



APRENDIZAJE DE MÁQUINA

MODELOS NO LINEALES I

AGENDA

01 Máquina de vectores de soporte

FALTA

02 K Vecinos

FALTA

03 Árboles de decisión

FALTA

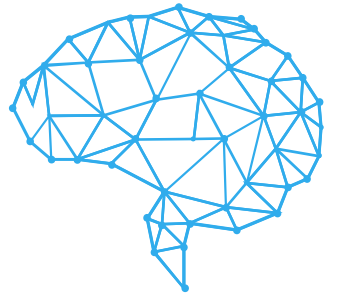




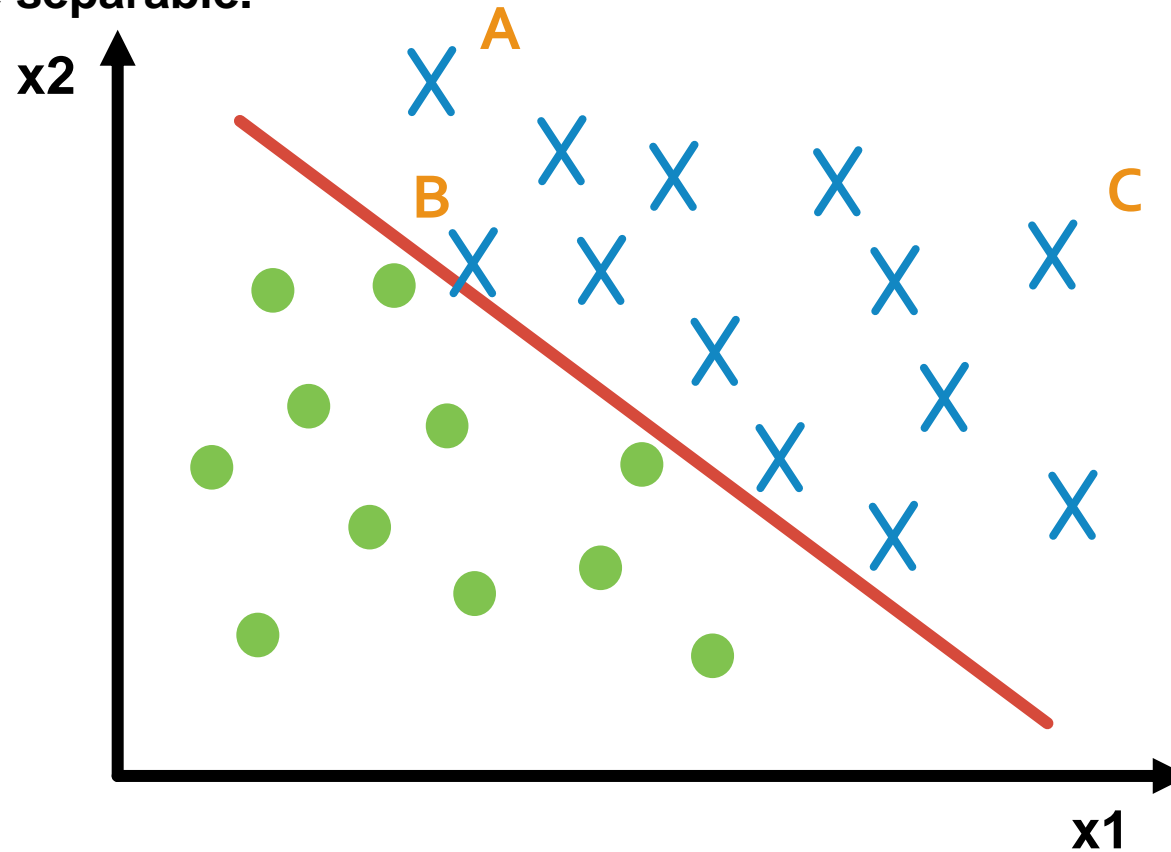
**MÁQUINA DE
VECTORES DE
SOPORTE**

M Á Q U I N A D E S V

I N T R O D U C C I Ó N



Se asume que tenemos un conjunto de datos de entrenamiento que es linealmente separable.



M Á Q U I N A D E S V

C A M B I O D E N O T A C I Ó N



Para **proceder**, se va a **realizar** un **cambio** de **notación**:

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

$$h_{w,b} \in \{-1, 1\}$$

Por lo tanto se tiene que:

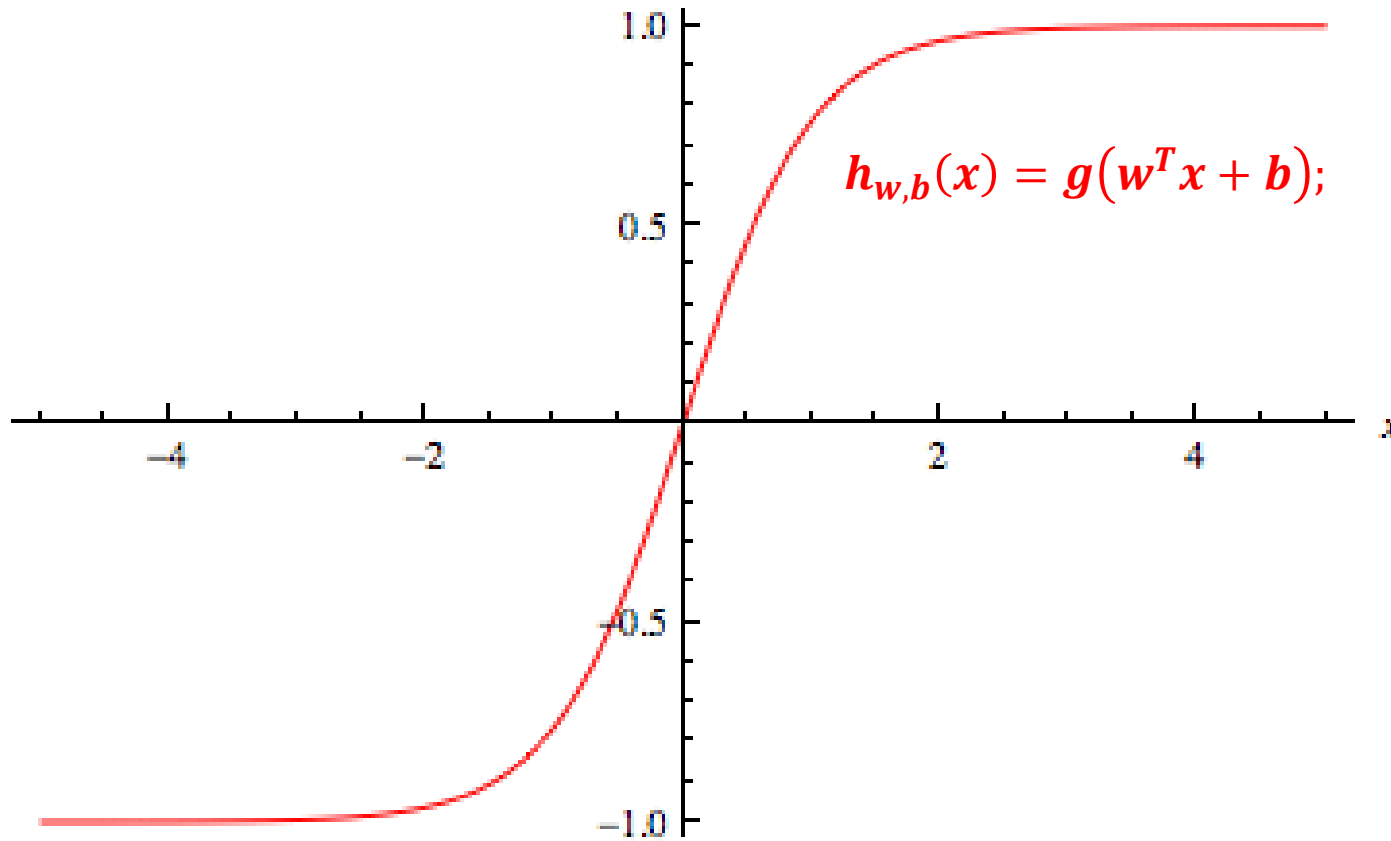
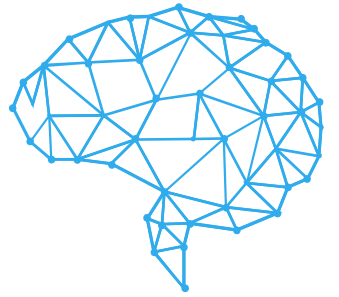
$$g(z) \begin{cases} 1 & \text{sí } z \geq 0 \\ -1 & \text{de otra forma} \end{cases}$$

$$h_{w,b}(x) = g(w^T x + b); w \in \mathbb{R}^n, x \in \mathbb{R}^n$$

(Se manipula el sesgo de manera independiente)

M Á Q U I N A D E S V

INTUICIÓN DE MARGEN



Se tienen **dos** cuestiones:

$$w^T x \rightarrow \infty$$
$$g(z) \rightarrow 1$$

$$w^T x \rightarrow -\infty$$
$$g(z) \rightarrow -1$$

M Á Q U I N A D E S V

M A R G E N F U N C I O N A L



Se define el **margen funcional respecto** a un **dato de entrenamiento** $(x^{(i)}, y^{(i)})$ **como:**

$$\widehat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

Un **margen funcional grande** representa una **predicción correcta y certera**.

Variable de respuesta $y^{(i)}$	Margen funcional $\widehat{\gamma}^{(i)}$
1	$w^T x^{(i)} \ll 0$ $\widehat{\gamma}^{(i)} \ll 0$
	$w^T x^{(i)} \gg 0$ $\widehat{\gamma}^{(i)} \gg 0$
-1	$w^T x^{(i)} \ll 0$ $\widehat{\gamma}^{(i)} \gg 0$
	$w^T \gg 0$ $\widehat{\gamma}^{(i)} \ll 0$

M Á Q U I N A D E S V

M A R G E N F U N C I O N A L



El problema del **margen funcional**, es que **solo depende** del **signo** de $y^{(i)}$ **y no** de la **magnitud** de $w^T x^{(i)} + b$

Es decir, puedo “**forzar arbitrariamente**” que la predicción sea la correcta tan solo escalando w .

$$\gamma^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

$\gamma^{(i)} \gg 0$ **sí y solo sí:**

$$w = kw$$

$$b = kb$$

$$k \gg 0$$

Una **solución** es **normalizar** el **vector** w (se verá más adelante):

$$\frac{w}{\|w\|_2}$$

M Á Q U I N A D E S V

M A R G E N F U N C I O N A L



En este caso, nos **centraremos** en el **dato** de **entrenamiento**, que **proporciona** el **peor margen funcional** de **todos**, es decir:

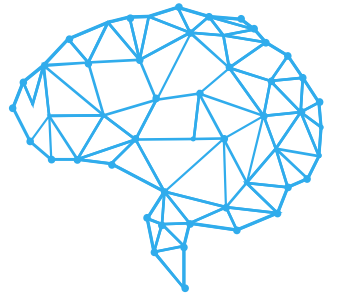
$$\hat{\gamma} = \min_{i=1,2,\dots,m} \gamma^{(i)}$$

En consecuencia el **margen funcional** se **define** como:

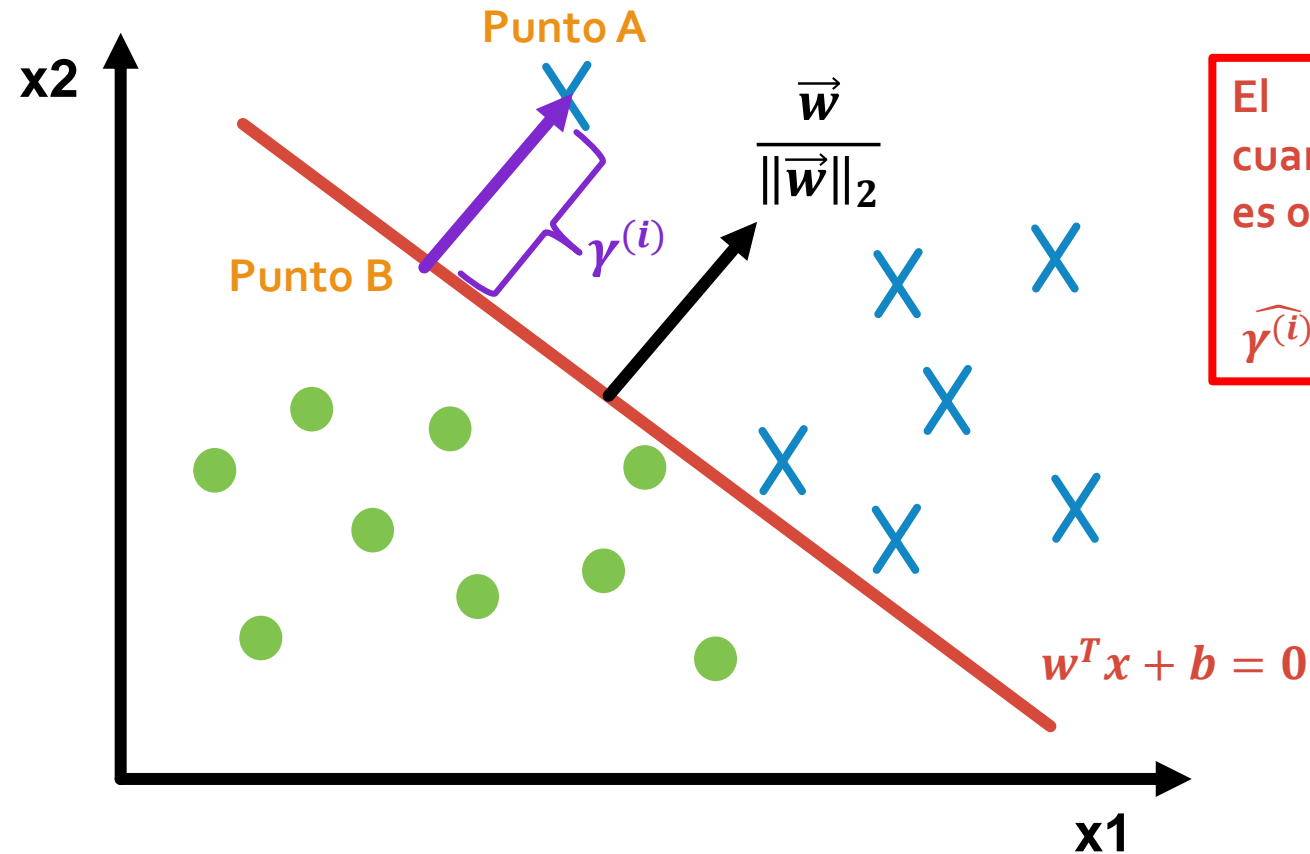
El mínimo de todos los márgenes funcionales contemplando todos los datos de entrenamiento $(x^{(i)}, y^{(i)})$

M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O



Para definir el **margen geométrico** $\gamma^{(i)}$ se **visualiza lo siguiente**:

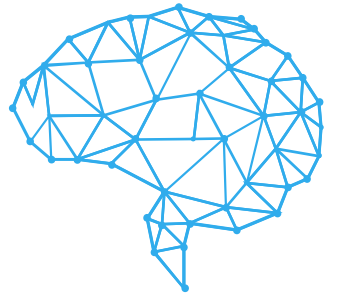


El hiperplano se forma cuando el margen funcional es 0

$$\widehat{\gamma^{(i)}} = y^{(i)}(w^T x^{(i)} + b) = 0$$

M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O



El vector w es **perpendicular** al **hiperplano** formado por $w^T x + b = 0$ por la **siguiente razón**:

