\vec{w}\|} = \frac{t_n (\vec{w}^T \vec{y}(\vec{x}_n) + b)}{\|\vec{w}\|}$$

$$t_n y(\vec{x}_n) > 0, \forall n$$

The max margin solution is given by

$$\arg \max_{\vec{w}, b} \left\{ \frac{1}{\|\vec{w}\|} \min_n [t_n (\vec{w}^T \vec{y}(\vec{x}_n) + b)] \right\} \quad (*)$$

finds  $\vec{x}_n$  closest to DB

Note that (\*) is invariant under  $\begin{cases} \vec{w} \rightarrow \kappa \vec{w} \\ b \rightarrow \kappa b \end{cases} \quad \forall \kappa > 0$

Use this freedom to choose

$$t_n (\vec{w}^T \vec{y}(\vec{x}_n) + b) = 1 \text{ for the closest } \vec{x}_n.$$

$$\text{Then } \underbrace{t_n (\vec{w}^T \vec{y}(\vec{x}_n) + b)}_{\text{canonical representation}} \geq 1, \forall n \quad (**)$$

$\vec{x}_n$ 's for which the LHS = 1 are called active (there is at least one), while all other  $\vec{x}_n$ 's are called inactive. Then (\*) becomes

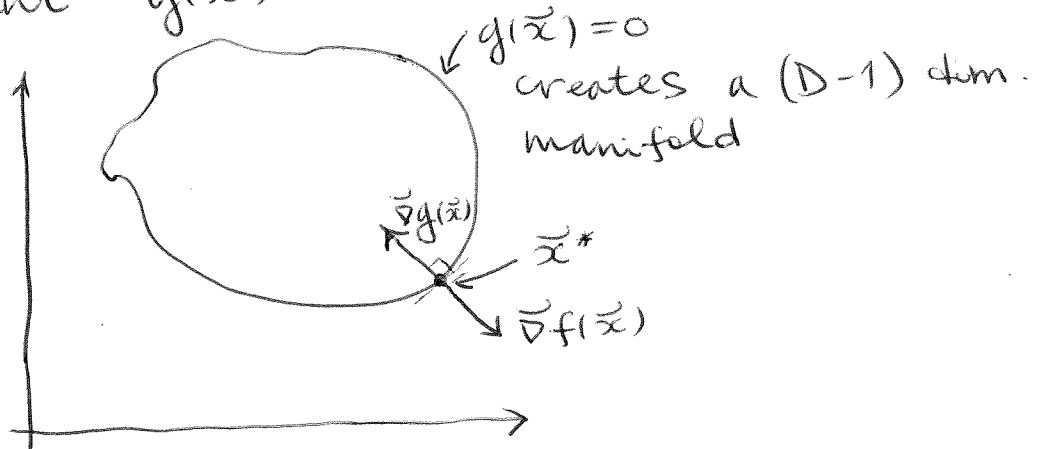
$$\arg \max_{\vec{w}, b} \left\{ \frac{1}{\|\vec{w}\|} \right\} \Rightarrow \arg \min_{\vec{w}, b} \left\{ \frac{1}{2} \|\vec{w}\|^2 \right\}$$

It's a quadratic programming problem: minimize a quadratic function subject to linear constraints (\*\*)

This is best approached using Lagrange multipliers.

Lagrange multipliers

Consider  $\vec{x} = \langle x_1, \dots, x_D \rangle$  subject to a constraint  $g(\vec{x}) = 0$



We want to maximize  $f(\vec{x})$  subject to  $g(\vec{x}) = 0$ .

Consider  $g(\vec{x} + \vec{\epsilon}) \approx g(\vec{x}) + \vec{\epsilon}^T \cdot \nabla g(\vec{x})$   
both points lie on the surface :  $g(\vec{x} + \vec{\epsilon}) = g(\vec{x}) = 0$

As  $\|\vec{\epsilon}\| \rightarrow 0$ ,  $\vec{\epsilon}^T \cdot \nabla g(\vec{x}) = 0 \Rightarrow \vec{\epsilon}$  is parallel to the surface in this limit, so  $\nabla g(\vec{x}) \perp$  surface.

We seek  $\vec{x}^*$  on the surface s.t.  $f(\vec{x})$  is maximized. At  $\vec{x}^*$ ,  $\nabla f(\vec{x}) \perp$  surface, otherwise we could move along the gradient & maximize  $f(\vec{x})$  further.

