# Support Vector Machines . Vladimir Vapnik 1994   ①
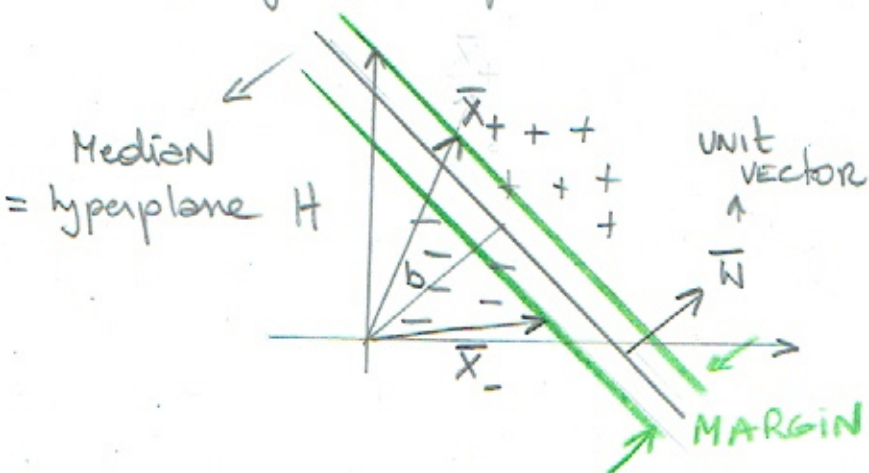
↳ classification algorithm (not so good for regression)
| VERY popular
| VERY powerful.



Median = hyperplane H

UNIT VECTOR

$\overline{W}$

MARGIN

**PERCEPTRON THEOREM**

If there is a hyperplane that can separate the classes there are ∞ many!

↓

which one is the BEST ?|

↓

**SVM** : use hyperplane that maximizes the margin between the classes

$$W^T x + b = 0 \quad \text{ON } H$$

| $\overline{W}$ IS $\perp H$
| $\overline{W}$ IS UNIT VECTOR.

ASSUME 2 points $x'$ AND $x''$ ON H.

$$\rightarrow \quad W^T x' + b = 0$$
$$\underline{W^T x'' + b = 0}$$
$$W^T \cdot (x' - x'') = 0 \quad \rightarrow \text{MEANING } W^T \text{ AND } (x' - x'')$$
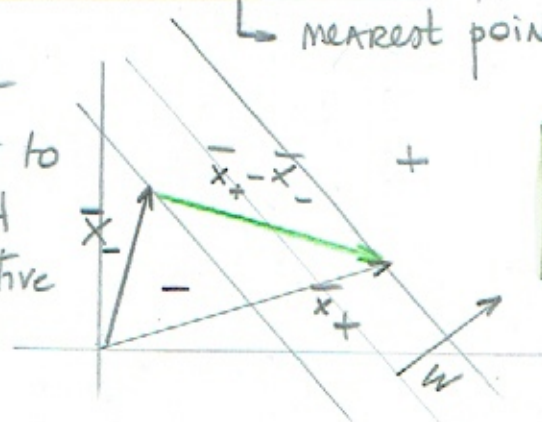ARE ORTHOGONAL.

$\rightarrow W^T$ IS ORTHOGONAL to any VECTOR on the plane

$\rightarrow W^T \perp H.$

distance of VECTOR $\begin{cases} \overline{X}_- \\ \overline{X}_+ \end{cases}$ to H ?

↳ nearest point to the plane !

$\overline{X}_+ , \overline{X}_-$ ARE vectors closest to the plane H in their respective classes.



$$\boxed{\text{MARGIN} = (\overline{X}_+ - \overline{X}_-) \cdot \frac{\overline{W}}{\|W\|}}$$

set : $\begin{cases} \bar{w}^T \cdot \bar{x}_+ + b \geq 1 \\ \bar{w}^T \cdot \bar{x}_- + b \leq -1 \end{cases} \rightarrow y_i(\bar{w} \cdot \bar{x} + b) \geq 1$    $\downarrow$ set    ②

$\phantom{set :}$ $\hookrightarrow y_i = +1$ in $+$
$\phantom{set : aaaaaaa}$ $y_i = -1$ in $-$

(1) set : $\boxed{y_i(\bar{w}^T \cdot \bar{x}_i + b) = 0}$

$\phantom{(1) set :}$ $\hookrightarrow$ IN GUTTER

(1) $\rightarrow$ $\bar{w}^T \cdot \bar{x}_+ = 1 - b \rightarrow \bar{x}_+ = \dfrac{1-b}{w}$

$\phantom{(1) \rightarrow}$ $\bar{w}^T \cdot \bar{x}_- \leq -1 - b \rightarrow \bar{x}_- = -\dfrac{(1+b)}{w}$

$\rightarrow$ MARGIN $= \left[\left(\dfrac{1-b}{w}\right) + \left(\dfrac{1+b}{w}\right)\right] \cdot \dfrac{\bar{w}}{\|w\|}$

$\phantom{\rightarrow MARGIN} = \dfrac{2}{\|w\|}$

GOAL $\rightarrow$ MAX $\dfrac{2}{\|w\|}$ $\iff$ MAX $\dfrac{1}{\|w\|}$ $\iff$ min $\|w\|$

$\phantom{GOAL \rightarrow MAX aaaaaaaaaaaaa}$ $\iff$ $\boxed{\text{min } \dfrac{1}{2}\|w\|^2}$

$\phantom{GOAL \rightarrow}$ subject to : $y_i(\bar{w}^T \cdot \bar{x} + b) \geq 1$

$\rightarrow$ $\underline{\text{CONSTRAINT optimization Problem}}$ $\rightarrow$ Lagrange
$\phantom{\rightarrow}$ $\llcorner$ INEQUALITY constraints $\phantom{aaaaaa}$ multipliers

$\phantom{\rightarrow aaaaaaa}$ $\hookrightarrow$ KKT (KARUSH – KUHN – TUCKER)

SVM : minimize $\dfrac{1}{2}\|w\|^2$ with condition that all points
$\phantom{SVM :}$ are correctly classified!

$$\min \frac{1}{2} \|w\|_2^2 \quad \text{with constraint} \quad y_i(w^T x_i + b) \geq 1$$

$$- \sum \alpha_i \left[ y_i(w^T x_i + b) - 1 \right]$$

* $\boxed{L(w, b, \alpha) = \frac{1}{2}\|w\|_2^2 - \sum \alpha_i \left( y_i(w^T x_i + b) - 1 \right)}$

MINIMIZE w, b
MAXIMIZE $\alpha$.

WITH $\alpha_i \geq 0$ (because of INEQUALITY constraint)

instead of MINIMIZING over w, b subject to constraints involving $\alpha$, we can MAXIMIZE over $\alpha$ subject to relations obtained previously for w and b.

$$\frac{\partial L}{\partial w} = \bar{w} - \sum \alpha_i y_i x_i \overset{\text{SET}}{=} 0 \text{ (VECTOR)}$$

$$\longrightarrow \boxed{\bar{w} = \sum \alpha_i y_i \bar{x}_i}$$

$$\frac{\partial L}{\partial b} = \boxed{\sum \alpha_i y_i \overset{\text{set}}{=} 0} \quad (**)$$

→ PUT IN (*)
→ free of w, b

(*) $L = \frac{1}{2} \left( \sum \alpha_i y_i \bar{x}_i \right) \left( \sum \alpha_j y_j x_j \right)$

$$- \left( \sum \alpha_i y_i \bar{x}_i \right) \cdot \left( \sum \alpha_j y_j \bar{x}_j \right)$$

$$- \sum \alpha_i y_i b \longrightarrow 0 \text{ AS PER } (**)$$

$$+ \sum \alpha_i$$

we can train a classifier in high D without computing $\bar{w}, b$ that define (***) H

MAX $\alpha_i \alpha_j$

$$\boxed{L(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j \, y_i y_j \left( \bar{x}_i \cdot \bar{x}_j \right)}$$

w, b free!

subject to: $\alpha_i \geq 0 \quad \sum \alpha_i y_i = 0$
KKT

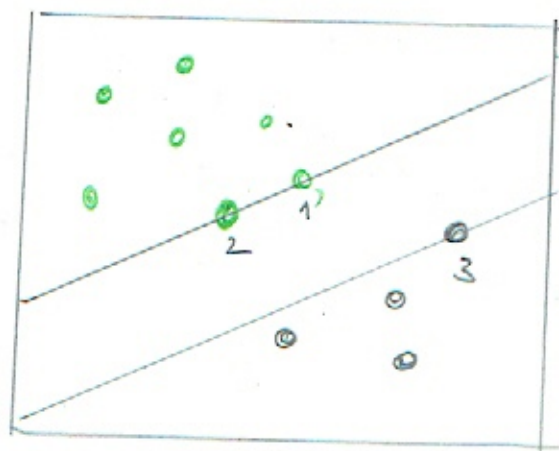↳ MAX only depends on on dot product of PAIRS of samples

SUPPORT VECTORS    a

IN GUTTER $\rightarrow \quad y_i(\bar{w}^T x_i + b) - 1 = 0$

$\Rightarrow \quad L(\alpha) \quad \boxed{\alpha > 0} \quad \Rightarrow \quad \bar{w} = \sum_{x_i = SV} \alpha_i y_i x_i$

→ All points <u>outside gutter</u> do not contribute
   to MAX $L(\alpha)$     $\hookrightarrow \alpha = 0$
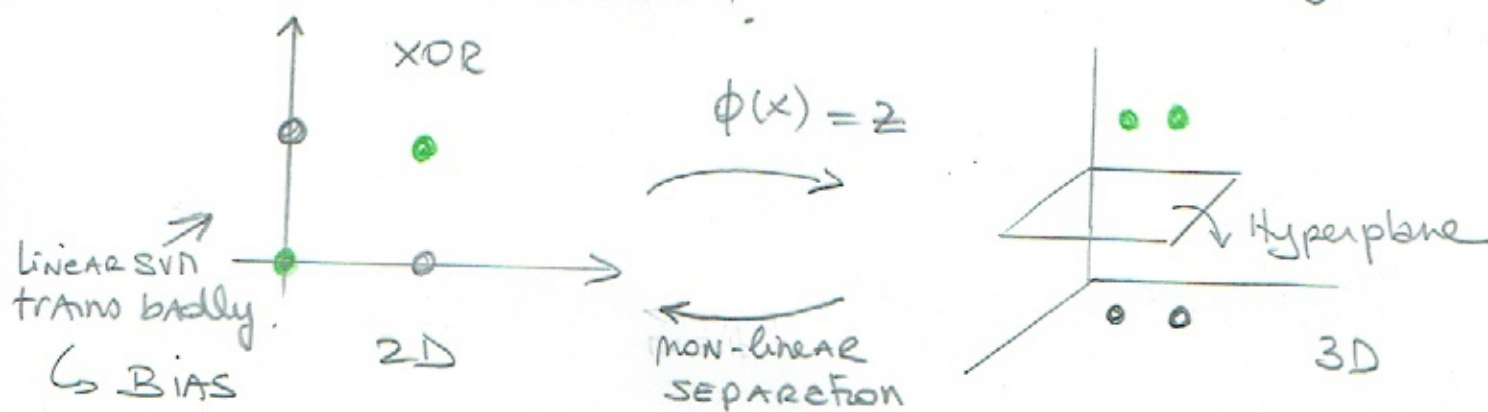
→ Robust WITH Respect to outliers!



3 SUPPORT VECTORS!
↓
Achieve the
MARGIN !

$\bar{w} = \sum_{SV} \alpha_i y_i x_1$

W has $d$ dimensions
<u>but</u>: only 3 $\alpha$s contribute
IN above example!

up to now we worked WITH
linearly SEPARABLE DATA.

→ what happens if the DATA IS NOT linearly
   separable in 2D?



XOR

$\phi(x) = z$

Hyperplane

LINEAR SVM
trains badly.
$\hookrightarrow$ BIAS

2D

non-linear
separation

3D

2D → 3D ⟹ $d\uparrow$ ⟹ computationally more expensive

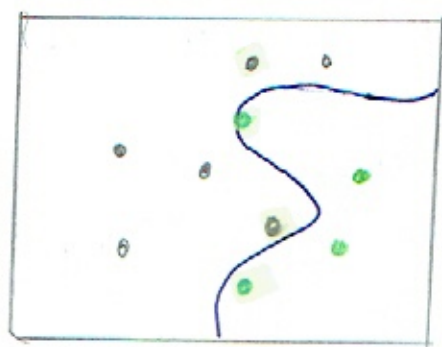$$\binom{***}{p3} \rightarrow \boxed{\angle(\alpha) = \sum x_i - \frac{1}{2}\sum_i\sum_j y_i y_j \alpha_i \alpha_j \; z_i^T z_j}$$

<u>what if</u> we go from 2D → $10^6$ D

* $\alpha$s ARE related to NUMBER of TRAINING SAMPLES AND NOT dependent ON dimension of $\bar{z}$

* $z_i^T z_j$ is not really A CONCERN even if dimension $z$ is $10^6$

<u>NOW</u>: ASSUME we found hyperplane in Z - SPACE

→ what happens IN 2D - space

↳ support vectors live IN z-space !

↓

we know which are the vectors that are support vectors! $\alpha \neq 0$



look for them in 2D

## KERNEL trick

$$\boxed{K(x,y) = \langle \phi(x), \phi(y)\rangle} \rightarrow \text{dot product} = \text{SCALAR}$$

$x, y \in \mathbb{R}^d$

↓ $\phi$ is function

$d \longrightarrow m$

with $m >>> d$

to calculate $\langle \phi(x), \phi(y)\rangle$ we would calculate $\phi(x), \phi(y)$ first and then do dot product

WITH KERNEL ⟹ no need to go to m-dim space !

example· $x = \begin{pmatrix} x_1 & 1 \\ x_2 & 2 \\ x_3 & 3 \end{pmatrix}$ $y = \begin{pmatrix} y_1 & 4 \\ y_2 & 5 \\ y_3 & 6 \end{pmatrix}$ $\in \mathbb{R}^{3 = d}$

$\phi(x) = (x_1 x_1, x_1 x_2, x_1 x_3, x_2 x_1, x_2 x_2, x_2 x_3, x_3 x_1, x_3 x_2$

$\boxed{k(x,y) = (\langle x, y \rangle)^2}$ $\qquad x_3 x_3)$

$\qquad\qquad\qquad\qquad \hookrightarrow$ 9 dot products

$\phi(x) = (1, 2, 3, 2, 4, 6, 3, 6, 9)$ ↙

$\phi(y) = (16, 20, 24, 20, 25, 30, 24, 30, 36)$ ↙

$k(x,y) = \langle \phi(x), \phi(y) \rangle = 16 + 40 + 72 + 40 + 100 + 180 + 72$

$\qquad\qquad\qquad\qquad\qquad + 180 + 324$

$\qquad\qquad\qquad\qquad = 1024$

$3D \longrightarrow 9D$ $\qquad\qquad\qquad\qquad$ SAME Result but much FASTER!

$k(x,y) = (4 + 10 + 18)^2 = 1024$

KERNEl functions

$k(x,y) = x^T y$ $\quad$ linear

$k(x,y) = (1 + x^T y)^P$ $\quad$ polynomiAL

$k(x,y) = e^{\frac{-(x-y)^2}{\sigma^2}}$ $\quad$ RBF (RAdiAl BASE Function)

VERy POPUlAR !
works out of the box
EVERy problem becomes
linearly separable