

# SVM<sub>A</sub>

## Separable case.



First, a remark:  $\vec{w}\vec{x} + b = 0$  defines a plane

$d_{\pm} = \frac{\vec{w}\vec{x} + b}{\|\vec{w}\|_2}$  is the distance between

the point  $\vec{x}$  and the plane  $(\vec{w}, b)$ . (signed distance)

• On one side of the plane,  $d_{\pm} \geq 0$ , on the other side, we have  $d_{\pm} \leq 0$ .

• If  $\vec{x}_n$  is well classified,  $\text{sign}(d_{\pm}) = \pm 1 = t_n$ .

• So, for well-classified points,  $d_{\pm}(\vec{x}_n) \cdot t_n = \frac{t_n(\vec{w}\vec{x}_n + b)}{\|\vec{w}\|_2} = d$  is the (positive) distance between  $\vec{x}_n$  and the plane.

$d = \frac{t_n(\vec{w}\vec{x}_n + b)}{\|\vec{w}\|} \geq 0$  is a distance (positive)

• Rescaling: the plane  $(\vec{w}, b)$  is the same as the plane encoded in  $(c\vec{w}, cb)$ , ( $\forall c \in \mathbb{R}^{+*}$  ( $c > 0$ )):

we can check this:  $d' = \frac{t_n(\vec{w} \cdot c\vec{x}_n + cb)}{c\|\vec{w}\|} = d$ .

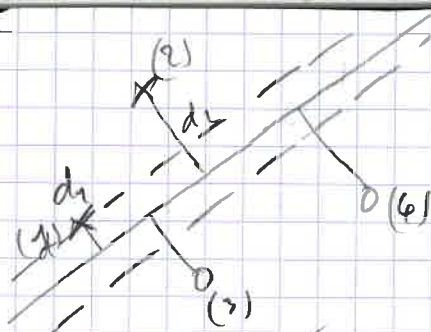
• So, we can always choose  $c$  (rescale  $\|\vec{w}\|$ ) such that we have  $t_n(\vec{w}\vec{x}_n + b) \geq 1, \forall n$  (the choice of  $\vec{w}$  is arbitrary)

• There are many solutions to correctly classify a linearly separable set of data.

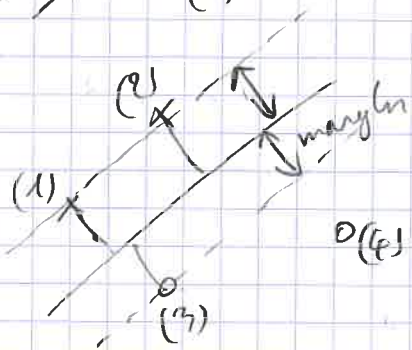
• We call "margin" the minimal distance between the points and the plane:

$$\text{margin} = \min_n (d(\vec{x}_n)) = \min_n \left( \frac{t_n(\vec{w}\vec{x}_n + b)}{\|\vec{w}\|_2} \right)$$





$d_1$  is minimal, so the margin is  $\text{margin} = d_1 = d(\vec{x}_1)$  in this example.



here,  $d_3$  is the smallest dist (almost equal to  $d_1$ ), so  $\text{margin} = d_3 = d(\vec{x}_3)$ .

This (b) seems better than the first one.

• The core idea of SVM (also re-named "Separators with Vast Margin") is to find the solution with the largest margin, which is expected to work better than those with small margin.

• So, we define the SVM sol<sup>o</sup> (among all the sol<sup>o</sup> to the prob of linear separation) as:

$$w_{SVM}^* = \underset{w, b}{\operatorname{argmax}} \left( \underbrace{\min_n \left( \underbrace{\frac{tn(w\vec{x}_n + b)}{\|w\|_2}}_{\text{distance}} \right)}_{\text{margin}} \right)$$

the largest margin (or its corresp.  $\vec{w}$ , rather)

From there, it's just maths (optimize<sup>o</sup>).



$$b, w_{svm}^* = \operatorname{argmax}_{w,b} \left( \min_n \left( \frac{k_n(\vec{w} \cdot \vec{x}_n + b)}{\|\vec{w}\|_2} \right) \right)$$

$$= \operatorname{argmax}_{w,b} \left( \frac{1}{\|\vec{w}\|_2} \min_n \left( \frac{k_n(\vec{w} \cdot \vec{x}_n + b)}{\|\vec{w}\|_2} \right) \right)$$

the distances  $d_n$  can be set to  $\geq 1$  by rescaling  $\vec{w}, b$  by a constant  $c > 0$ .

$$= \operatorname{argmax}_{w,b \text{ such that } k_n(\vec{w} \cdot \vec{x}_n + b) \geq 1, \forall n} \left( \frac{1}{\|\vec{w}\|_2} \right)$$

because  $\frac{1}{x^2}$  is  $\downarrow$  for all  $x > 0$ ,  $\operatorname{argmax}$  becomes  $\operatorname{argmin}$

$$= \operatorname{argmin}_{w,b} \left( \frac{1}{2} \|\vec{w}\|_2^2 \right)$$

$$= \operatorname{argmin}_{w,b} \left( \frac{1}{2} \|\vec{w}\|_2^2 - \sum_n \alpha_n (k_n(\vec{w} \cdot \vec{x}_n + b) - 1) \right)$$

$$= \operatorname{argmin}_{w,b} (\mathcal{L})$$

The  $\alpha_n$  are positive, they are Lagrange multipliers.