Thus  $\vec{\nabla}g \uparrow \downarrow \vec{\nabla}f$  or  $\vec{\nabla}g \uparrow \uparrow \vec{\nabla}f$ . In both cases,

$$(***) \quad \vec{\nabla}f + \lambda \vec{\nabla}g = 0 \quad \text{for some } \lambda \neq 0$$

↑  
Lagrange multiplier

Introduce  $L(\vec{x}, \lambda) = f(\vec{x}) + \lambda g(\vec{x})$ , then

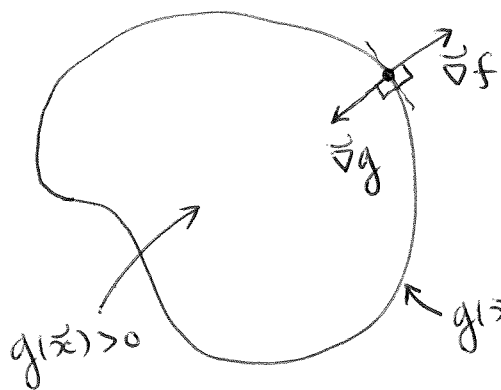
$\vec{\nabla}L = 0$  gives (\*\*\*) and

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow g(\vec{x}) = 0, \quad \text{original constraint}$$

Now, consider  $g(\vec{x}) \geq 0$  inequality constraint

If  $g(\vec{x}) > 0 \Rightarrow \vec{\nabla}f(\vec{x}) = 0$ , or  
 "constraint inactive"  $\vec{\nabla}L = 0$  with  $\lambda = 0$

If  $g(\vec{x}) = 0 \Rightarrow \vec{\nabla}f + \lambda \vec{\nabla}g = 0$  as before, for some  $\lambda \neq 0$ .  
 "constraint active"



Note that  $\vec{\nabla}f$  should point away from the  $g(\vec{x}) > 0$  region since otherwise we could improve in it by stepping inside the region

So either  $f(\vec{x})$  is maximized inside the region or on the boundary. In the latter case,  $\vec{\nabla}f$  points away from the region.  
 Note that  $\vec{\nabla}g$  points inward since  $g(\vec{x}) > 0$  inside the boundary

This means that

$$\vec{\nabla} f = -\lambda \vec{\nabla} g \quad \text{for some } \underline{\underline{\lambda > 0}}$$

In either case,  $\lambda g(\vec{x}) = 0$ .

Thus, maximizing  $f(\vec{x})$  subject to  $g(\vec{x}) \geq 0$  is the same as finding

$$\begin{cases} \frac{\partial L}{\partial \vec{x}} = 0, \\ \nabla L = 0 \end{cases}, \text{ subject to}$$

$$\begin{cases} g(\vec{x}) \geq 0, \\ \lambda \geq 0, \\ \lambda g(\vec{x}) = 0 \end{cases} \quad \text{KKT conditions}$$

If we want to minimize  $f(\vec{x})$  subject to  $g(\vec{x}) \geq 0$ , we minimize

$$L(\vec{x}, \lambda) = f(\vec{x}) - \lambda g(\vec{x}) : \begin{cases} \vec{\nabla} L = 0 \\ \frac{\partial L}{\partial \lambda} = 0 \end{cases}$$

subject to  $\lambda \geq 0$

To solve the constrained quadratic programming problem, we introduce

$$J(\vec{w}, b, \vec{a}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\vec{w}^T \cdot \vec{\Psi}(\vec{x}_n) + b) - 1\},$$

where  $\vec{a} = (a_1, \dots, a_N)$  &  $a_n \geq 0, \forall n$  are Lagrange multipliers.

$$\begin{cases} \frac{\partial J}{\partial \vec{w}} = 0 \Rightarrow \vec{w} = \sum_n a_n t_n \vec{\Psi}(\vec{x}_n), \\ \frac{\partial J}{\partial b} = 0 \Rightarrow \sum_n a_n t_n = 0. \end{cases}$$

Then  $J \rightarrow \tilde{J}(\vec{a}) = \sum_n a_n + \frac{1}{2} \sum_{n,m} a_n t_n a_m t_m \times$   
 $\times \vec{\Psi}^T(\vec{x}_n) \cdot \vec{\Psi}(\vec{x}_m) - \sum_{n,m} a_n t_n a_m t_m \vec{\Psi}^T(\vec{x}_m) \cdot \vec{\Psi}(\vec{x}_n) -$   
 $- b \sum_n a_n t_n = \sum_n a_n - \frac{1}{2} \sum_{n,m} a_n a_m t_n t_m \underbrace{\vec{\Psi}^T(\vec{x}_n) \cdot \vec{\Psi}(\vec{x}_m)}_{k(\vec{x}_n, \vec{x}_m)},$   
 subject to  $\begin{cases} a_n \geq 0, \forall n \\ \sum_n a_n t_n = 0 \end{cases}$

We need to maximize  $\tilde{J}(\vec{a})$ .

