

# SUPPORT VECTOR MACHINES

**Santiago Alférez**

**APRENDIZAJE AUTOMÁTICO DE MÁQUINA**

**MACC – UR – EICT**

# Contenido

- Hiperplano
- Clasificación usando un hiperplano
- El clasificador de margen máximo
- El clasificador mediante vectores de soporte
- Ampliación de descriptores (polinomial)
- El truco del kernel
- Máquinas de vectores de soporte
- Clasificación con más de dos clases
- Implementación
- Máquinas de vectores de soporte para regresión

## Máquinas de vectores de soporte

El problema de clasificación binaria se aproxima directamente:

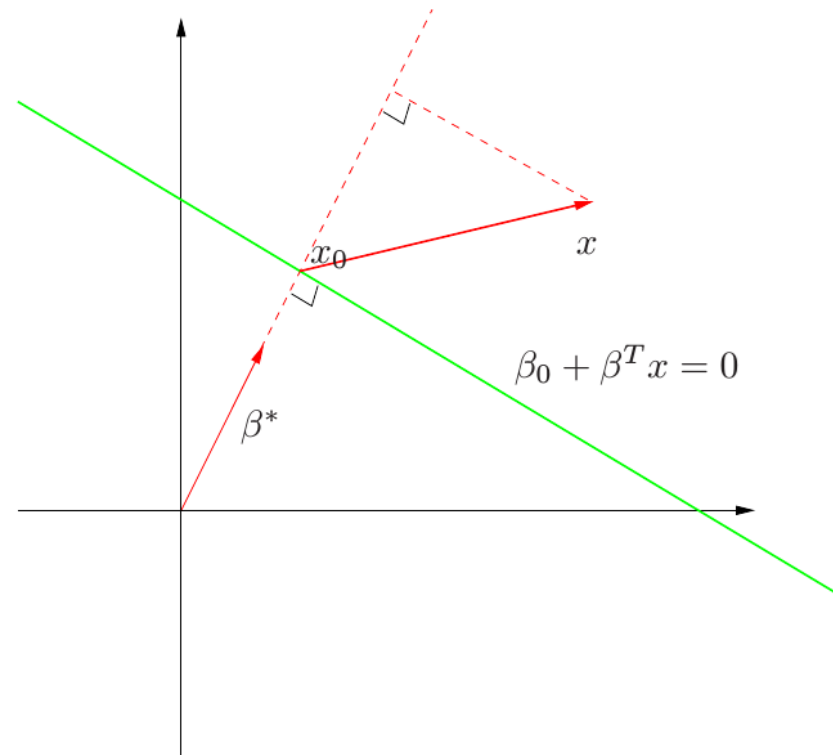
Al encontrar un plano que separe las clases en el espacio de descriptores.

Si no es posible, se puede hacer de dos formas:

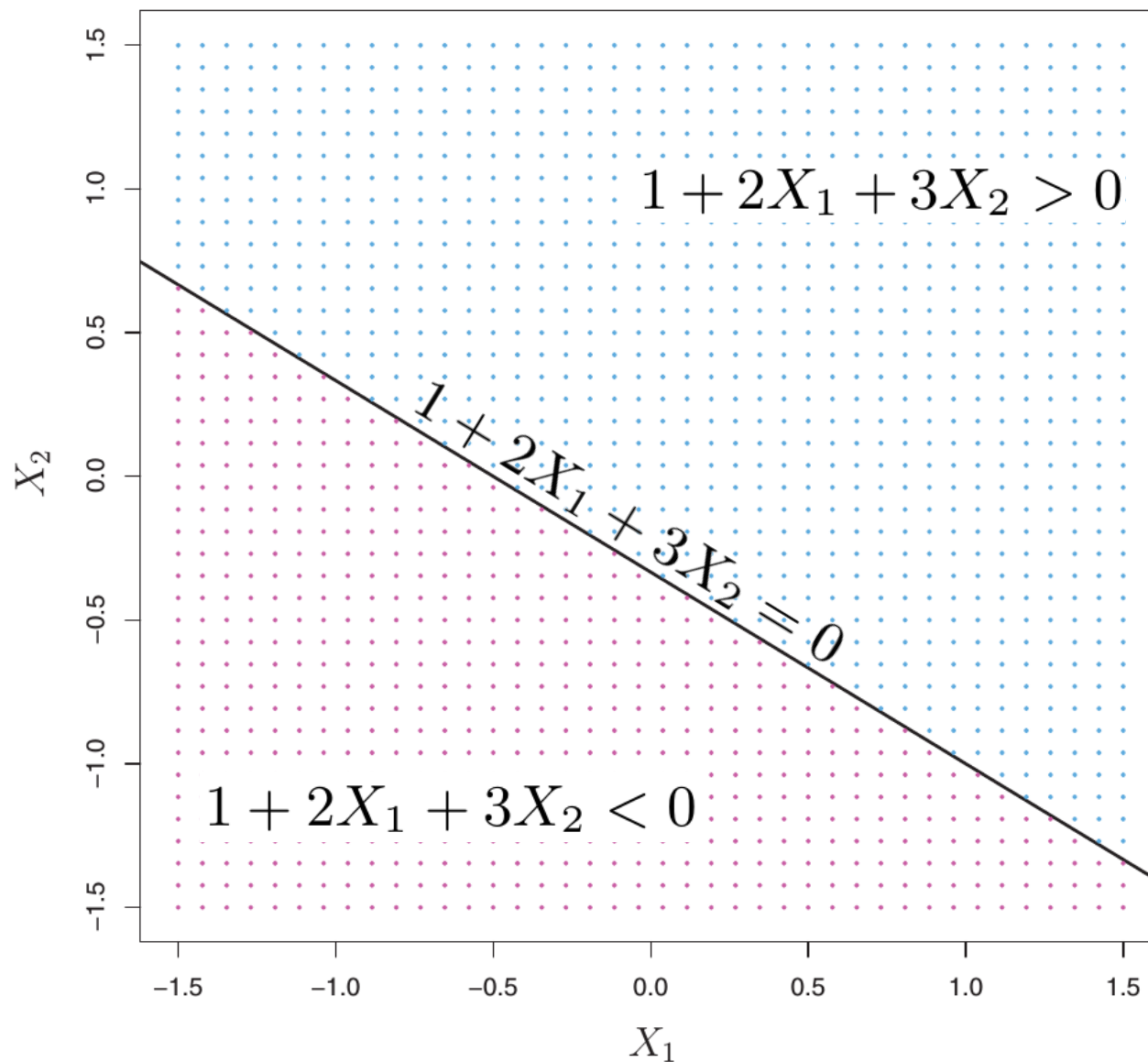
- Suavizando lo que quiere decir “que separe”, y
- Enriqueciendo y aumentando el espacio de descriptores tal que la separación sea posible.

# ¿Qué es un hiperplano?

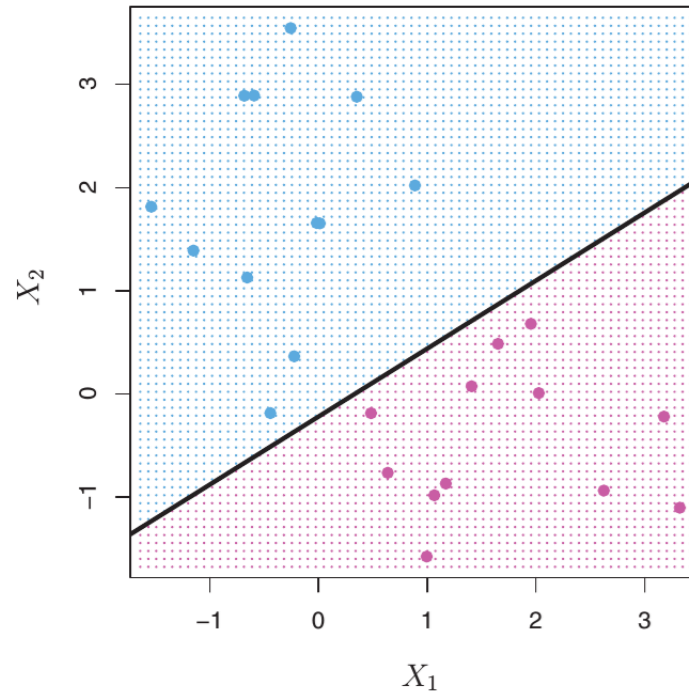
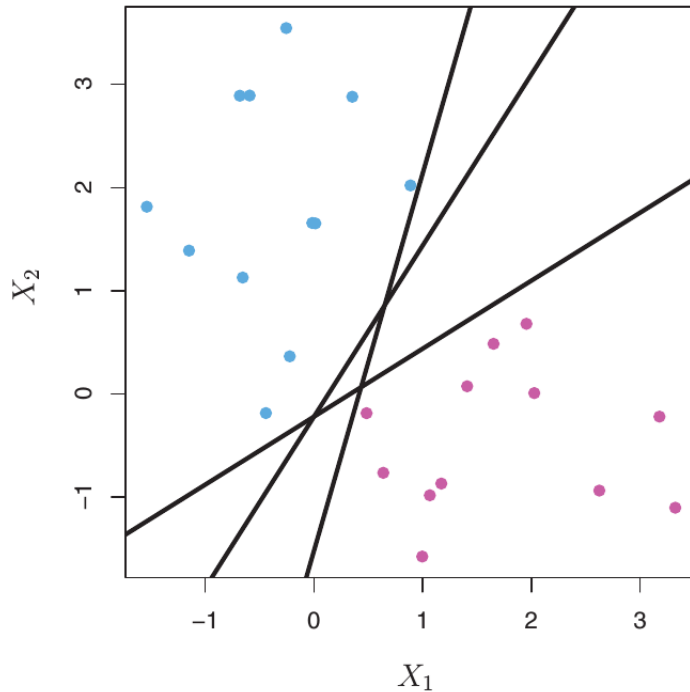
- ✓ Un hiperplano en  $p$  dimensiones es un subespacio plano de dimensión  $p - 1$ .
- ✓ En general, la ecuación para un hiperplano tiene la forma:
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0$$
- ✓ En  $p = 2$  dimensiones, un hiperplano es una línea.
- ✓ Si  $\beta_0 = 0$ , el hiperplano pasa por el origen, de lo contrario no.
- ✓ El vector  $\beta = (\beta_1, \beta_2, \dots, \beta_p)$  se denomina vector normal, apunta en una dirección ortogonal a la superficie del hiperplano.