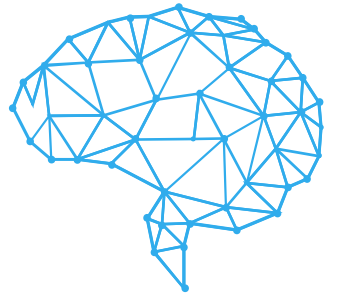
Sí $b = 0$ entonces **se tiene que**:

$$\begin{aligned}w^T x + b &= 0 \\w^T x &= 0\end{aligned}$$

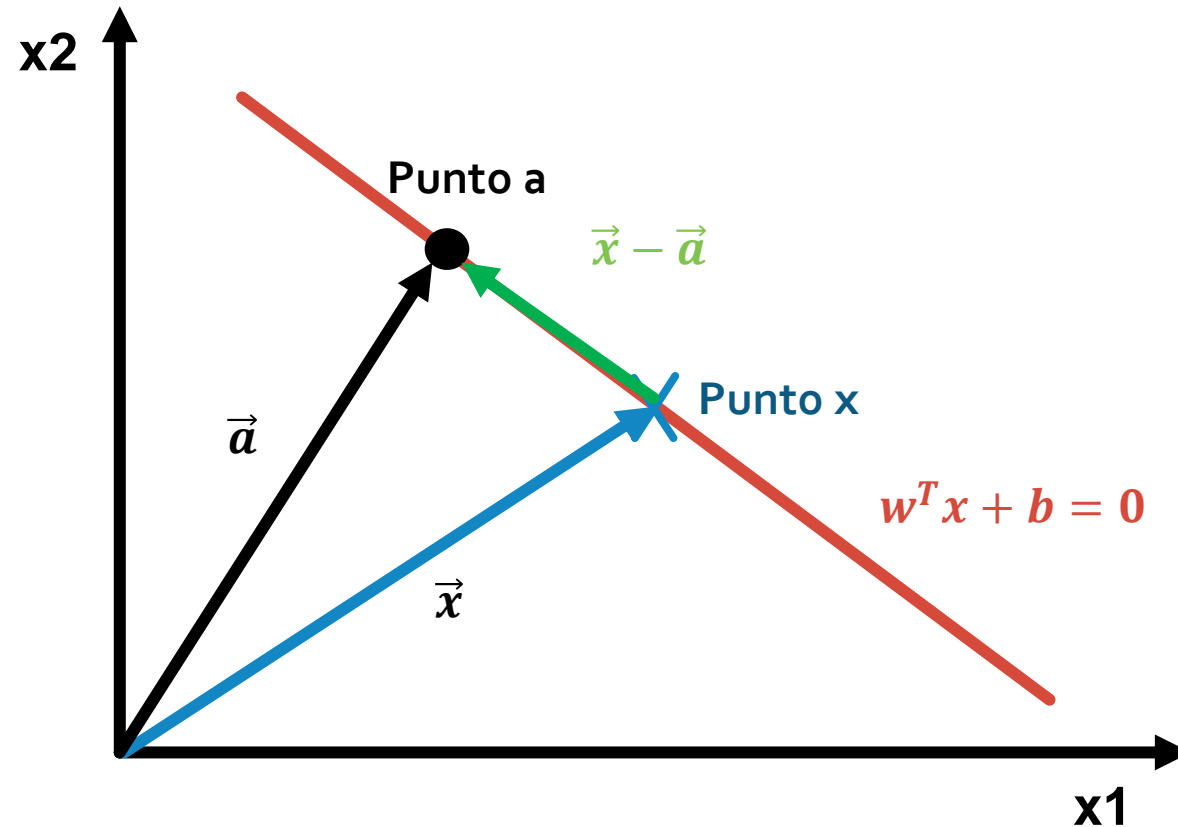
Por lo tanto para que $w^T x = 0$ **sea posible**, **ambos vectores** deben ser **perpendiculares**

M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O

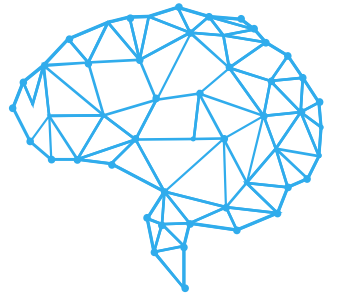


Para ver el caso de $b \neq 0$, se toma un punto a en el plano, definido por el vector \vec{a} :



M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O



Por lo tanto, la **ecuación** del **plano** **también** se puede **expresar** como:

$$w^T x + b = 0$$

$$w^T (\vec{x} - \vec{a}) = 0$$

Desarrollando se tiene:

$$w^T \vec{x} - w^T \vec{a} = 0$$

La **única manera** de que **esto** se **cumpla** es:

$$w, \vec{x} \text{ sean } \perp$$

$$w, \vec{a} \text{ sean } \perp$$

M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O

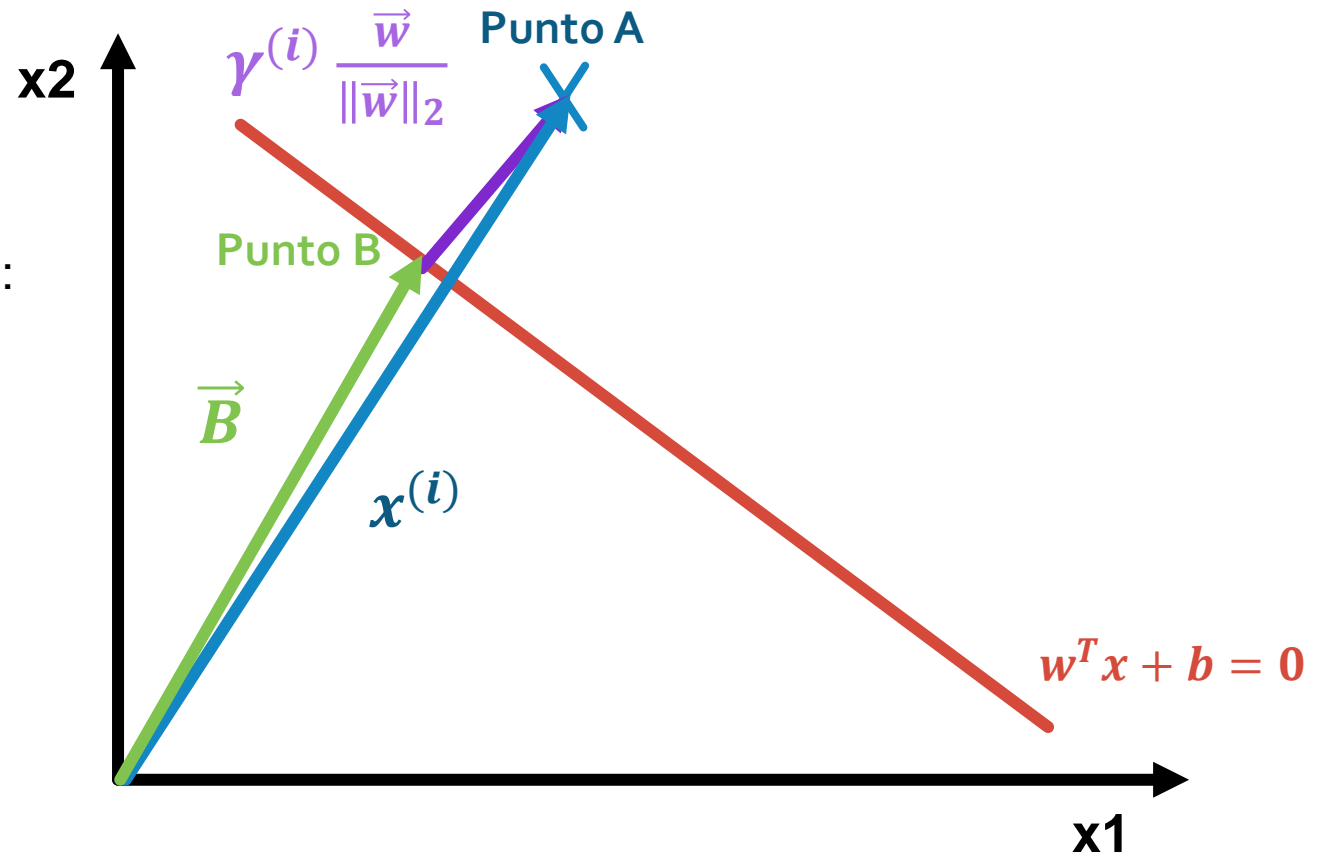


Una vez demostrado que \vec{w}, \vec{x} son \perp , se **define** lo siguiente para **obtener el margen geométrico**:

$$\vec{B} + \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|_2} = \vec{x}^{(i)}$$

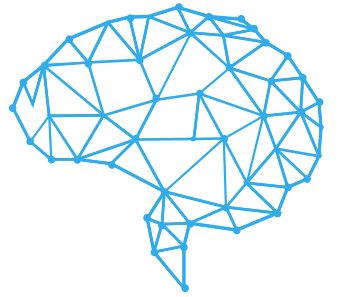
Se obtiene el **punto B** en el **hiperplano**:

$$\vec{B} = \vec{x}^{(i)} - \gamma^{(i)} \frac{\vec{w}}{\|\vec{w}\|_2}$$



M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O



Como el punto \vec{B} se encuentra en el hiperplano, se debe satisfacer la ecuación del hiperplano:

$$\mathbf{w}^T \vec{B} + b = 0$$

$$\mathbf{w}^T \left(\mathbf{x}^{(i)} - \gamma^{(i)} \frac{\vec{\mathbf{w}}}{\|\vec{\mathbf{w}}\|_2} \right) + b = 0$$

Se calcula el margen geométrico:

$$\gamma^{(i)} = \frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\vec{\mathbf{w}}\|_2}$$

M Á Q U I N A D E S V

M A R G E N G E O M É T R I C O



Para tener en **cuenta** el **signo**, de manera general, el **margen geométrico** estaría **dado por**:

$$\gamma^{(i)} = y^{(i)} \left[\frac{\mathbf{w}^T \mathbf{x}^{(i)} + b}{\|\vec{\mathbf{w}}\|_2} \right]$$

Nos damos cuenta que el **margen geométrico** es **invariante** a la **escala** de w por la **normalización**.

De igual forma, **contemplando todos** los **datos** de **entrenamiento**, se define el **peor escenario**:

$$\gamma = \min_{i=1,2,\dots,m} \gamma^{(i)}$$

M Á Q U I N A D E S V

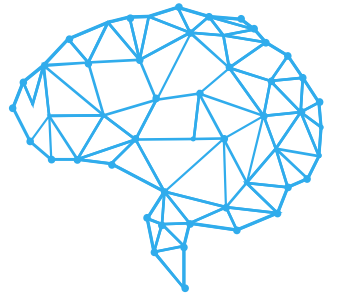
MARGEN GEOMÉTRICO Y FUNCIONAL



Relacionando **ambos márgenes**, **geométrico** y **funcional**, se tiene que:

$$\gamma^{(i)} = \frac{\widehat{\gamma^{(i)}}}{\|\vec{w}\|_2}$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Una vez **definidos ambos márgenes**, se puede expresar el **objetivo de clasificación** como un **problema de maximización del margen geométrico**.

$$\max_{\gamma, w, b} \gamma$$

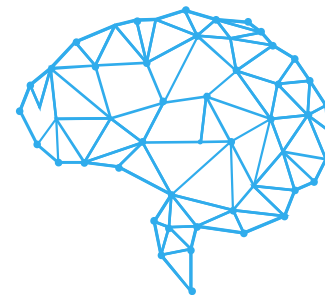
Sujeto a:

$$y^{(i)}(w^T x^{(i)} + b) \geq \gamma; i = 1, \dots, m$$

$$\|w\| = 1$$

Con $\|w\| = 1$ garantizamos que el **margen geométrico** es igual al **margen funcional**.

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Debido a que la **restricción** $\|w\| = 1$ define un **problema de optimización NO convexo** (**encontrar pesos w en una esfera**) se transforma el problema para **eliminar la restricción**:

$$\max_{\gamma, w, b} \frac{\hat{\gamma}}{\|w\|}$$

Sujeto a:

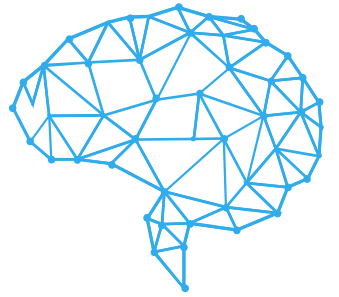
$$y^{(i)}(w^T x^{(i)} + b) \geq \hat{\gamma}; i = 1, \dots, m$$

De todas formas se **sigue** teniendo un **problema NO convexo**, pero ahora éste se **encuentra en la función objetivo**.

$$\frac{\hat{\gamma}}{\|w\|}$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Si ahora **imponemos** la **restricción de escalamiento**, que determina que el **margen funcional** debe ser $\hat{\gamma} = 1$.

$$\max_{\gamma, w, b} \frac{1}{\|w\|}$$

Sujeto a:

$$y^{(i)}(w^T x^{(i)} + b) \geq \hat{1}; i = 1, \dots, m$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Nos damos cuenta que la **maximización** es lo **mismo** que un **problema de minimización**:

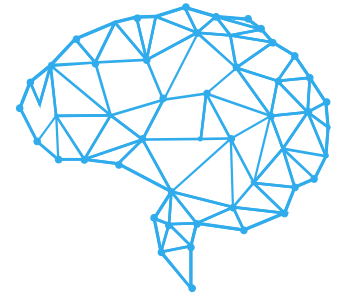
$$\max_{\gamma, w, b} \frac{1}{\|w\|} = \min_{\gamma, w, b} \frac{1}{2} \|w\|^2$$

Sujeto a:

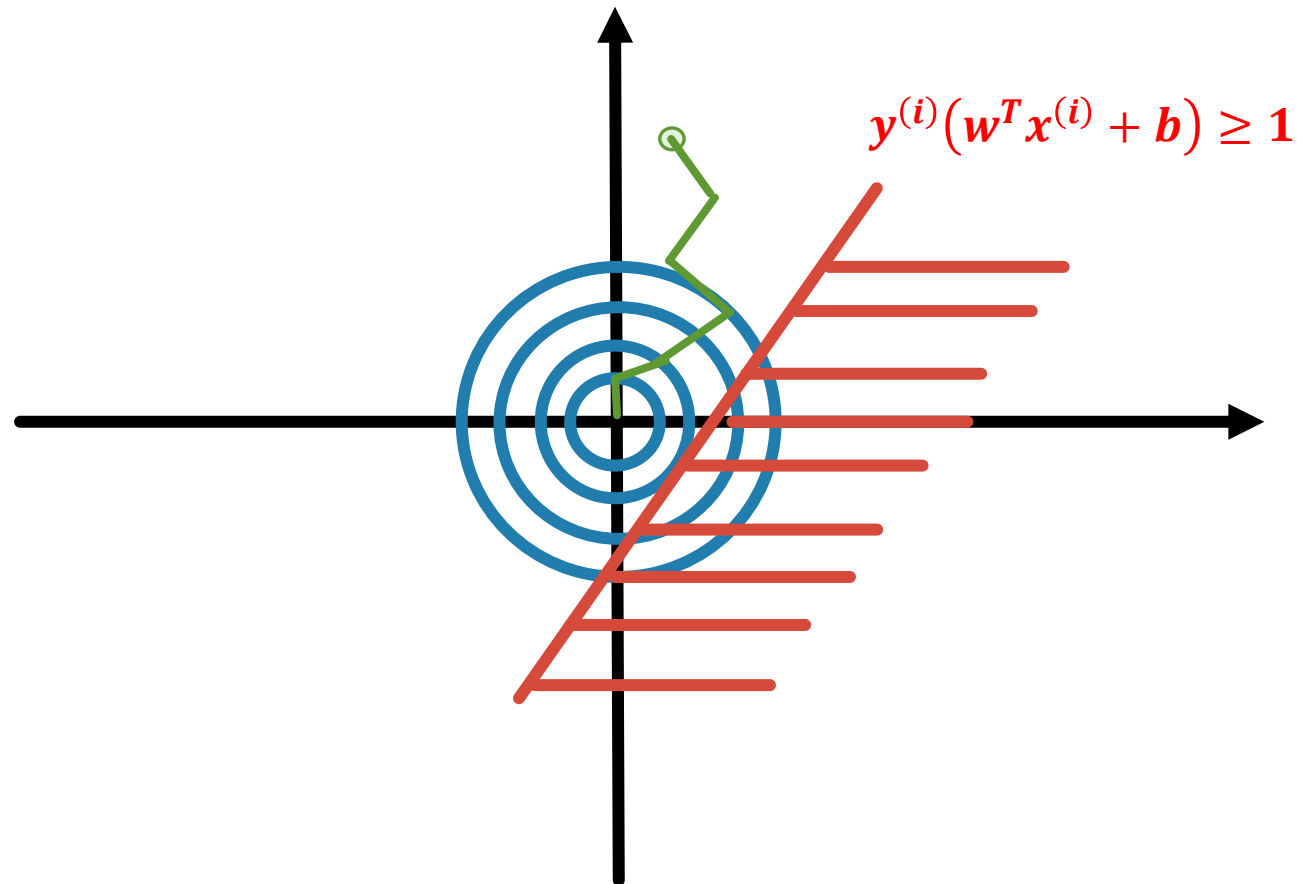
$$y^{(i)}(w^T x^{(i)} + b) \geq 1; i = 1, \dots, m$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Ejemplo gráfico:



M Á Q U I N A D E S V



MULTIPLICADORES DE LAGRANGE

Ahora veremos como **solucionar** problemas de **optimización** sujetos a **restricciones**. Consideremos el **siguiente** problema de **optimización**:

$$\min_w f(w)$$

Sujeto a:

$$h_i(w) = 0; i = 1, \dots, l$$

Este **problema** se **puede** **solucionar** con **multiplicadores** de **Lagrange**. El **lagrangiano** estaría dado por:

$$L(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

donde β_i son los **multiplicadores** de **Lagrange**.

M Á Q U I N A D E S V

MULTIPLICADORES DE LAGRANGE



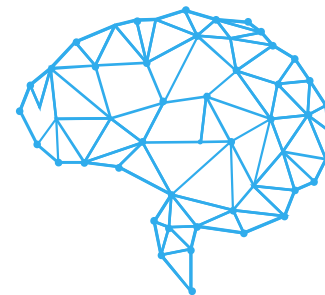
Para encontrar w y β :

$$\frac{\partial L}{\partial w_i} = 0$$

$$\frac{\partial L}{\partial \beta_i} = 0$$

M Á Q U I N A D E S V

P R O B L E M A P R I M A L



Se define un **problema de optimización** más **general**, que considere **tanto igualdades** como **desigualdades**: el **Problema Primal**

$$\min_w f(w)$$

Sujeto a:

$$h_i(w) = 0; i = 1, \dots, l$$

$$g_i(w) \leq 0; i = 1, \dots, k$$

El **lagrangiano** estaría dado por (donde α_i y β_i son los **multiplicadores de Lagrange**):

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

M Á Q U I N A D E S V

P R O B L E M A P R I M A L



Se va a **definir el problema primal** (donde P se refiere a **Primal**)

$$\theta_P(w) = \max_{\alpha, \beta: \alpha_i \geq 0} L(w, \alpha, \beta)$$

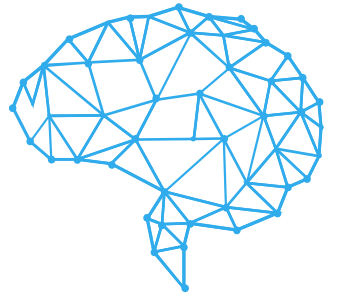
$$\theta_P(w) = \max_{\alpha, \beta: \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Sí w llega a **violar cualquiera** de las **dos restricciones** $h_i(w) \neq 0$ y $g_i(w) > 0$:

$$\theta_P(w) = \infty$$

M Á Q U I N A D E S V

P R O B L E M A P R I M A L



Por el **contrario**, si w satisface **ambas restricciones** $h_i(w) = 0$ y $g_i(w) \leq 0$:

$$\theta_P(w) = f(w)$$

En resumen:

$$\theta_P(w) = \begin{cases} f(w) & \text{si se satisfacen} \\ \infty & \text{no se satisfacen} \end{cases}$$

M Á Q U I N A D E S V

P R O B L E M A P R I M A L



Ahora se define la **minimización** de $\theta_P(w)$:

$$\min_w \theta_P(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} L(w, \alpha, \beta)$$

Nos damos cuenta que **volvemos** al **problema inicial** de **minimizar** respecto a w .

De manera específica, el valor primal p^* es:

$$p^* = \min_w \theta_P(w)$$

M Á Q U I N A D E S V

P R O B L E M A D U A L



Vamos a ver un **problema semejante**, pero ahora **minimizamos primero respecto a w** :

$$\theta_D(\alpha, \beta) = \min_w L(w, \alpha, \beta)$$

El subíndice **D** indica que estamos definiendo el **problema dual** de esta manera:

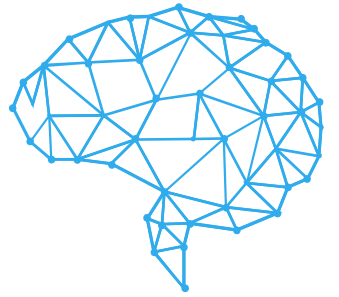
$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w L(w, \alpha, \beta)$$

Vemos que **solo** se **cambia** el **orden** de **minimización** y **maximización**. El **valor dual** está **definido** como:

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \theta_D(\alpha, \beta)$$

M Á Q U I N A D E S V

PROBLEMA DUAL y PRIMAL



A partir de **pruebas matemáticas** se llega a la **conclusión** de que:

$$d^* \leq p^*$$

$$\max \min \leq \min \max$$

Para **ciertas condiciones** se cumple **a veces** que:

$$d^* = p^*$$

M Á Q U I N A D E S V

C O N D I C I O N E S K K T



Para que se cumpla $d^* = p^*$ se deben cumplir ciertas condiciones. **Antes de definirlas, vamos a realizar las siguientes suposiciones:**

1. Se define que las funciones f y g_i son **convexas** (sí y solo sí la **Hessiana** de la **función** es **positiva semi-definida**).
2. Las funciones h_i son **afines (lineales)**.
3. Las funciones g_i son **factibles**, es decir $\exists w: g_i(w) < 0, \forall i$.

Dadas las suposiciones deben existir w^*, α^*, β^* , tal que w^* **satisfaga el problema primal**, así como α^*, β^* **satisfagan el problema dual**.

Es decir **que se cumpla:**

$$d^* = p^*$$

M Á Q U I N A D E S V

C O N D I C I O N E S K K T



Para esto, w^* , α^* , β^* deben **satisfacer** las **condiciones Karush-Kuhn-Tucker (KKT)**:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, n$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

M Á Q U I N A D E S V

C O N D I C I O N E S K K T



Sí existen parámetros w^*, α^*, β^* que **satisfagan** las **condiciones Karush-Kuhn-Tucker (KKT)**, se **garantiza** la **solución** para **ambos** problemas, el **primal** y el **dual**.

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, n$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



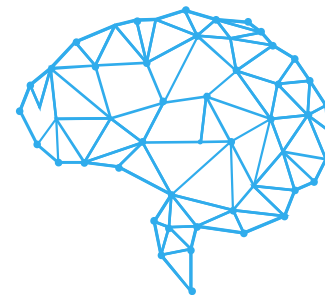
Regresando al problema del clasificador marginal, se sabe que teníamos el siguiente problema de optimización (primal):

$$\max_{w,b} \frac{1}{\|w\|} = \min_{w,b} \frac{1}{2} \|w\|^2$$

Sujeto a:

$$y^{(i)}(w^T x^{(i)} + b) \geq 1; i = 1, \dots, m$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Transformando la restricción en base a la nomenclatura usada en la dualidad de Lagrange:

$$\max_{w,b} \frac{1}{\|w\|} = \min_{w,b} \frac{1}{2} \|w\|^2$$

Sujeto a:

$$g_i(w, b) = -y^{(i)}(w^T x^{(i)} + b) + 1 \leq 0; i = 1, \dots, m$$

Observamos, que **solo** existe **una desigualdad**, por lo que **solo** tendremos los **multiplicadores de Lagrange** α_i .

De igual forma, **no solo** tenemos el **parámetro** w , pero **también** el **parámetro** b .

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Construimos el Lagrangiano como:

$$L(\mathbf{w}, \mathbf{b}, \alpha) = f(\mathbf{w}) + \sum_{i=1}^m \alpha_i g_i(\mathbf{w})$$

$$L(\mathbf{w}, \mathbf{b}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i [-\mathbf{y}^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + \mathbf{b}) + \mathbf{1}]$$

$$L(\mathbf{w}, \mathbf{b}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i [-\hat{\mathbf{y}}^{(i)} + \mathbf{1}]$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Debido a la **condición de complementariedad (restricción activa)** en **KKT**:

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

Para los **multiplicadores de Lagrange** $\alpha_i > 0$ se tiene que:

$$g_i(w) = [-\hat{\gamma}^{(i)} + 1] = 0$$

Por lo que el **margen funcional** $\gamma^{(i)}$ tendrá que ser **igual a 1** para el dato de entrenamiento i :

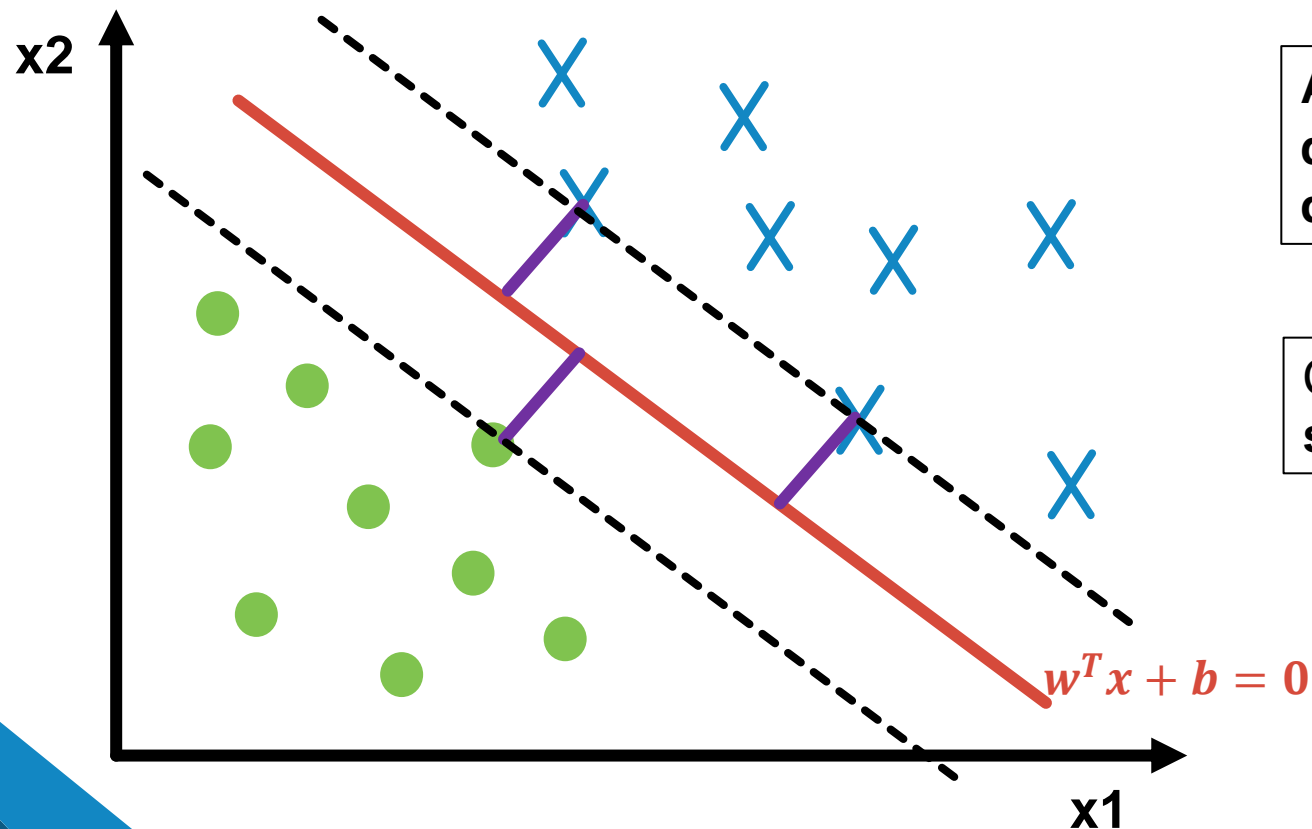
$$\hat{\gamma}^{(i)} = 1$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Lo que esto significa es que los **datos de entrenamiento** con margen $\gamma^{(i)} = 1$ serán los **más cercanos al hiperplano** y por lo general sus $\alpha_i \neq 0$ porque $g_i = 0$.



A estos datos de entrenamiento, que cumplen con $\gamma^{(i)} = 1$, se les denominarán **vectores de soporte**.

Como hay **pocos vectores de soporte**, la mayoría de $\alpha_i = 0$.

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Se **desarrollará** ahora la **versión dual** del **problema**:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1]$$

Se **define** $\theta_D(\alpha)$:

$$\theta_D(\alpha) = \min_{w, b} L(w, b, \alpha)$$

Se **minimiza** respecto a w :

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i [y^{(i)}(x^{(i)})] = 0$$
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Se minimiza respecto a b :

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Se tienen **ambas** restricciones:

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Se **sustituye** el mínimo $\mathbf{w} = \sum_{i=1}^m \alpha_i \mathbf{y}^{(i)} \mathbf{x}^{(i)}$ en el **Lagrangiano**:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i [\mathbf{y}^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1]$$

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \alpha_i [\mathbf{y}^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1]$$

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \left(\sum_{i=1}^m \alpha_i \mathbf{y}^{(i)} \mathbf{x}^{(i)} \right)^T \sum_{j=1}^m \alpha_j \mathbf{y}^{(j)} \mathbf{x}^{(j)} - \sum_{i=1}^m \alpha_i [\mathbf{y}^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1]$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

$$L(w, b, \alpha) = \frac{1}{2} \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} - \sum_{i=1}^m \alpha_i \left[y^{(i)} \left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b \right) - 1 \right]$$

$$L(w, b, \alpha) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \alpha_i y^{(i)} b + \sum_{i=1}^m \alpha_i$$

$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)} + \sum_{i=1}^m \alpha_i$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Se **sustituye** ahora la **restricción** de b de la última ecuación $\sum_{i=1}^m \alpha_i y^{(i)} = 0$:

$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)} + \sum_{i=1}^m \alpha_i$$

$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - 0 + \sum_{i=1}^m \alpha_i$$

$$\min_{w,b} L(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)}$$

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Recordamos la **forma** del **problema dual**, para el **caso** del **clasificador marginal**:

$$\max_{\alpha} W(\alpha) = \max_{\alpha} \min_{w, b} L(w, b, \alpha)$$

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

Sujeto a:

$$\alpha_i \geq 0, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Podemos **observar** que el **problema dual** **satisface** las **condiciones KKT**.