To classify new datapoints, need to evaluate the sign of

$$y(\vec{x}) = \vec{w}^T \cdot \vec{\Psi}(\vec{x}) + b = \sum_n a_n t_n \vec{\Psi}^T(\vec{x}_n) \cdot \vec{\Psi}(\vec{x}) + b =$$

$$= \sum_n a_n t_n k(\vec{x}_n, \vec{x}) + b \quad (*)$$

The KKT conditions are:

$$\begin{cases} a_n \geq 0, \\ t_n y(\vec{x}_n) - 1 \geq 0, \\ a_n \{t_n y(\vec{x}_n) - 1\} = 0 \end{cases} \Rightarrow \begin{array}{l} \text{for each } n, \\ a_n = 0 \text{ OR } t_n y(\vec{x}_n) = 1. \end{array}$$

If  $a_n = 0$ , (\*) is unaffected, so only the points for which  $t_n y(\vec{x}_n) = 1$  matter. These are called support vectors, and lie on the max margin hyperplanes.

If we have  $\vec{a}$  by solving the quadratic programming problem, we can find  $b$ :

$$t_n \left( \sum_{m \in S} a_m t_m k(\vec{x}_n, \vec{x}_m) + b \right) = 1, \quad \forall n \in S$$

$\uparrow$   
 set of support vectors of size  $N_s$

Average over all  $n \in S$  for numerical stability:

$$b = \frac{1}{N_s} \sum_{n \in S} \left[ t_n - \sum_{m \in S} a_m t_m k(\vec{x}_n, \vec{x}_m) \right]$$

Finally,  $\vec{w} = \sum_{n \in S} a_n t_n \vec{g}(\vec{x}_n)$ .

The max margin classifier can be described by an error function (to be minimized)

$$\sum_{n=1}^N E_{\infty}(y(\vec{x}_n) t_n - 1) + \lambda \|\vec{w}\|^2, \quad \text{where}$$

$$E_{\infty}(z) = \begin{cases} 0, & z \geq 0 \\ \infty, & z < 0 \end{cases} \quad \lambda > 0, \text{ otherwise its precise value is not important.}$$



Note that max margin classifiers <sup>↳ SVMs</sup> can draw non-linear boundaries in the  $\bar{x}$ -space, due to use of non-linear kernels (cf. Fig. 7.9 in Bishop).

### Overlapping class distributions

So far, we've assumed linear separability in some (potentially unknown explicitly) feature space  $\bar{y}(\bar{x})$ . We need to modify the SVM to allow for training point misclassification in order to relax this assumption.

Introduce slack variables:

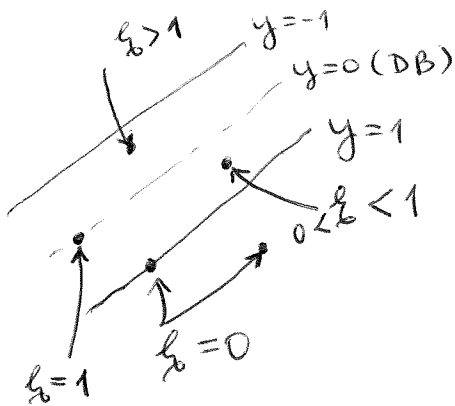
$$\xi_n \geq 0 \quad n=1, \dots, N$$

s.t.  $\xi_n = 0$  for all points on or inside the correct margin boundary &

$$\xi_n = |t_n - y(\bar{x}_n)| \quad \text{for } \underbrace{\text{misclassified}}_{\text{all points on the wrong side of the margin boundary}}$$

linear penalty

Points on DB have  $\xi_n = 1$ , and all points with  $\xi_n > 1$  are misclassified:



Now,  $\underbrace{t_n y(\bar{x}_n) \geq 1}_{\text{hard margin constraint}}, \forall n$

is replaced with

$$t_n y(\bar{x}_n) \geq 1 - \xi_n, \forall n$$

$$\xi_n \geq 0$$

soft margin constraint