## Ejemplo de un hiperplano



# Clasificación usando un hiperplano separador



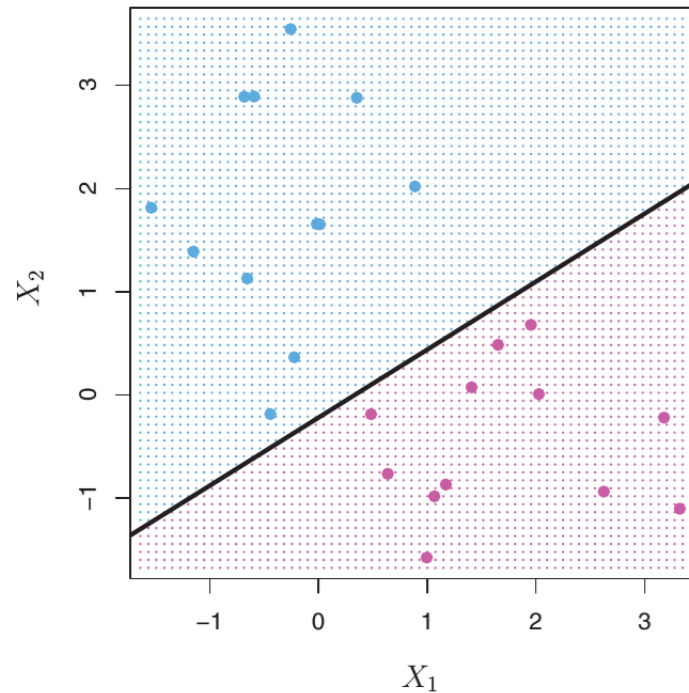
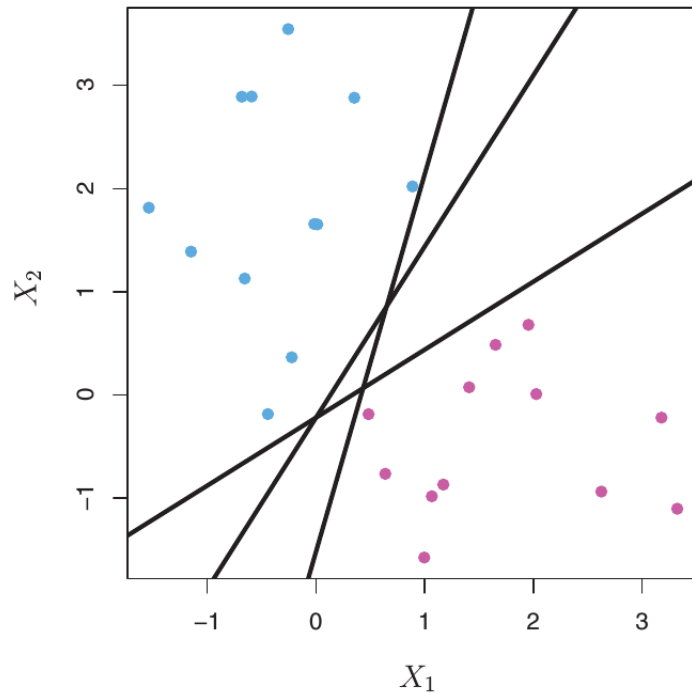
●  $\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} > 0$  if  $y_i = 1$

$\mathbf{x}^T \boldsymbol{\beta} + \beta_0$

$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} < 0$  if  $y_i = -1$  ●

$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) > 0$

# Clasificación usando un hiperplano separador

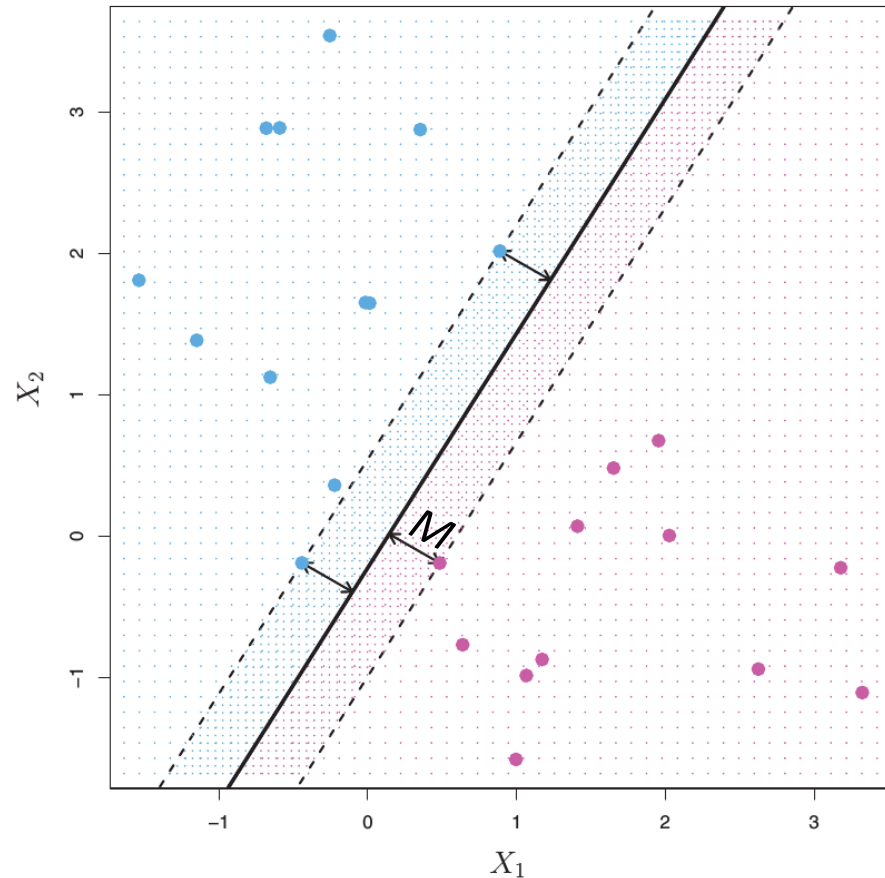


## Predicción

$$x^* \rightarrow f(x^*) = \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_p x_p^* \rightarrow \begin{cases} f(x) > 0 \Rightarrow x \in \text{clase 1} \\ f(x) < 0 \Rightarrow x \in \text{clase -1} \end{cases}$$

# El clasificador de margen máximo

El **margen** es la “calle” máxima que se puede colocar entre dos observaciones de diferentes clases.



El hiperplano con máximo margen es aquel que soluciona:

$$\max_{\beta_0, \beta_1, \dots, \beta_p} M$$

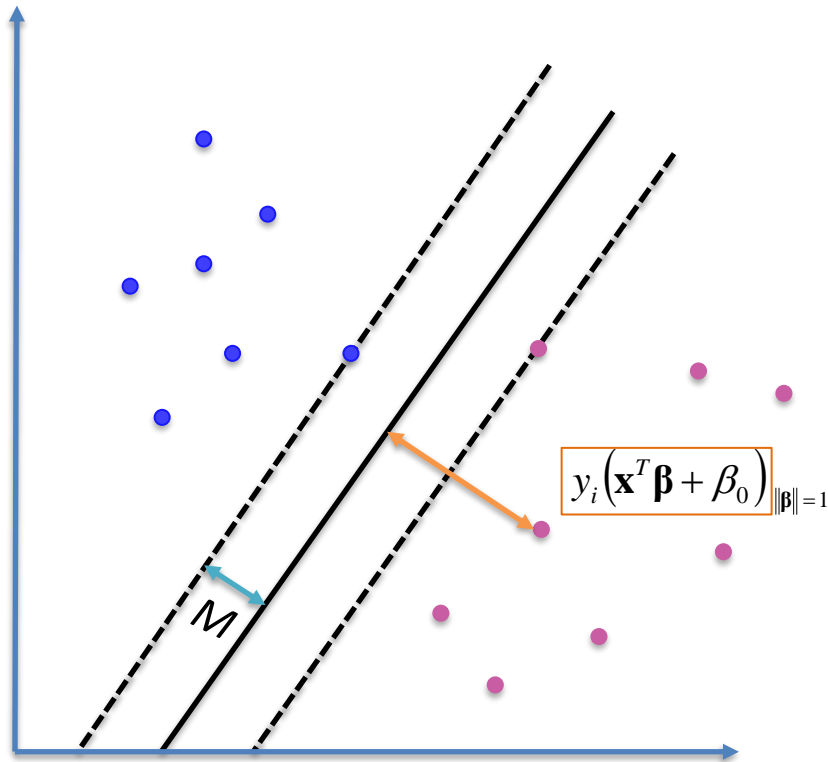
$$\text{sujeto a : } \begin{cases} \sum_{j=1}^n \beta_j^2 = 1 \\ y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip1}) \geq M \\ \forall i = 1, \dots, N \end{cases}$$

$$\max_{\beta_0, \beta, \|\beta\|=1} M$$
$$\text{sujeto a : } y_i (\mathbf{x}^T \boldsymbol{\beta} + \beta_0) \geq M$$



# El clasificador de margen máximo

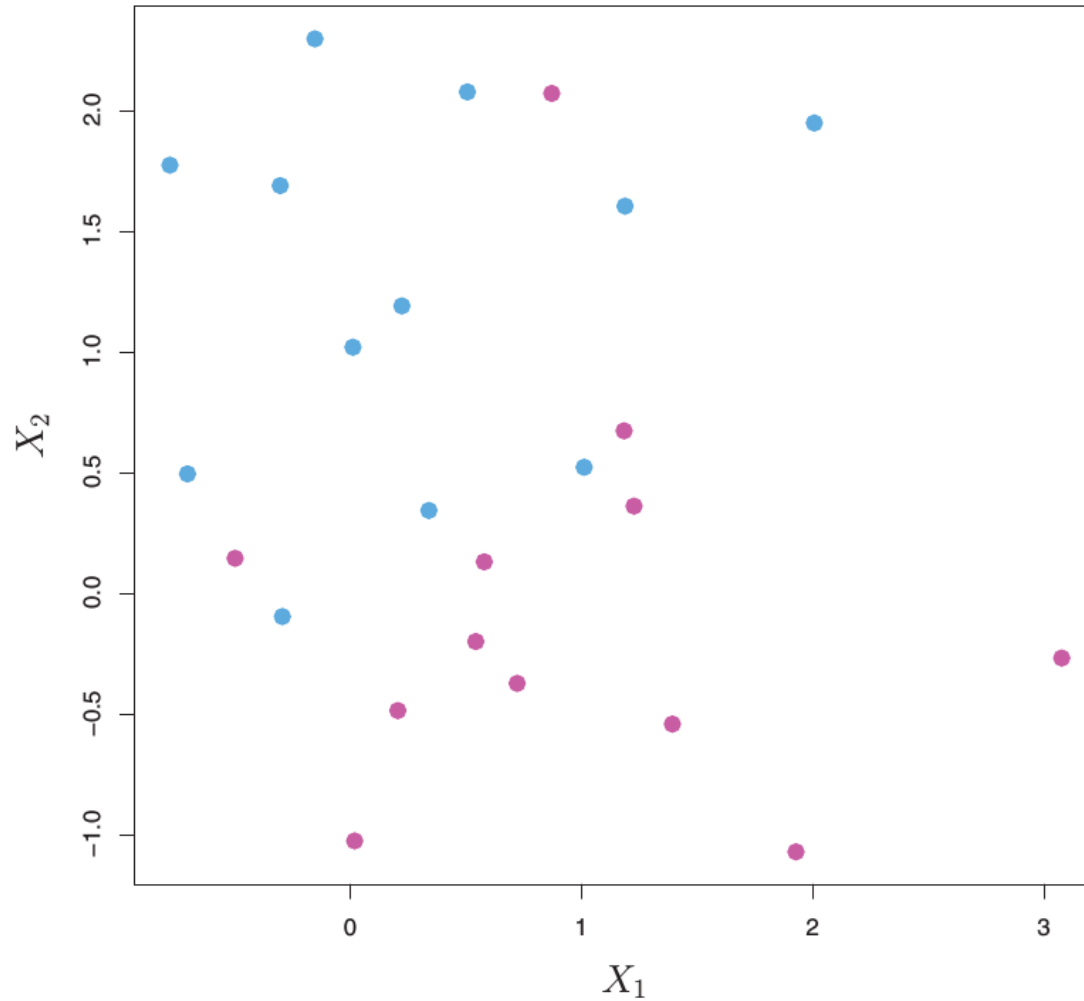
El **margen** es la “calle” máxima que se puede colocar entre dos observaciones de diferentes clases.



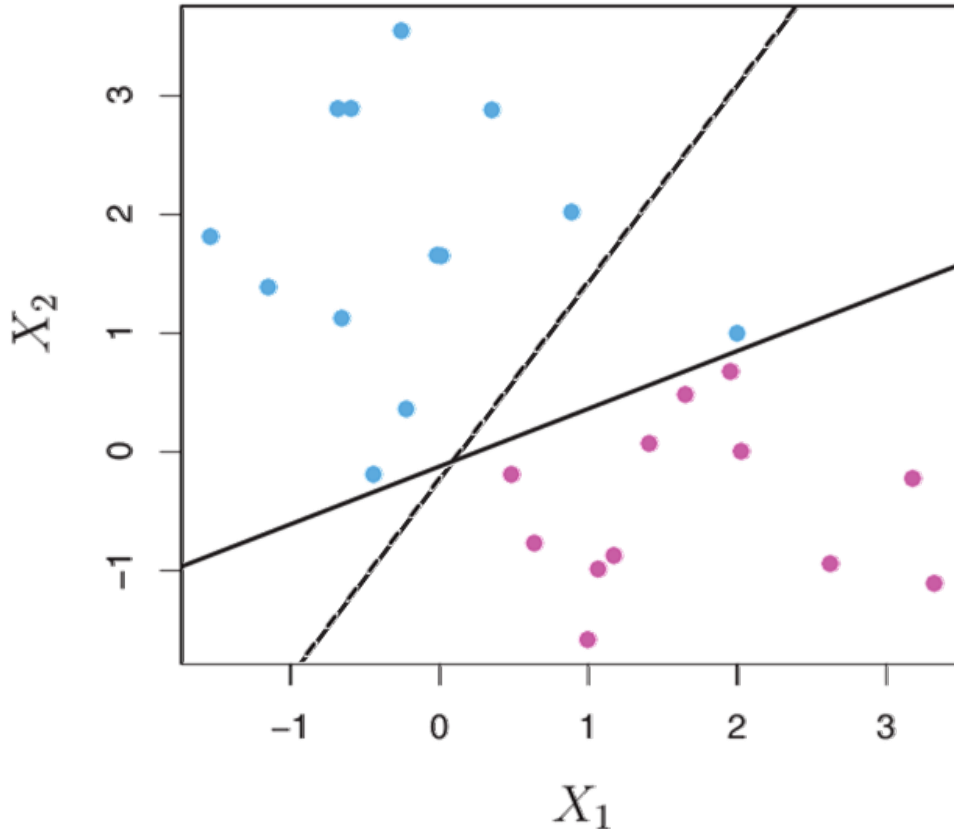
El hiperplano con máximo margen es aquel que soluciona:

$$\begin{aligned} & \max_{\beta_0, \boldsymbol{\beta}, \|\boldsymbol{\beta}\|=1} M \\ & \text{sujeto a : } y_i(\mathbf{x}^T \boldsymbol{\beta} + \beta_0) \geq M \end{aligned}$$

## Datos (linealmente) no separables

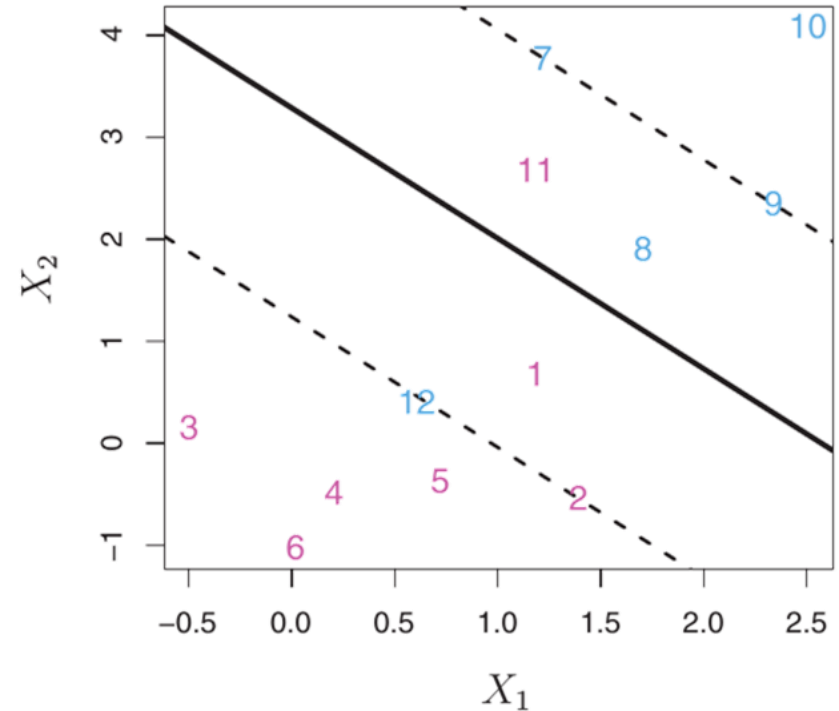
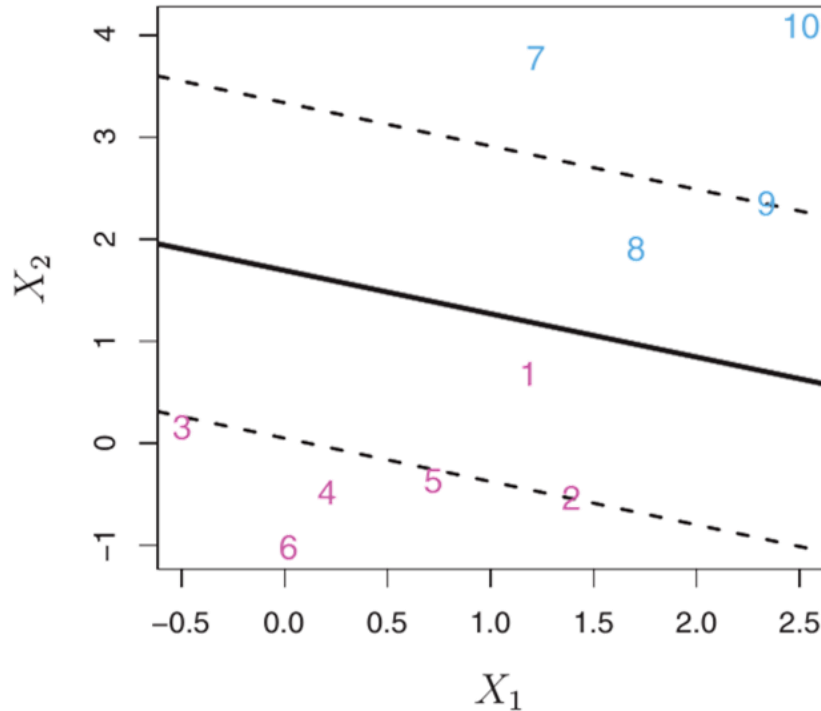


# Nuevos datos en el clasificador de margen máximo



- El clasificador de margen máximo es sensible a nuevos datos.
- El margen máximo es muy delgado, es decir, la clasificación es poco confiable.
- Es deseable un clasificador que no separe las dos clases perfectamente con el fin de:
  - ✓ Tener una gran robustez a datos individuales
  - ✓ Mejor clasificación de la *mayoría* (no la totalidad) de los datos del training.

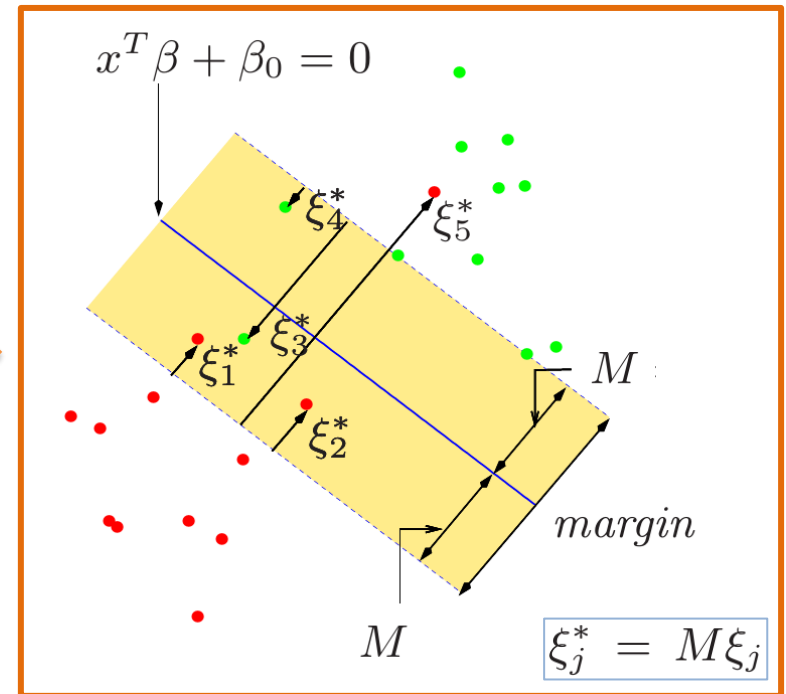
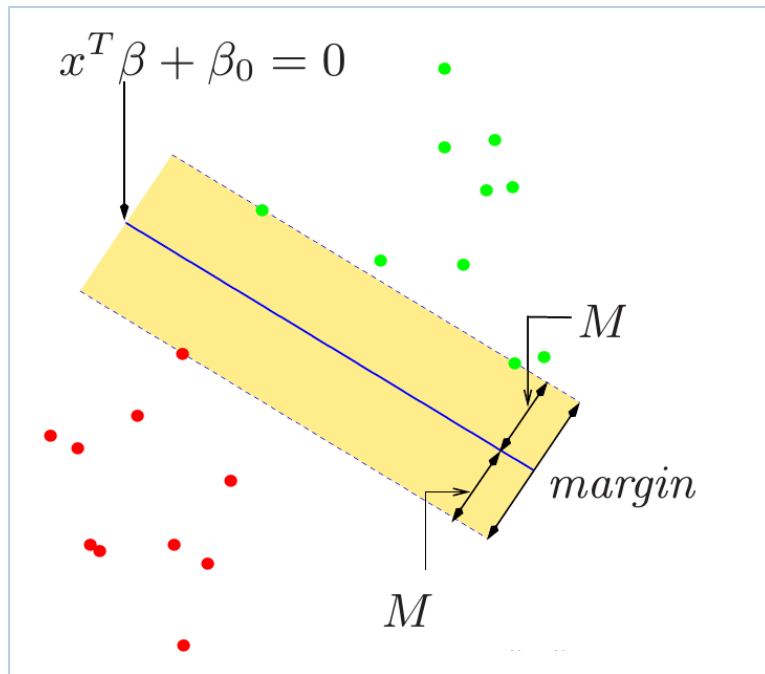
# Clasificador mediante vectores de soporte



- ✓ 3, 4, 5 y 6 están en el lado correcto del margen.
- ✓ 2 está en el margen y 1 está en el lado incorrecto del margen.
- ✓ 7 y 10 están en el lado correcto del margen.
- ✓ 9 está en el margen.
- ✓ 8 está en el lado incorrecto del margen.
- ✓ No hay observaciones en el lado incorrecto del hiperplano.

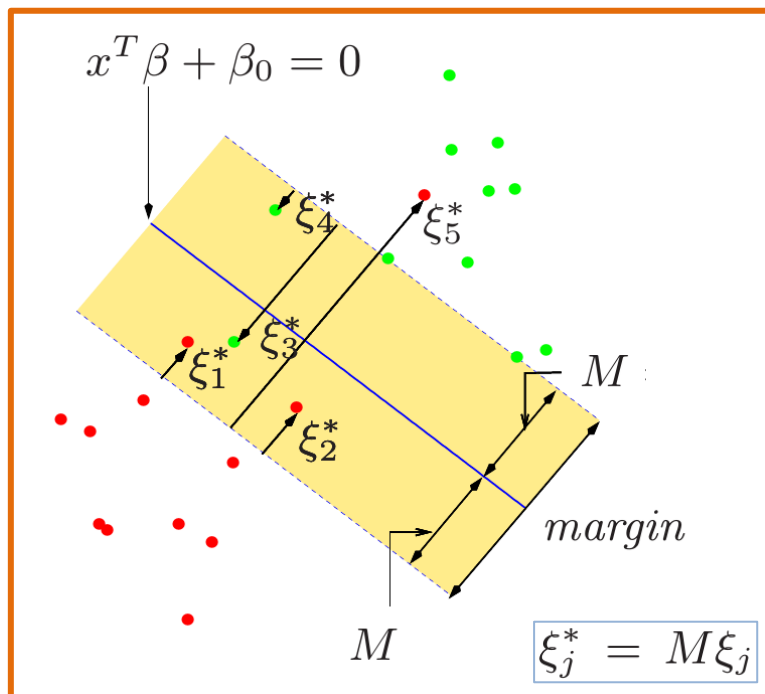
- ✓ Hay dos puntos adicionales, 11 y 12.
- ✓ Estos dos datos en el lado correcto del hiperplano y en el lado incorrecto del margen.

# Detalles del clasificador mediante vectores de soporte



$$\begin{aligned} & \max_{\beta_0, \beta, \xi, \|\beta\|=1} M \\ \text{sujeto a : } & \begin{cases} y_i (\mathbf{x}^T \beta + \beta_0) \geq M (1 - \xi_i) \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq \text{constante} = \frac{1}{C'} \end{cases} \end{aligned}$$

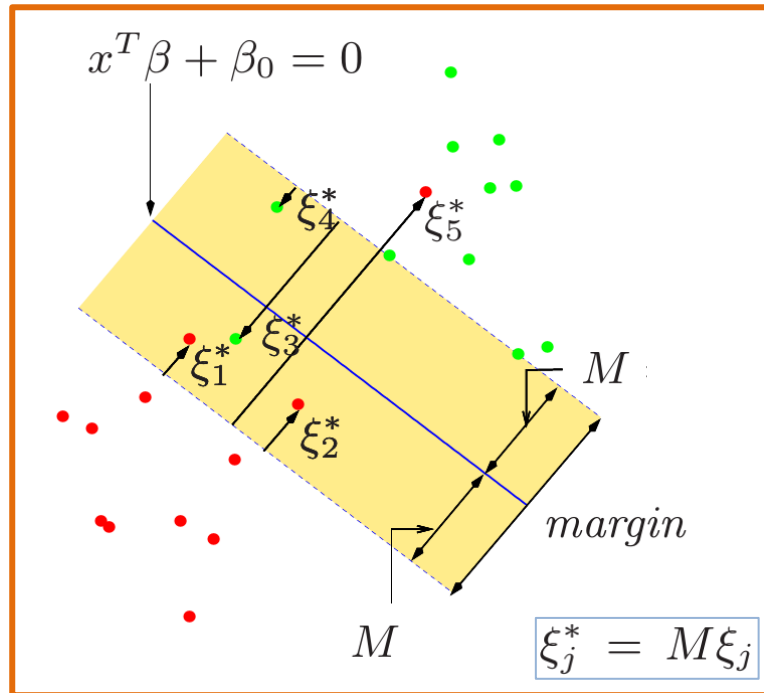
# Detalles del clasificador mediante vectores de soporte



- Si  $\xi_i = 0$  entonces el dato  $i$  está sobre el lado correcto del margen.
- Si  $\xi_i > 0$  entonces el dato  $i$  está sobre el lado incorrecto del margen.
- Si  $\xi_i > 1$  entonces el dato  $i$  está sobre el lado incorrecto del hiperplano.
- Si  $C' \rightarrow \infty$  (*presupuesto* = 0), no pueden haber violaciones al margen:  $\xi_1 = \dots = \xi_n = 0$
- Para el resto de casos, el número máximo de datos que pueden estar en el lado incorrecto del hiperplano está limitado por el *presupuesto* (máximo  $1/C'$  datos).
- Tanto como el *presupuesto* se incremente ( $C'$  disminuya), habrá más tolerancia respecto a las violaciones al margen (el margen se ensanchará).
- Tanto como el *presupuesto* disminuya ( $C'$  aumente), habrá menos tolerancia y el margen disminuirá.

$$\begin{aligned} & \max_{\beta_0, \beta, \xi, \|\beta\|=1} M \\ \text{sujeto a : } & \begin{cases} y_i(\mathbf{x}^T \beta + \beta_0) \geq M(1 - \xi_i) \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq \text{presupuesto} = \frac{1}{C'} \end{cases} \end{aligned}$$

# Detalles del clasificador mediante vectores de soporte



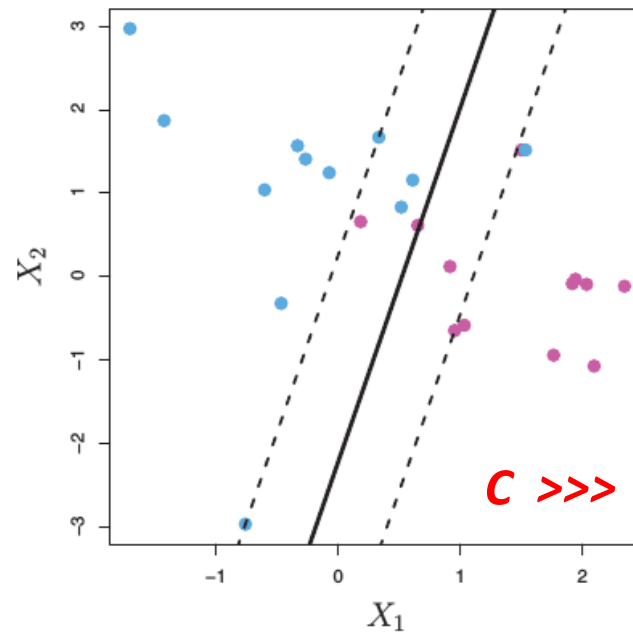
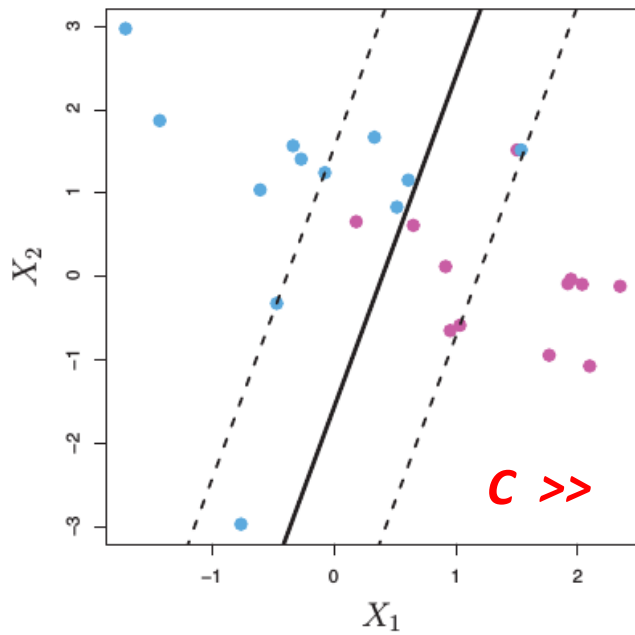
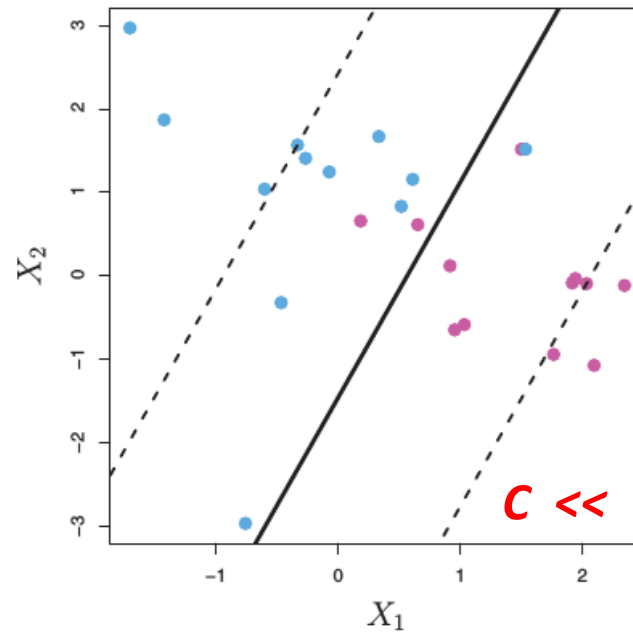
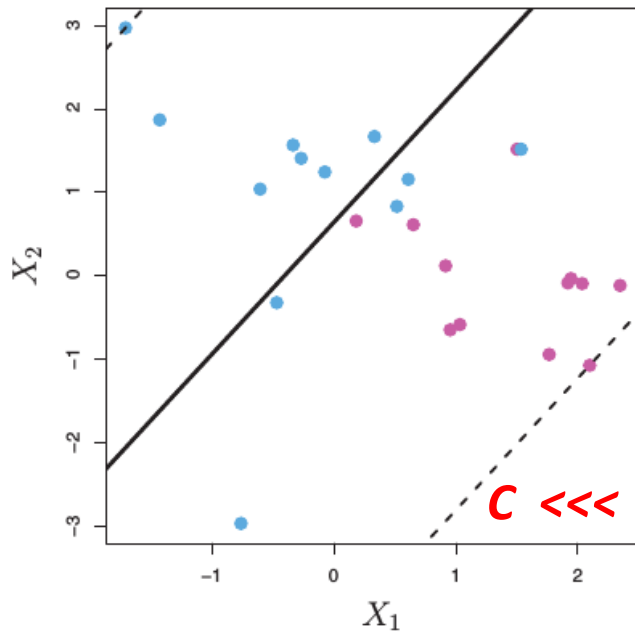
$$\begin{aligned} & \max_{\beta_0, \beta, \xi, \|\beta\|=1} M \\ \text{sujeto a : } & \begin{cases} y_i (\mathbf{x}^T \boldsymbol{\beta} + \beta_0) \geq M (1 - \xi_i) \\ \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq \frac{1}{C'} \end{cases} \end{aligned}$$

La solución a este problema de **optimización**



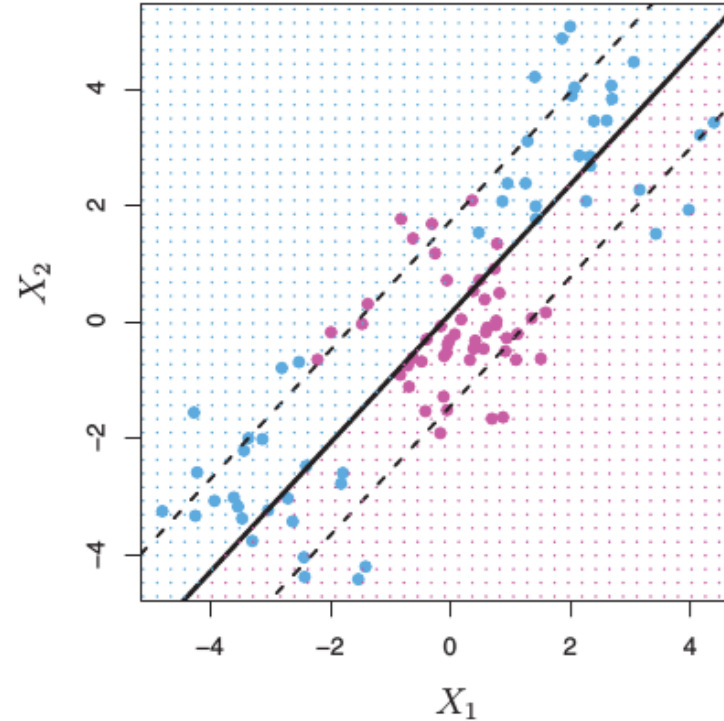
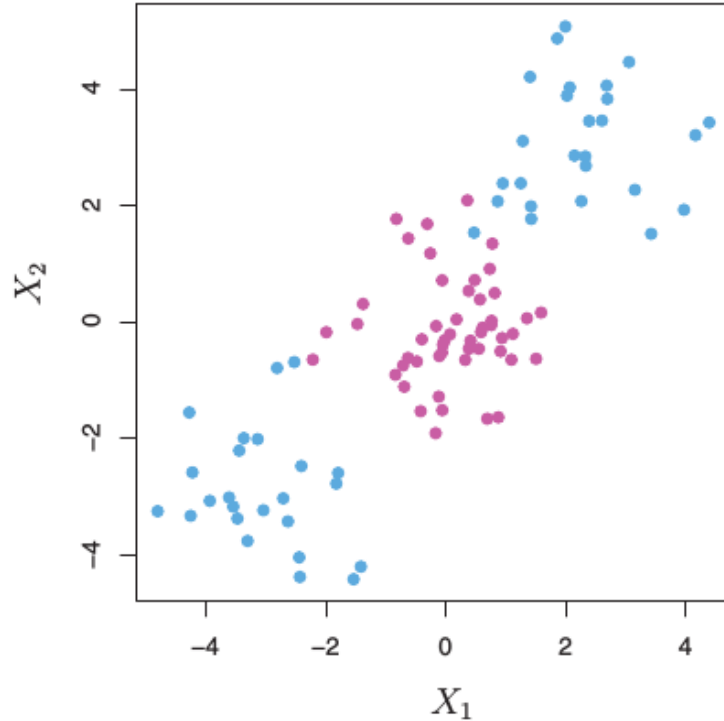
- ✓ Sólo los datos que se encuentran sobre el margen o que violan el margen afectan a la determinación del hiperplano. Es decir, los datos que se encuentran en el lado correcto del margen no afectan al clasificador.
- ✓ Los datos que se encuentran en el margen o en el lado incorrecto del margen son los **vectores de soporte**.

# C es un hiperparámetro





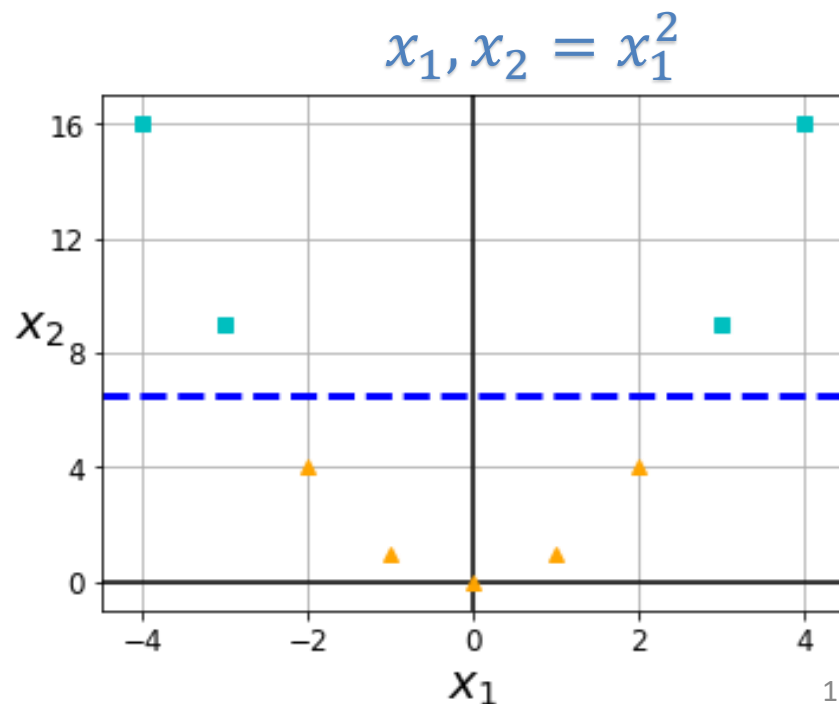
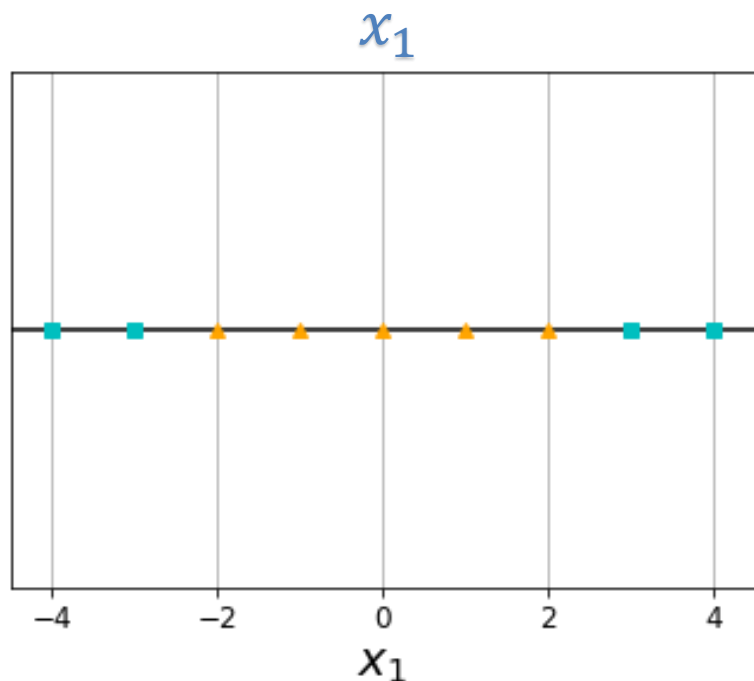
## ¿Qué sucede cuando hay fronteras de decisión no lineales?



- Algunas veces una frontera de decisión lineal simplemente no trabaja, independiente del valor de  $C$ .
- ¿Qué se puede hacer? -> Máquinas de vectores de soporte

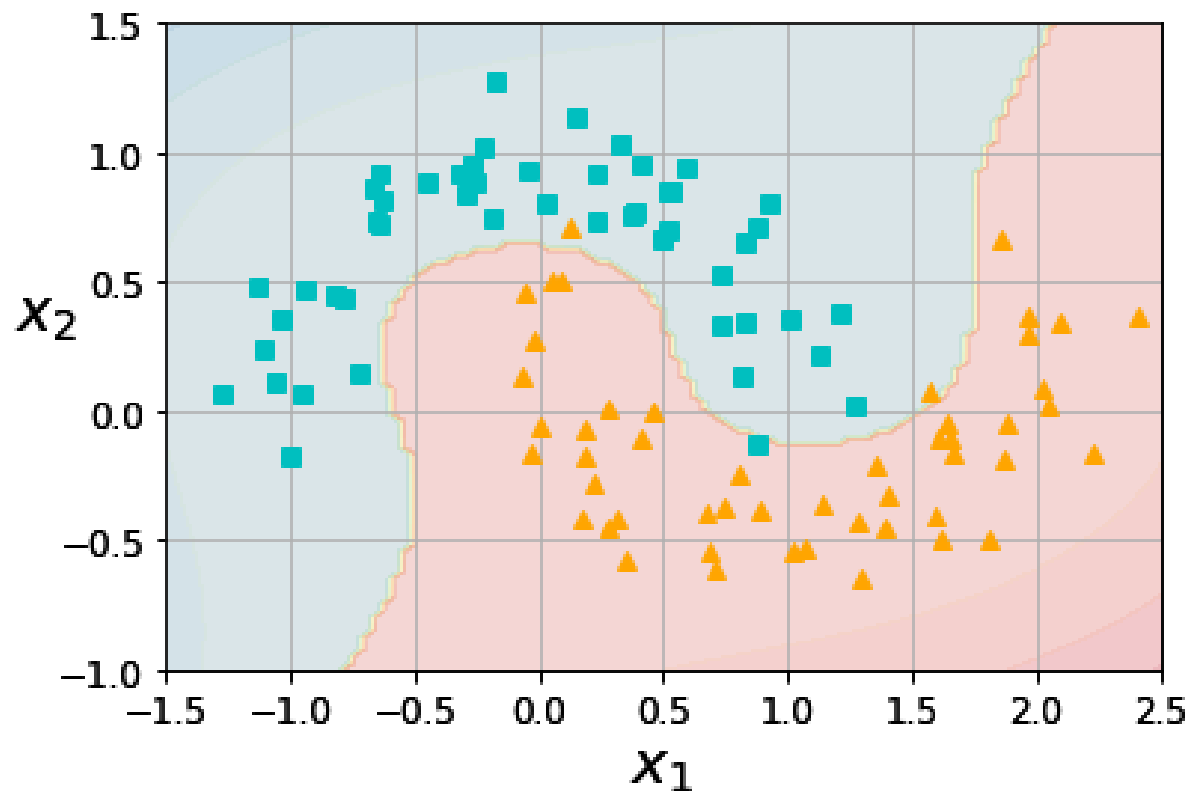
## Ampliación de los descriptores (feature expansion)

- Incrementar el espacio de descriptores incluyendo transformaciones, por ejm:  $X_1^2, X_1^3, X_1X_2, X_1X_2^2, \dots$ . Se va de un espacio  $p$ -dimensional a uno  $m$ -dimensional, donde  $m > p$ .
- Se ajusta el clasificador mediante vectores de soporte en el espacio ampliado.
- Esto resulta en una frontera de decisión no lineal en el espacio (de descriptores) original.



## Ampliación de los descriptores: polinomio cúbico

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2 + \beta_4 X_2^2 + \beta_5 X_1 X_2 + \beta_6 X_1^3 + \beta_7 X_2^3 + \beta_8 X_1 X_2^2 + \beta_9 X_1^2 X_2 = 0$$




- Se ha usado una ampliación cúbica polinomial
- Se incrementa de 2 a 9 variables
- El clasificador mediante vectores de soporte en el espacio ampliado soluciona el problema en el espacio original (de baja dimensión).

# Producto interno (escalar) y vectores de soporte

Producto interno:  $\langle \mathbf{x}_i | \mathbf{x}_{i'} \rangle = \sum_{j=1}^p x_{ij} x_{i'j}$

- El clasificador mediante vectores de soporte puede representarse como:

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i \langle \mathbf{x} | \mathbf{x}_i \rangle$$


- Para estimar los  $n$  parámetros  $\alpha_1, \alpha_2, \dots, \alpha_n$  y  $\beta_0$ , se requieren  $\binom{n}{2}$  productos internos  $\langle \mathbf{x}_i | \mathbf{x}_{i'} \rangle$  entre todos los pares de datos de entrenamiento.
- Resulta que la mayoría de  $\alpha_i$  pueden ser cero (los que no están relacionados a los vectores de soporte):

$$f(\mathbf{x}) = \beta_0 + \sum_{i \in S} \alpha_i \langle \mathbf{x} | \mathbf{x}_i \rangle$$

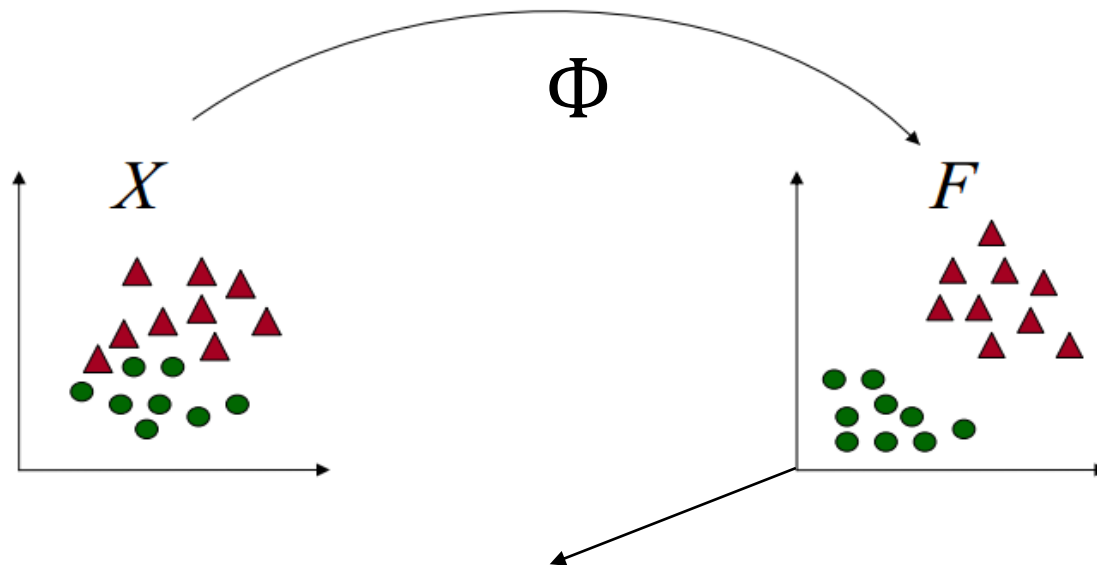
Donde  $S$  es el conjunto (de vectores) de soporte de los índices  $i$  tal que  $\alpha_i > 0$ .

## El truco del kernel (kernel trick)

- Se mapean los datos a un espacio de descriptores ampliado vía una función de mapeo  $\Phi(x)$ .
- El mapeo puede evaluarse a través de una función kernel (teorema de Mercer):

$$K(x_i, x_{i'}) = \langle \Phi(x_i) | \Phi(x_{i'}) \rangle$$

- Se construye una función lineal en el nuevo espacio ampliado de descriptores.



## Máquinas de vectores de soporte (SVM)

- Si se pueden calcular los productos internos entre los datos, se puede ajustar un clasificador SVM.
- Algunas funciones kernel especiales pueden hacer esto directamente, por ejm:

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = (\gamma \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle + r)^d$$

calcula el productor interno necesario para un polinomio de dimensión  $d$ , que normalmente necesitaría calcular un número  $\binom{p+d}{d}$  de funciones base.

- La solución ahora tiene la forma:

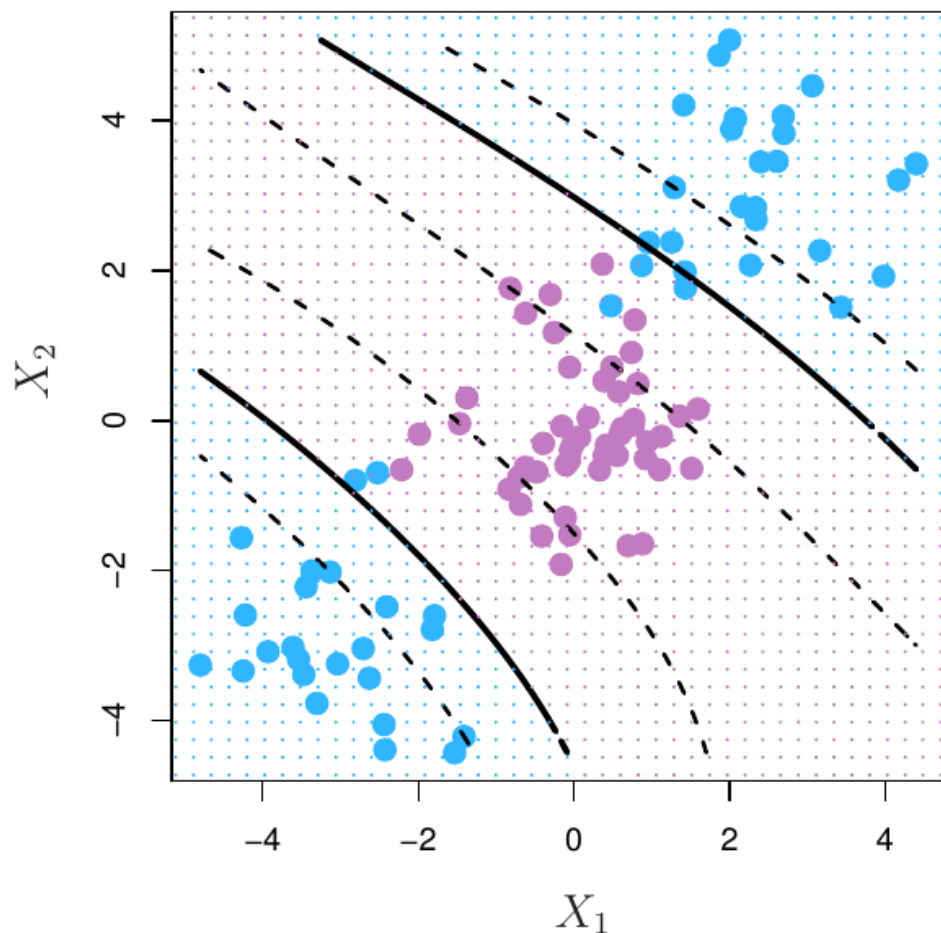
$$f(x) = \beta_0 + \sum_{i \in S}^n \hat{\alpha}_i K(\mathbf{x}_i, \mathbf{x}_{i'})$$

## Máquinas de vectores de soporte (SVM)

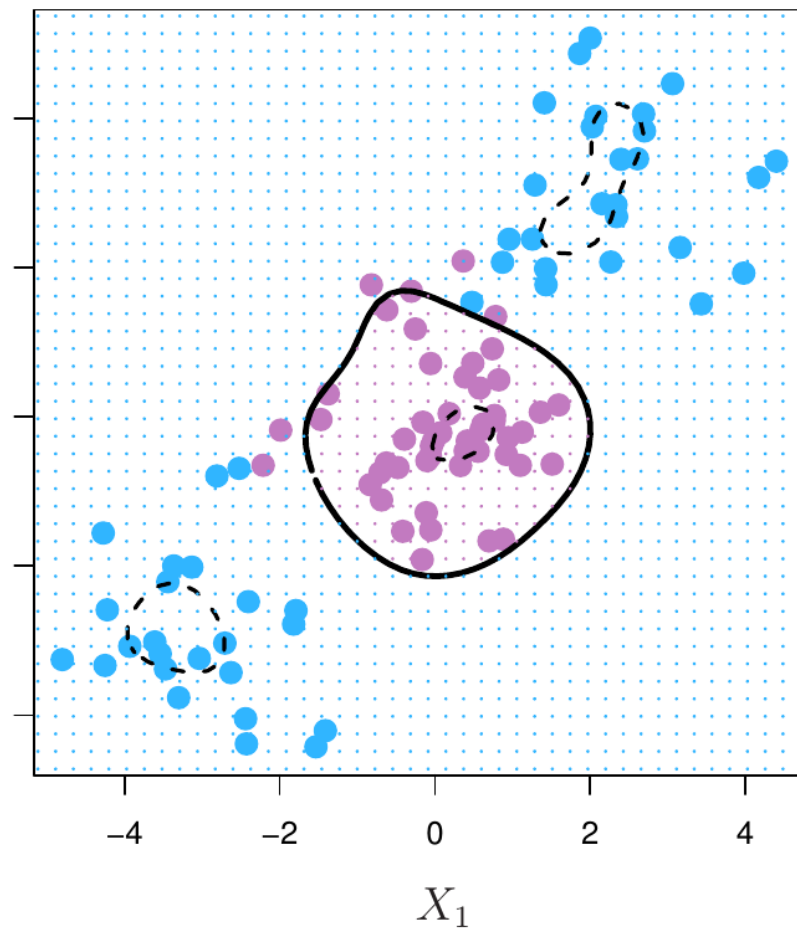
Algunos tipos usuales de funciones kernel son:

- Lineal:  $K(\mathbf{x}_i, \mathbf{x}_{i'}) = \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle$
- Polinomial:  $K(\mathbf{x}_i, \mathbf{x}_{i'}) = (\gamma \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle + r)^d$
- Gaussiano (RBF):  
$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_{i'}\|^2)$$
- Sigmoide:  
$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \tanh(\gamma \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle + r)$$

# Máquinas de vectores de soporte (SVM)



Kernel polinomial de grado 3

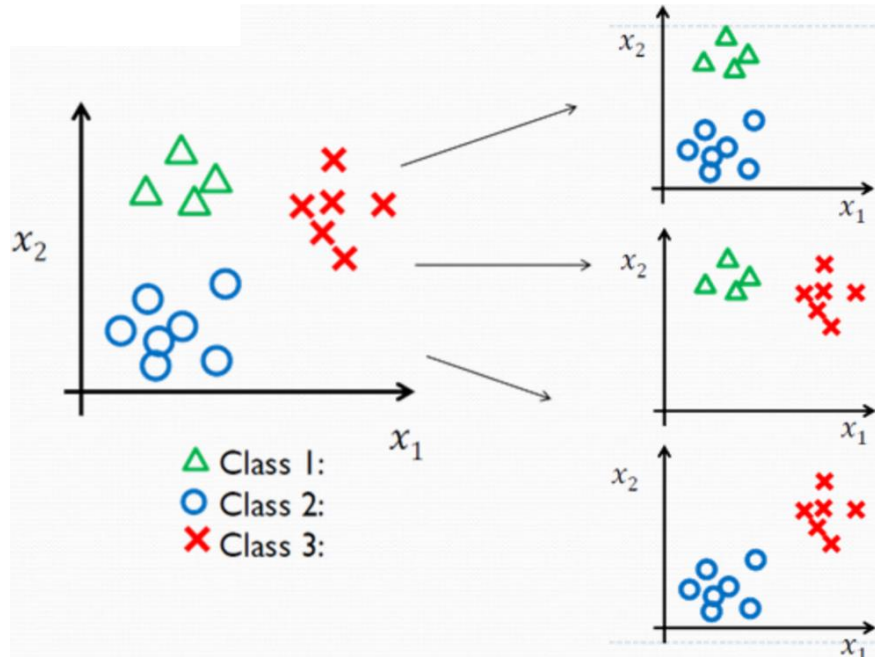


Kernel Gaussiano



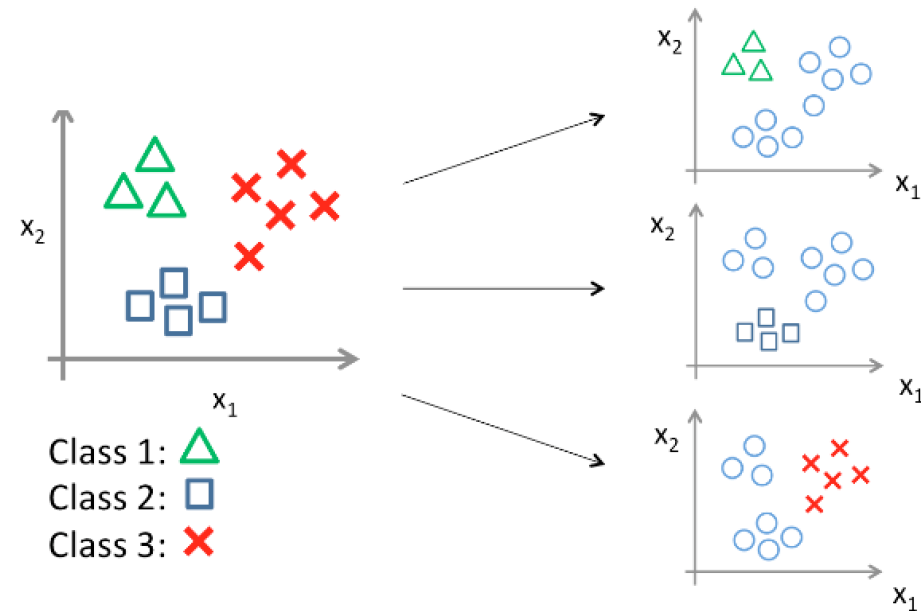
# SVM con más de dos clases

## One vs one



- Se escoge para cada observación la clase que gane la mayoría de clasificaciones
- $\binom{K}{2}$  clasificadores

## One vs All (one vs rest)



- Se escoge para cada observación la clase con la **mayor** probabilidad.
- Un clasificador por clase.

## Implementación con scikit-learn

Class	Time complexity	Out-of-core support	Scaling required	Kernel trick
LinearSVC	$O(m \times n)$	No	Yes	No
SGDClassifier	$O(m \times n)$	Yes	Yes	No
SVC	$O(m^2 \times n)$ to $O(m^3 \times n)$	No	Yes	Yes

# Implementación con scikit-learn

```
from sklearn.svm import LinearSVC
```

Librería: **liblinear**

Muy rápida, pero solo kernel lineal

```
1 iris = datasets.load_iris()
2 X = iris["data"]
3 y = (iris["target"])
4 svm_clf = Pipeline((
5     ("scaler", StandardScaler()),
6     ("linear_svc", LinearSVC(C=1)),
7 ))
8
9 svm_clf.fit(X, y)
```

# Implementación con scikit-learn

Librería: **libsvm**

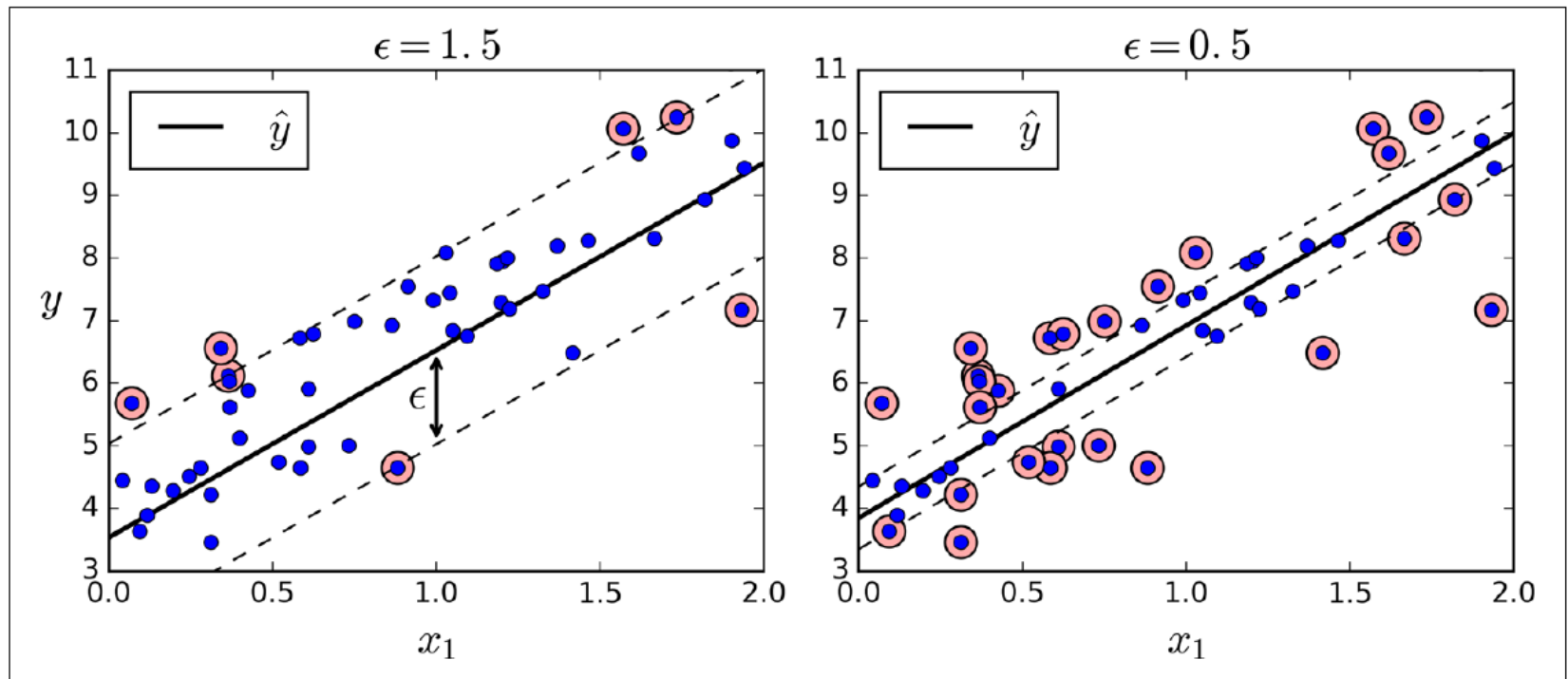
No tan rápida

Diferentes kernels

```
1  from sklearn.svm import SVC
2
3  iris = datasets.load_iris()
4  X = iris["data"]
5  y = (iris["target"])
6
7  svm_clf = Pipeline((
8      ("scaler", StandardScaler()),
9      ("RBF_svc", SVC(C=10, kernel='rbf', gamma=0.5)),
10 ))
11
12 svm_clf.fit(X, y)
```

## Support vector regression (SVR)

- El método de SVM puede ampliarse para resolver problemas de regresión. Este método se llama Regresión mediante vectores de soporte (SVR).



# Implementación con scikit-learn

Librería: **libsvm**

No tan rápida

Diferentes kernels

```
1 from sklearn.svm import SVR
2 X = [[0, 0], [2, 2]]
3 y = [0.5, 2.5]
4 clf = SVR()
5 clf.fit(X, y)
6
7
8 # clf.predict([[1, 1]])
9 clf
```

```
SVR(C=1.0, cache_size=200, coef0=0.0, degree=3, epsilon=0.1, gamma='auto',
    kernel='rbf', max_iter=-1, shrinking=True, tol=0.001, verbose=False)
```