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, n$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

SE SATISFACE $d^* = p^*$

max min = min max

PODEMOS RESOLVER EL PROBLEMA PRIMAL RESOLVIENDO EL DUAL

M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Recordemos que **ya** tenemos la **expresión** de w , pero está en **función** de α_i :

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

Cómo se **satisface** $d^* = p^*$ podemos **resolver** el **problema dual** para **obtener** α_i :

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

Finalmente, con las α_i podemos obtener b **resolviendo** el **problema primal** con la w obtenida.

M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

Para obtener b sabemos que la **ecuación** del **plano** es $w^T x^{(i)} + b = 0$, donde:

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

Despejamos b :

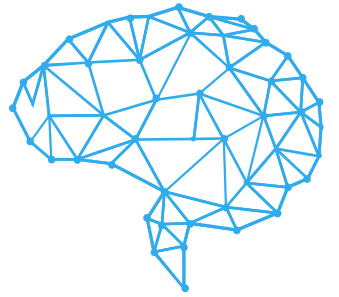
$$b = -w^T x^{(i)}$$

Como $\alpha_i = 0$ para **vectores** que **no son** de **soporte**, esto **solo aplica** para **vectores** de **soporte** (en este caso, **vectores** que están **más cerca** del **hiperplano**):

$$\max_{i:y^{(i)}=-1} -w^T x^{(i)}$$

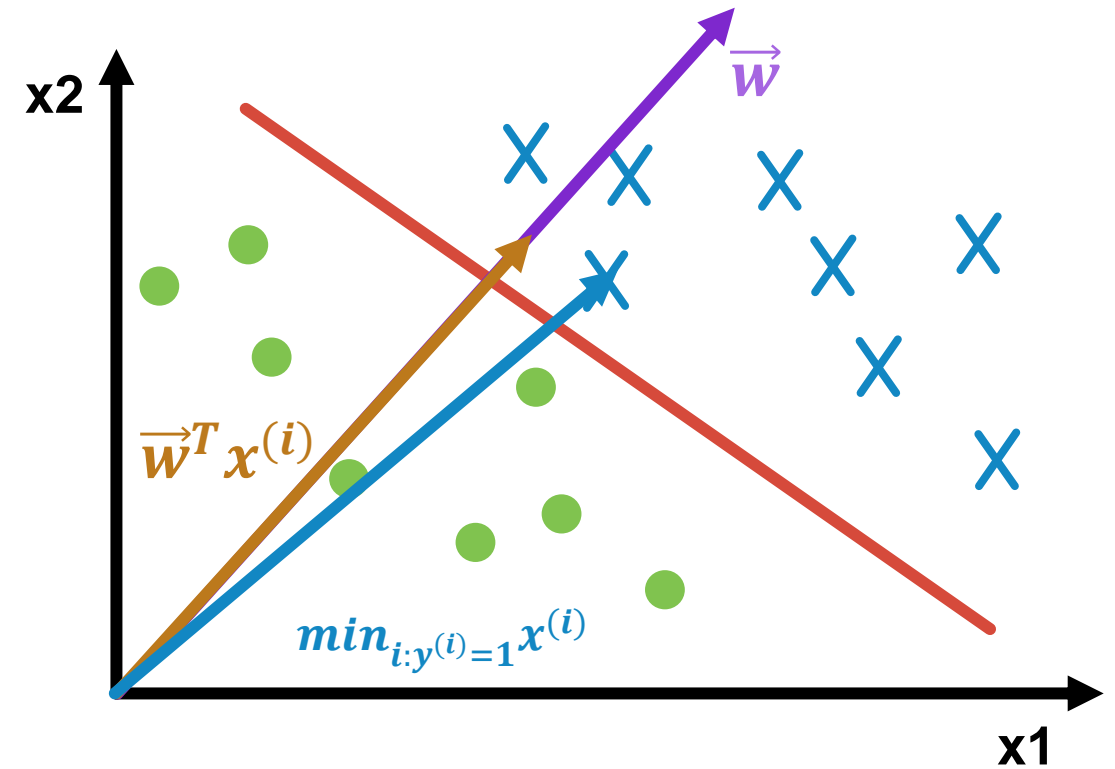
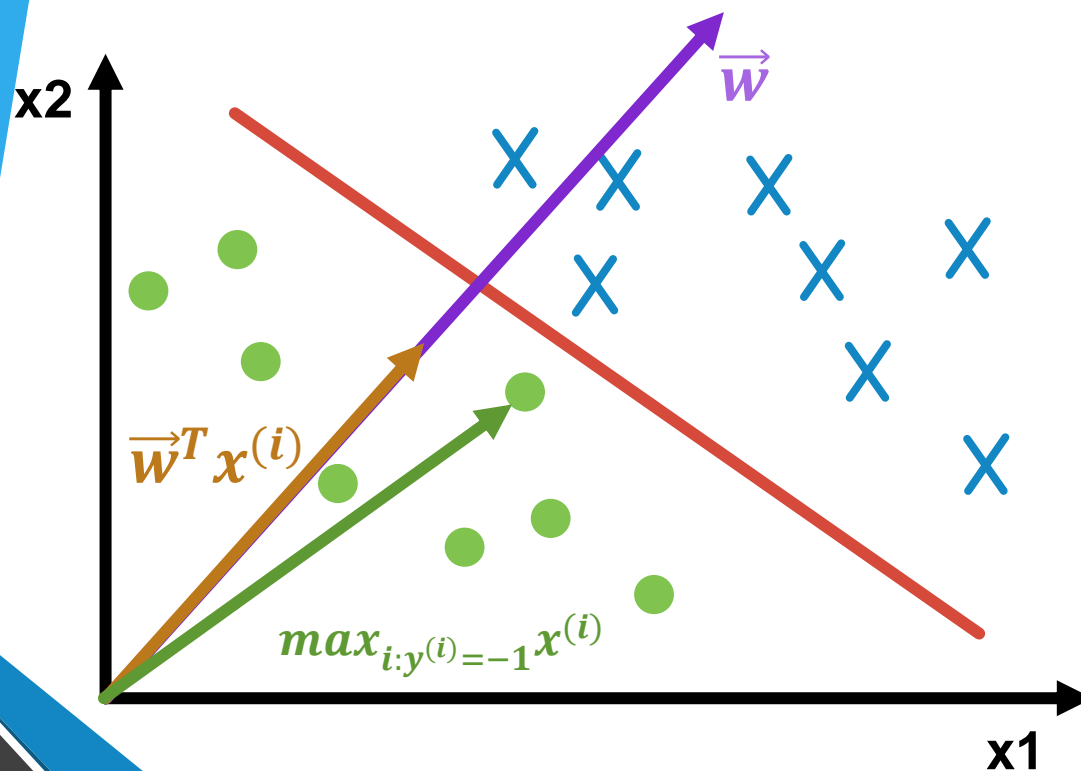
$$\min_{i:y^{(i)}=1} w^T x^{(i)}$$