We can now solve for  $\vec{w}, b$ :

$$\vec{\nabla}_w \mathcal{L} = \vec{0} \Leftrightarrow \frac{\vec{w}}{2} - \sum_n \alpha_n k_n(\vec{x}_n) = \vec{0}$$

$$\Leftrightarrow \vec{w} = \sum_n \alpha_n k_n \vec{x}_n$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Leftrightarrow \sum_n \alpha_n k_n = 0$$

We now need to find the  $\alpha_n$ , so we just moved (re-wrote) the problem. At its  $\vec{\nabla} \mathcal{L} = \vec{0}$  point (extremum),  $\mathcal{L}$  can be written: (we replace  $\vec{w}$  with its expression in  $\alpha_n$ )

$$\mathcal{L} = \frac{1}{2} \|\vec{w}\|_2^2 - \sum_n \vec{w} \cdot \alpha_n k_n \vec{x}_n + b \sum_n \alpha_n k_n + \sum_n \alpha_n$$

$$= \frac{1}{2} \vec{w} \cdot \vec{w} - \vec{w} \cdot \vec{w} + b \sum_n \alpha_n k_n + \sum_n \alpha_n$$



$$= -\frac{1}{2} \left( \sum_n \alpha_n x_n t_n \right) \left( \sum_m \alpha_m x_m t_m \right) + 0 + \sum_n \alpha_n$$

$$= \sum_n \alpha_n - \frac{1}{2} \sum_n \sum_m \left( \alpha_n \alpha_m t_n t_m \underbrace{\vec{x}_n \cdot \vec{x}_m}_{K(\vec{x}_n, \vec{x}_m)} \right)$$

↑  
dual form.

↙  $\vec{x}_n \cdot \vec{x}_m$  can be replaced with  $K(\vec{x}_n, \vec{x}_m)$

Note: this is a convex problem, in  $\alpha$  !!

↓  
This is then the kernelized SVM  
Note that we needed this dual form.

At prediction time, we would have:

$$y^{\text{pred}} = \text{sign}(\vec{w} \cdot \vec{x}_{\text{test}}) = \text{sign} \left( \sum_n \alpha_n t_n \underbrace{\vec{x}_n \cdot \vec{x}_{\text{test}}}_{K(\vec{x}_n, \vec{x}_{\text{test}})} \right)$$

we can insert  $K(\vec{x}_n, \vec{x}_{\text{test}})$  here

• The goal is then to find the best  $\alpha_n$ .

KKT: 1) primal pbm  $y \cdot d' \geq 1$ ; 2)  $\alpha \geq 0$ ; 3) Gradient = 0;

4) for inequalities, either  $\alpha = 0$  or the constraint is saturated,  $\alpha \cdot (t \cdot d' - 1) = 0$

• The KKT conditions tell us things about the  $\alpha_n$  that do define the  $\vec{w}_{\text{SVM}}^*$ ,  $\alpha_n^* = \underset{\alpha_1, \dots, \alpha_N}{\text{argmin}} (L)$

$$\text{KKT: } \begin{cases} \alpha_n \geq 0 \\ t_n d'(\alpha_n) - 1 \geq 0 \quad \Leftrightarrow \quad t_n (\vec{w} \cdot \vec{x}_n + b) \geq 1 \\ \alpha_n (t_n d'(\alpha_n) - 1) = 0 \quad (\Rightarrow) \quad \alpha_n (t_n (\vec{w} \cdot \vec{x}_n + b) - 1) = 0 \end{cases}$$

(  $d'(\alpha_n) = \frac{\vec{w} \cdot \vec{x}_n + b}{1}$  )

i.e the  $\alpha_n$  are Lagrange multipliers,  $\alpha_n \geq 0$

i.e each inequality is either saturated ( $t_n d'(\alpha_n) - 1 = 0$ ),  
or the  $\alpha_n$  is 0 (practically ignored).  
↳ point  $\vec{x}_n$  is on the margin

There are 2 kind of points:

$$\begin{cases} \alpha_n = 0, & t_n d'(\alpha_n) - 1 > 0 : \text{points well. class outside the margin} \\ \alpha_n > 0, & \text{---} = 0 : \text{on the margin} \end{cases}$$

→ These are the support vectors!



Note that the sol<sup>0</sup> is  $w = \sum_n \alpha_n k_n x_n = \sum_{S.V.}$   
a linear combinat<sup>n</sup> of the support vectors only!

$b$  can be computed: