

Máquinas de vectores soporte (Support Vector Machine)

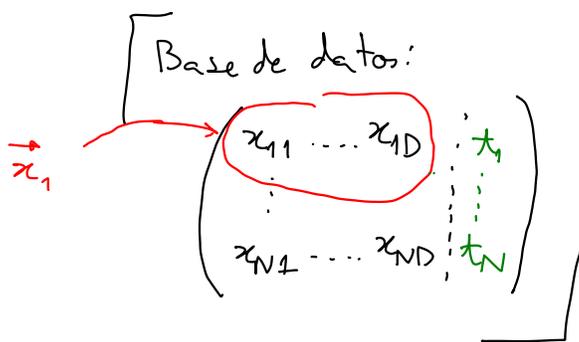
Dataset: Conjunto de datos clasificados en 2 clases

$$\mathcal{L} = \left\{ (\vec{x}_m, t_m) \mid m=1, \dots, N \right\}$$

↑
FEATURES

↑
TARGET

↑
NÚMERO DE INSTANCIAS
(FILAS DE LA B.O.)



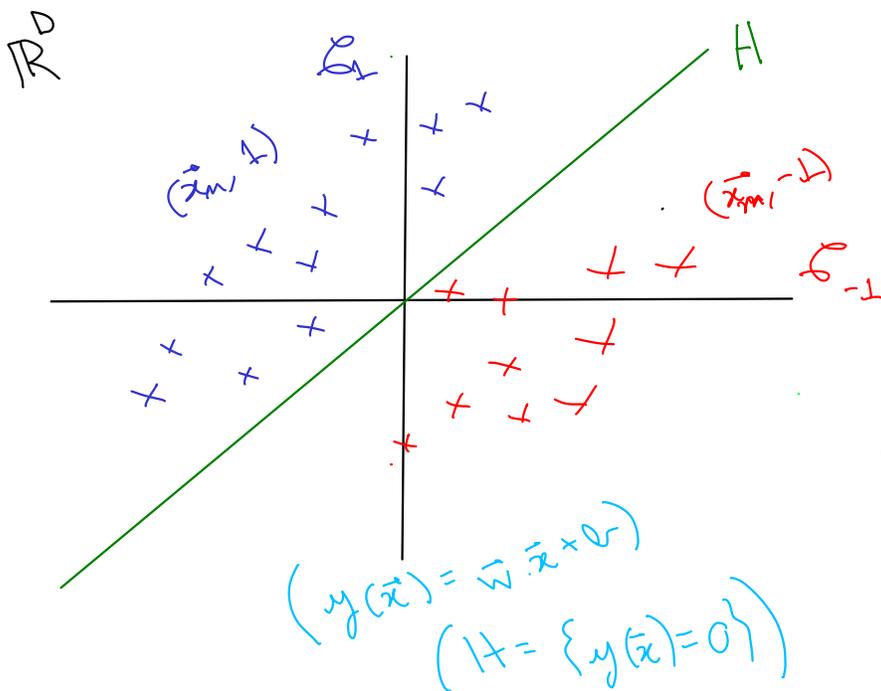
$$\vec{x}_m \in \mathbb{R}^D$$

CLASE \mathcal{L}_1

$$t_m \in \{1, -1\}$$

CLASE \mathcal{L}_{-1}

$$\mathcal{L} = \mathcal{L}_1 \cup \mathcal{L}_{-1}$$



$$H = \left\{ \underbrace{\vec{w}}_{\mathbb{R}^D} \cdot \underbrace{\vec{x}}_{\mathbb{R}^D} + \underbrace{b}_{\mathbb{R}} = 0 \right\}$$

$$(w_1 x_1 + \dots + w_n x_n + b = 0)$$

Un hiperplano es un hiperplano separador si:

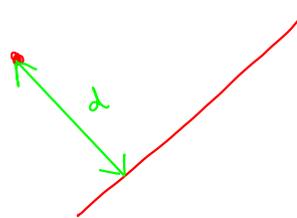
$$\mathcal{L}_1 \in \{y(\vec{x}) > 0\}$$

$$\mathcal{L}_{-1} \in \{y(\vec{x}) < 0\}$$

es decir,

$$t_m \cdot y(\vec{x}_m) > 0 \quad \forall m$$

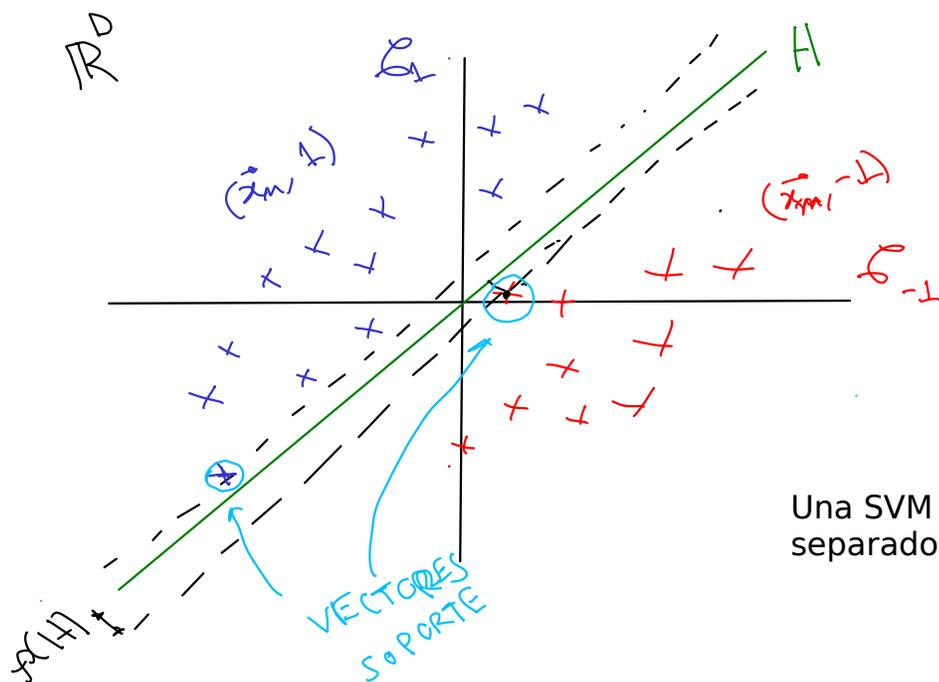
$$t_m \cdot y(\vec{x}_m) > 1$$



El margen de un hiperplano separador H es la distancia del conjunto de datos a H:

$$p(H) = \min_m d(\bar{x}_m, H)$$

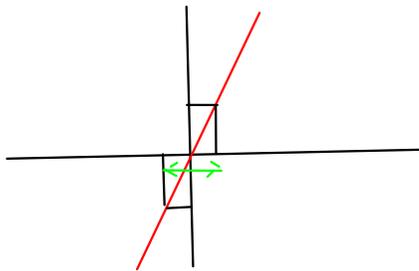
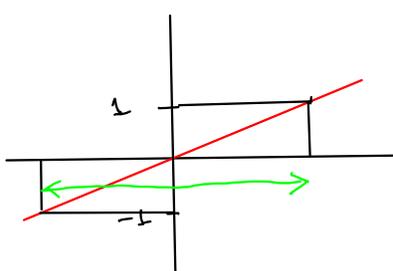
↑
MARGEN



Una SVM calcula un hiperplano separador maximizando su margen.

Es decir, busquemos \vec{w} tal que maximice $p(\{\vec{w} \cdot \vec{x} + b = 0\})$
 sujeto a $\boxed{\tan(\vec{w} \cdot \vec{x} + b) \geq 1}$.
 ← COND. DE HIPERPLANO SEPARADOR

Observación: $\|\vec{w}\|$ da la pendiente del hiperplano



A mayor pendiente, menor margen

Nuestro problema es equivalente a:

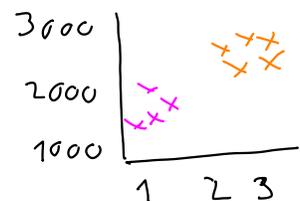
$$\left\{ \begin{array}{l} \text{minimizar} \\ \text{sujeta a} \end{array} \right. \quad \frac{1}{2} \|\vec{w}\|^2 = \frac{1}{2} \vec{w} \cdot \vec{w} = w_1^2 + \dots + w_D^2$$
$$k_m(\vec{w} \cdot \vec{x}_m + b) \geq 1 \quad \forall m$$

PROBLEMAS DE LA SVM:

- Sensible a outliers



- Sensible a escalas



- Estamos asumiendo que el conjunto es

LINEALMENTE SEPARABLE

(es decir, que existe un hiperplano separador)

Posibles soluciones para el caso no linealmente separable:

(Para un caso que es "casi" linealmente separable):

SOFT MARGIN

- Introducimos VARIABLES DE HOLGURA ("SLACK")

$$h_m \geq 0$$

• Sustituimos las condiciones $t_n y(\vec{x}_n) \xrightarrow{\forall n} t_n y(\vec{x}_n) \geq 1 - \ell_n$
 $m = 1, \dots, N$

• Sustituyo la función de coste $\frac{1}{2} \|\vec{w}\|^2 \xrightarrow{} \frac{1}{2} \|\vec{w}\|^2 + C \sum_{n=1}^N \ell_n$

$$\begin{cases} \min \frac{1}{2} \|\vec{w}\|^2 + C \sum_{n=1}^N \ell_n \\ \text{s.a. } t_n y(\vec{x}_n) + \ell_n \geq 1 \\ \ell_n \geq 0 \end{cases}$$

C es un HIPERPARÁMETRO que controla la penalización que se impone a las variables de holgura

C alto: Margen pequeño con pocas violaciones
 C bajo: Margen ancho con muchas violaciones

SVM NO LINEAL

El problema dual

Asociado al problema (PRIMAL) $\left\{ \begin{array}{l} \min \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.a. } t_n (\vec{w} \cdot \vec{x}_n + b) \geq 1 \quad \forall n \end{array} \right.$

Hay un problema DUAL equivalente $\left\{ \begin{array}{l} \text{minimizar } \frac{1}{2} \sum_{n=1}^N a_n t_n t_n \boxed{\vec{x}_n \cdot \vec{x}_n} \\ \text{s.a. } a_n \geq 0 \quad \forall n = 1, \dots, N \\ \sum_{n=1}^N a_n t_n = 0 \end{array} \right.$
 $-\sum_{n=1}^N a_n$

\hat{a} óptimos $\xrightarrow{} \left\{ \begin{array}{l} \hat{\vec{w}} = \sum_{n=1}^N \hat{a}_n t_n \vec{x}_n \\ \hat{b} = \frac{1}{\#S} \sum_{n \in S} t_n - \hat{\vec{w}} \cdot \vec{x}_n \end{array} \right.$

COSTES COMPUTACIONALES

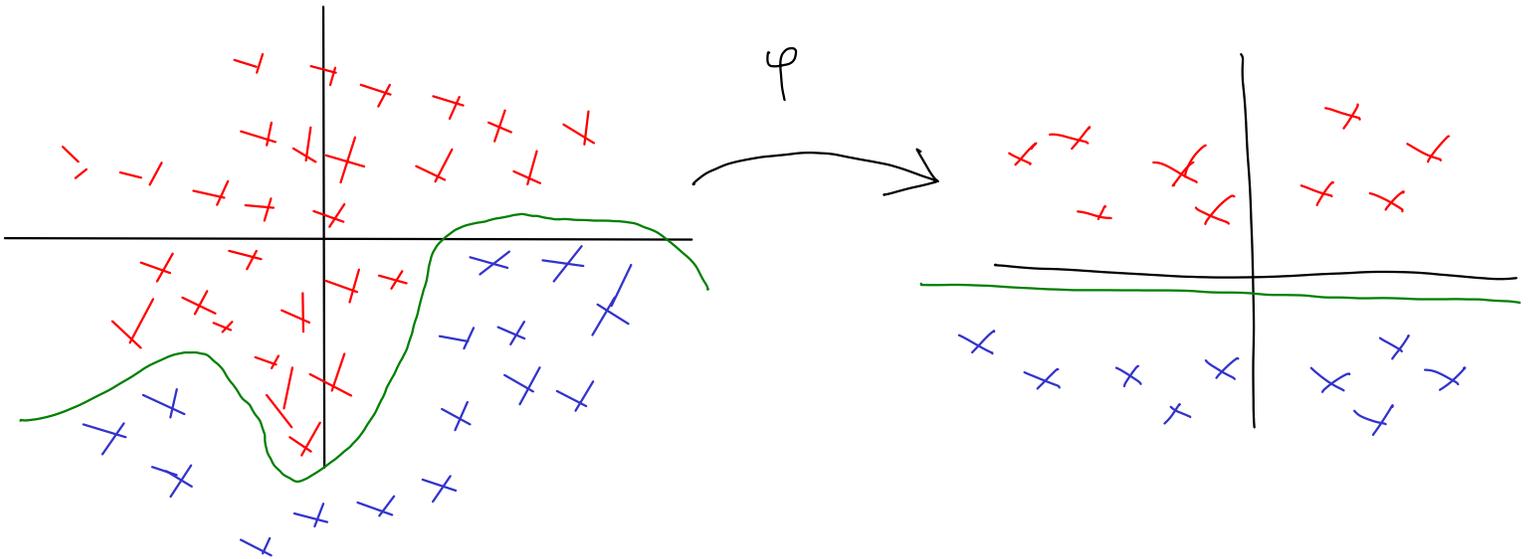
PRIMAL $\rightarrow O(D^3)$ $D = \# \text{ FEATURES}$

DUAL $\rightarrow O(N^3)$ $N = \# \text{ INSTANCIAS}$

$$S = \{n \mid \hat{a}_n > 0\}$$

Pero para lo que realmente nos interesa el problema dual es para poder usar el "kernel trick".

NUCLEOS



Supongamos que tenemos un training set $\{\vec{x}_n\} \subset \mathbb{R}^2$ y le aplicamos una transformación polinomial de grado 2

$$\mathbb{R}^2 \xrightarrow{\varphi} \mathbb{R}^3$$

$$(x, y) \mapsto (x^2, \sqrt{2}xy, y^2)$$

para "linealizar" el conjunto de datos.

Sucede la siguiente

$$\boxed{\varphi(\vec{a}) \cdot \varphi(\vec{b}) = (a_1^2, \sqrt{2}a_1a_2, a_2^2) \cdot (b_1^2, \sqrt{2}b_1b_2, b_2^2) = a_1^2b_1^2 + 2a_1a_2b_1b_2 + a_2^2b_2^2 = (a_1b_1 + a_2b_2)^2 = (\vec{a} \cdot \vec{b})^2}$$

En el problema dual, basta reemplazar $\vec{x}_i \cdot \vec{x}_n$ por $(\vec{x}_i \cdot \vec{x}_n)^2$
 ¡no hace falta aplicar la transformación φ !

$K(\vec{a}, \vec{b}) = (\vec{a} \cdot \vec{b})^2$ se llama NÚCLEO POLINOMIAL DE GRADO 2 (O CUADRÁTICO)

En general: $K: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$
 $(\vec{x}_n, \vec{x}_m) \mapsto K(\vec{x}_n, \vec{x}_m)$

es un núcleo si existe una transformación

$\varphi: \mathbb{R}^D \rightarrow \mathbb{R}^P$ tal que

$$K(\vec{x}_n, \vec{x}_m) = \varphi(\vec{x}_n) \cdot \varphi(\vec{x}_m)$$

En resumen, utilizar un núcleo en la función de coste es equivalente a hacer una transformación del dataset.

NÚCLEOS COMUNES:

• Lineal: $K(\vec{a}, \vec{b}) = \vec{a} \cdot \vec{b}$

• Polinomial: $K(\vec{a}, \vec{b}) = (\gamma \cdot \vec{a} \cdot \vec{b} + r)^d$
(de grado d)

• Gaussianos: $K(\vec{a}, \vec{b}) = \exp(-\gamma \|\vec{a} - \vec{b}\|^2)$
o radial

• Sigmoide: $K(\vec{a}, \vec{b}) = \tanh(\gamma \vec{a} \cdot \vec{b} + r)$

RED NEURONAL

Se usa, pero no es un verdadero núcleo