

# Máquinas de Vectores Soporte (SVM)

**Alfredo Cuesta Infante**

E. T. S. Ingeniería Informática  
Universidad Rey Juan Carlos

Master Univ. en Visión Artificial  
**Reconocimiento de Patrones**

## Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

## Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

## Clasificación no lineal

- El truco del Kernel

Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

Clasificación no lineal

- El truco del Kernel

## Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

## Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

## Clasificación no lineal

- El truco del Kernel

### Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

### Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

### Clasificación no lineal

- El truco del Kernel

## ¿Qué clasificador es mejor?

- ▶ Los clasificadores  $C_1$ ,  $C_2$  y  $C_3$  clasifican PERO...  
¿qué ocurrirá cuando llegue un nuevo ejemplo próximo a la superficie de decisión?
- ▶ Cuanto mayor sea el **margen** menor es el **riesgo** de FPs y FNs

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

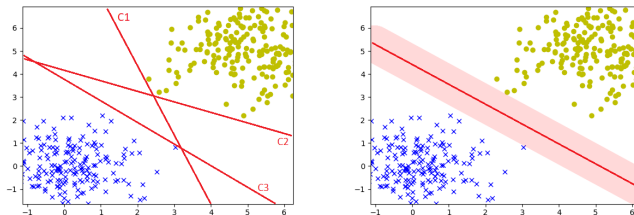
Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel



**Figura:** (Der.) Tres clasificadores que no cometen ningún error en el conjunto de entrenamiento pero seguramente funcionen mal con datos nuevos. (Izq.) Este clasificador que maximiza el margen entre ambas clases. [Fuente: Original de A. Cuesta]

“La justicia inflexible es frecuentemente la injusticia más grande”

Publio Terencio Africano (194 a.C. – 159 a.C.)

- ▶ Si permitimos algunos ejemplos dentro de los márgenes, entonces podemos aumentarlos y encontrar clasificadores con mayor potencial generalizador

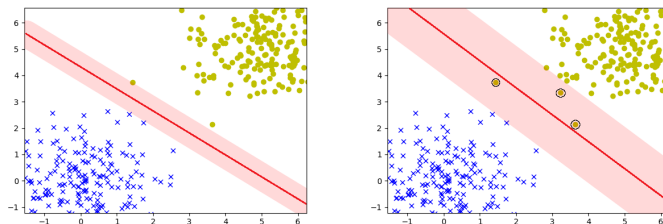


Figura: (Der.) Márgenes *hard*. (Izq.) márgenes *soft* [Fuente: Original de A. Cuesta]

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel

- ▶ Con SVM **no vamos a necesitar** transformar el vector de características a otro espacio.
- ▶ En su lugar *alguien* de ese espacio vendrá al nuestro con el resultado de las operaciones.



## Intuición

Riesgo estructural  
Márgenes  
Truco del Kernel

## Clasificación lineal

Notación  
Márgenes hard  
Márgenes soft  
El problema dual

## Clasificación no lineal

El truco del Kernel

Intuición

Riesgo estructural  
Márgenes  
Truco del Kernel

Clasificación lineal

Notación  
Márgenes hard  
Márgenes soft  
El problema dual

Clasificación no lineal

El truco del Kernel

## Notación

- Ahora los vectores  $\mathbf{w}$  y  $\mathbf{x}$  no estarán ‘expandidos’, es decir:

$$\mathbf{w} = (w_1, w_2, \dots, w_m)^T, \quad \mathbf{x} = (x_1, x_2, \dots, x_m)^T.$$

- Dado que hemos sacado el término independiente del vector de pesos, la expresión del discriminante es ahora:

$$\hat{t} = \begin{cases} +1 & \text{si } \mathbf{w}^T \mathbf{x} + b > 0 \\ -1 & \text{si } \mathbf{w}^T \mathbf{x} + b < 0 \end{cases} \quad (1)$$

donde  $b$  es el término independiente y, al igual que en la semana pasada, las etiquetas son  $+1$  y  $-1$  porque serán ‘matemáticamente convenientes’.

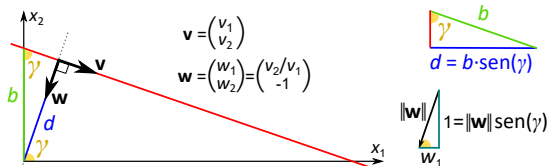
## Recordar también que en 2D

- Una recta se puede definir con un punto y un vector director  $\mathbf{v}$
- Un hiperplano se puede definir con un punto y un vector normal  $\mathbf{w}$
- El vector director y el normal son ortogonales:  $\mathbf{v}^T \mathbf{w} = \mathbf{w}^T \mathbf{v} = 0$
- La ecuación implícita de una recta:  $w_1 x_1 + w_2 x_2 + b = \mathbf{w}^T \mathbf{x} = 0$



## Planteamiento geométrico del problema

### 1. Calcular la distancia de la superficie de decisión al origen



Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel

**Figura:** El clasificador lineal viene dado por un vector director  $\mathbf{v}$  y un punto de corte  $b$ . El vector  $\mathbf{w}$  es perpendicular a  $\mathbf{v}$  y sirve para calcular la distancia de la recta al origen. Combinando todo podemos representar  $d$  en función de  $\mathbf{w}$  y  $b$ . [Fuente: Original de A. Cuesta]

En definitiva, la distancia del clasificador al origen es

$$d = \frac{b}{\|\mathbf{w}\|}$$



## Planteamiento geométrico del problema

### 3. En definitiva:

- ▶ Cuanto más pequeña sea la norma del vector de pesos  $\|\mathbf{w}\|$ , mayor será el margen  $\varepsilon$ ; es decir, buscamos

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{w}\|$$

- ▶ Las dos restricciones impuestas se pueden escribir en una sola:

$$\left. \begin{array}{ll} \text{Si } \hat{t}^{(i)} = +1 & \text{entonces } \mathbf{w}^{(i)T} \mathbf{x} + b \geq +1 \\ \text{Si } \hat{t}^{(i)} = -1 & \text{entonces } \mathbf{w}^{(i)T} \mathbf{x} + b < -1 \end{array} \right\} =$$

$$\left( \hat{t}^{(i)} \right) \left( \mathbf{w}^{(i)T} \mathbf{x} + b \right) \geq 1.$$

- ▶ Por tanto el problema de optimización es:

$$\text{minimizar } \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{sujeto a } \left( \hat{t}^{(i)} \right) \left( \mathbf{w}^{(i)T} \mathbf{x} + b \right) \geq 1,$$

$$\text{para } i = 1, 2, \dots, m.$$

$$\text{Recordar que } \|\mathbf{w}\| = \mathbf{w}^T \mathbf{w}$$

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel

## Introducción de variables *Slack*

- ▶ Permitimos que haya ejemplos en el margen utilizando un peso  $\zeta$  que incorporamos a la función de coste

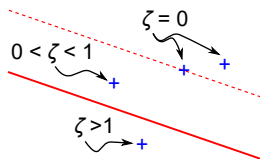


Figura: *Slack variables* para ejemplos que violan la restricción. [Fuente: Original de A. Cuesta]

### 1. Compromiso

- ▶ Queremos ampliar el margen
- ▶ Queremos minimizar el número de ejemplos que hay dentro del margen

Para gestionarlo se añade el hiperparámetro  $C$ :

$$\min_{\mathbf{w}, b, \zeta} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \zeta^{(i)}$$

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel

## Introducción de variables *Slack*

### 2. Modificación de las restricciones

$$\left(\hat{\mathbf{t}}^{(i)}\right)\left(\mathbf{w}^{(i)T} \mathbf{x} + b\right) \geq 1 - \zeta^{(i)}$$

### 3. En definitiva

$$\begin{array}{ll} \underset{\mathbf{w}, b, \zeta}{\text{minimizar}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \zeta^{(i)}, \\ \text{sujeto a} & \left\{ \begin{array}{l} \left(\hat{\mathbf{t}}^{(i)}\right)\left(\mathbf{w}^{(i)T} \mathbf{x} + b\right) \geq 1 - \zeta^{(i)} \\ \zeta^{(i)} \geq 0 \\ i = 1, 2, \dots, m. \end{array} \right. \end{array} \quad (2)$$

- ▶  $C$  es una medida de cuanto queremos evitar que haya ejemplos en el margen. Cuanto más pequeño, menos queremos evitarlo, es decir más *permisivos* somos.
- ▶ Problema de optimización cuadrática con restricciones !!
- ▶ Solución mediante **multiplicadores de Lagrange**

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel

## El problema dual

## Multiplicadores de Lagrange

- Técnica mediante la cual el problema se formula en términos de unas nuevas variables  $\alpha^{(i)}$  con restricciones más sencillas.

$$\begin{aligned} & \underset{\alpha}{\text{minimizar}} && \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^{(i)} \alpha^{(j)} \mathbf{t}^{(i)} \mathbf{t}^{(j)} \mathbf{x}^{(i)T} \mathbf{x}^{(j)} - \sum_{i=1}^m \alpha^{(i)} \\ & \text{sujeto a} && \alpha^{(i)} > 0 \\ & && i = 1, 2, \dots, m. \end{aligned} \quad (3)$$

- ▶ Cada ejemplo  $\mathbf{x}^{(i)}$  tiene un multiplicador de Lagrange  $\alpha^{(i)}$ .
- ▶ El problema se puede resolver mediante métodos computacionales.
- ▶ Una vez resuelto, algunos (muchos) multiplicadores se anulan. Los ejemplos asociados a los multiplicadores distintos de cero son los **vectores soporte**.
- ▶ El vector de pesos óptimo es:

$$\mathbf{w}^* = \sum_{i=1}^m \alpha^{*(i)} \mathbf{t}^{(i)} \mathbf{x}^{(i)} ; \quad b^* = \frac{1}{n_s} \sum_{i: \alpha^{*(i)} > 0} \left( (1 - t^{(i)}) (\mathbf{w}^{*T} \mathbf{x}^{(i)}) \right) \quad (4)$$

$n_s$  es el número de vectores soporte, o sea el número de  $\alpha^{*(i)} > 0$ .

## Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

## Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

## Clasificación no lineal

- El truco del Kernel

### Intuición

- Riesgo estructural
- Márgenes
- Truco del Kernel

### Clasificación lineal

- Notación
- Márgenes hard
- Márgenes soft
- El problema dual

### Clasificación no lineal

- El truco del Kernel

## El truco del Kernel

- ▶ En la formulación dual, la función de coste depende del **producto escalar de cada dos vectores transformados**
- ▶ Pero si en vez de tener la transformación  $\phi$  tuvieramos una función  **$K$  (Kernel)** que nos devuelva el resultado de dicho producto escalar en el espacio al que nos lleva  $\phi$ , ¿podríamos resolver el problema?

$$\begin{aligned} \min_{\alpha} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^{(i)} \alpha^{(j)} t^{(i)} t^{(j)} \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) - \sum_{i=1}^m \alpha^{(i)} \\ \Downarrow \\ \min_{\alpha} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha^{(i)} \alpha^{(j)} t^{(i)} t^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \sum_{i=1}^m \alpha^{(i)} \end{aligned}$$

- ▶ Aparentemente **NO** porque la solución  $\{\mathbf{w}^*, b^*\}$  depende de  $\phi(\mathbf{x}^{(i)})$

$$\mathbf{w}^* = \sum_{i=1}^m \alpha^{*(i)} t^{(i)} \phi(\mathbf{x}^{(i)}) ; \quad b^* = \frac{1}{n_s} \sum_{i: \alpha^{*(i)} > 0} \left( (1 - t^{(i)}) (\mathbf{w}^{(*)T} \phi(\mathbf{x}^{(i)})) \right)$$

- ▶ Pero ¿para qué queremos saber  $\{\mathbf{w}^*, b^*\}$ ?

¿No es para clasificar nuevos ejemplos?

→ Vamos a ver qué aspecto tiene la función discriminante

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel



## El truco del Kernel

Sea  $\mathbf{z}$  un nuevo ejemplo. Entonces su etiqueta estimada es:

$$\begin{aligned}\hat{t} &= \mathbf{w}^{*T} \phi(\mathbf{z}) + b^* = \left( \sum_{i=1}^m \alpha^{*(i)} t^{(i)} \phi(\mathbf{x}^{(i)}) \right)^T \phi(\mathbf{z}) + b^* \\ &= \sum_{i: \alpha^{*(i)} > 0} \alpha^{*(i)} t^{(i)} \left( \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{z}) \right) + b^* \\ &= \sum_{i: \alpha^{*(i)} > 0} \alpha^{*(i)} t^{(i)} K(\mathbf{x}^{(i)}, \mathbf{z}) + b^*\end{aligned}$$

### ► ¡¡ Podemos clasificar ejemplos sin necesitar $\phi$ !!

Necesitamos:

- Resolver el problema dual planteado con un Kernel para obtener los multiplicadores no nulos
- Utilizar la expresión de arriba para clasificar nuevos ejemplos, donde el término independiente óptimo  $b^*$  es

$$b^* = \frac{1}{n_s} \sum_{i: \alpha^{*(i)} > 0} \left( (1 - t^{(i)}) \sum_{j: \alpha^{*(j)} > 0} (\alpha^{*(j)} t^{(j)} K(\mathbf{x}^{(j)}, \mathbf{z})) \right)$$

# Algunos Kernels importantes

Nombre	$K(\mathbf{a}, \mathbf{b}) =$	Hiperparámetros
Lineal	$\mathbf{a}^T \mathbf{b}$	— — —
Polinomial	$(\gamma \mathbf{a}^T \mathbf{b} + r)^d$	$\gamma, r, d$
Sigmoide	$\tanh(\gamma \mathbf{a}^T \mathbf{b} + r)$	$\gamma, r$
RBF Gaussiano	$\exp(-\gamma \ \mathbf{a} - \mathbf{b}\ ^2)$	$\gamma$

*RBF = Radial Basis Function*

Intuición

Riesgo estructural

Márgenes

Truco del Kernel

Clasificación lineal

Notación

Márgenes hard

Márgenes soft

El problema dual

Clasificación no lineal

El truco del Kernel