M Á Q U I N A D E S V



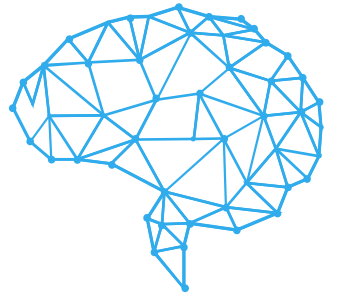
CLASIFICADOR MARGINAL ÓPTIMO

Se representan ambos vectores $\max_{i:y^{(i)}=-1} -w^T x^{(i)}$ y $\min_{i:y^{(i)}=1} w^T x^{(i)}$:



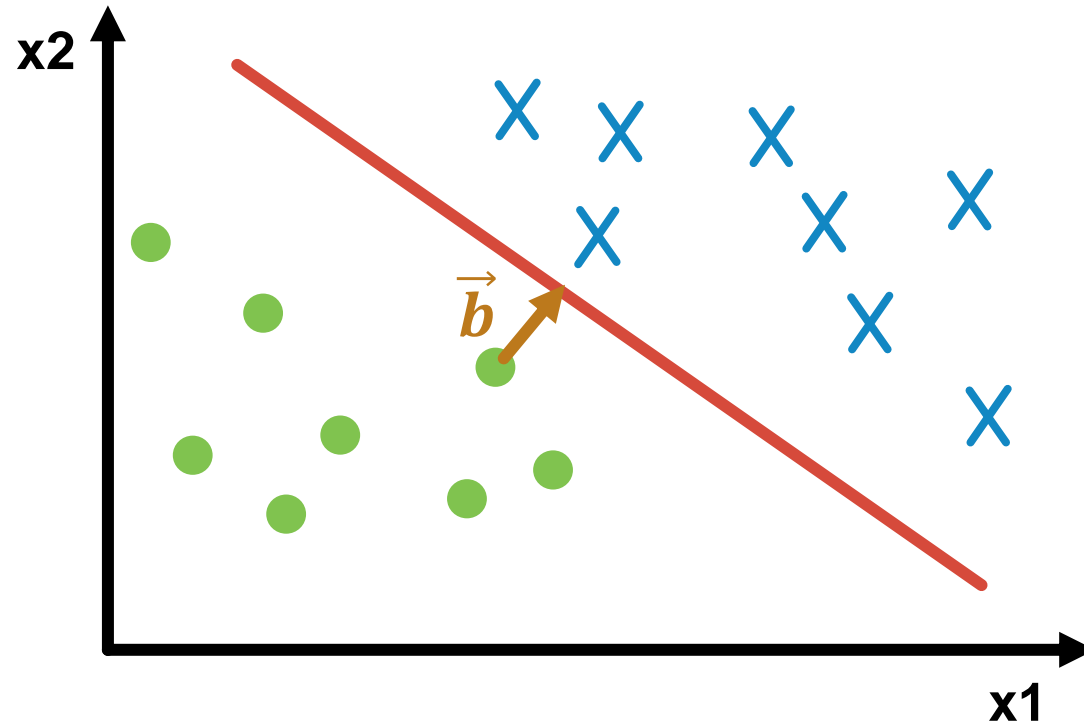
M Á Q U I N A D E S V

CLASIFICADOR MARGINAL ÓPTIMO



Por lo tanto ***b*** sería:

$$b^* = -\frac{\max_{i:y(i)=-1} w^{*T} x^{(i)} + \min_{i:y(i)=1} w^{*T} x^{(i)}}{2}.$$



M Á Q U I N A D E S V



CLASIFICADOR MARGINAL ÓPTIMO

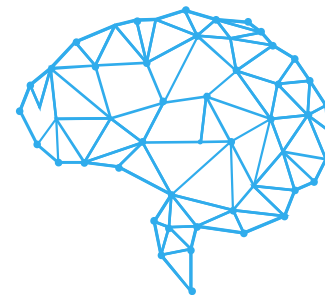
Cómo ya **calculamos** todos los **parámetros**, vamos a **expresar** las **predicciones** en función de **productos puntos**:

$$h_{w,b}(x) = g(w^T x + b) = g\left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}\right)^T x + b\right)$$

$$h_{w,b}(x) = g\left(\sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b\right)$$

Como $\alpha_i = 0$ para **vectores** que **no son** de **soporte**. Es decir, para **hacer** una **nueva predicción** se **calcularían** las **proyecciones solamente** con los **vectores de soporte**.

M Á Q U I N A D E S V
K E R N E L S



Para **comenzar** vamos a **realizar** un **mapeado** de las **características** $x^{(i)}$ a un **espacio dimensional más alto**. Vamos a ver el **caso** con **una** sola **característica** $x^{(i)} \in \mathbb{R}$:

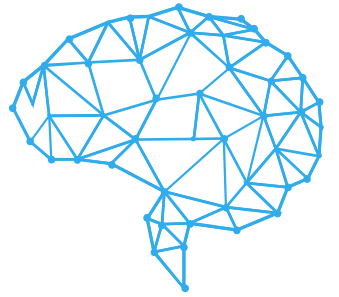
$$\phi(x^{(i)}) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

De esta manera, podemos **cambiar todo** lo que hemos hecho en el **algoritmo** de la siguiente **forma**:

$$\langle x^{(i)}, x^{(j)} \rangle \rightarrow \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle = \phi(x^{(i)})^T \phi(x^{(j)})$$

M Á Q U I N A D E S V

K E R N E L S



Al **producto punto** de las **características** en el **espacio altamente dimensional** se le llama **Kernel**:

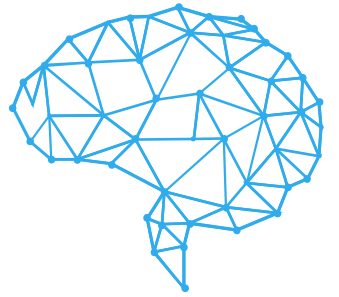
$$K(x, z) = \phi(x)^T \phi(z)$$

Por lo tanto, el **algoritmo** va a **calcular** las **características** en **este espacio** de **mayores dimensiones**.

En principio, si ϕ **realiza** un mapeo a un **espacio altamente dimensional** pareciera que **calcular** $K(x, z)$ sería **muy costoso** computacionalmente. De hecho, sería **imposible** si se **llegara** a realizar un **mapeo** en **infinitas dimensiones**.

M Á Q U I N A D E S V

K E R N E L S



Aún así, existen **casos especiales** donde el **costo computacional** no se **dispara**. Son estos **casos** los que nos interesan para **poder incluirlos** en la **máquina de vectores de soporte**.

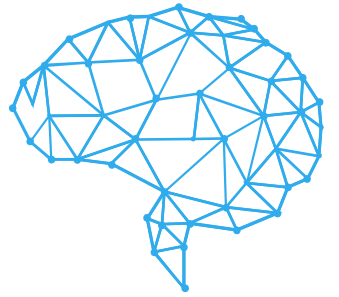
EJEMPLO:

Se tienen **dos vectores** $x, z \in \mathbb{R}^n$ por lo tanto, el **kernel** estaría dado por:

$$K(x, z) = (x^T z)^2$$

$$K(x, z) = \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right)$$

M Á Q U I N A D E S V
K E R N E L S



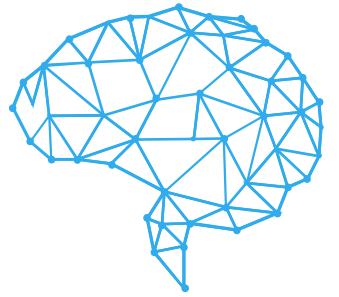
EJEMPLO:

$$K(x, z) = \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right)$$

$$K(x, z) = \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j$$

$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j$$

M Á Q U I N A D E S V
K E R N E L S



EJEMPLO:

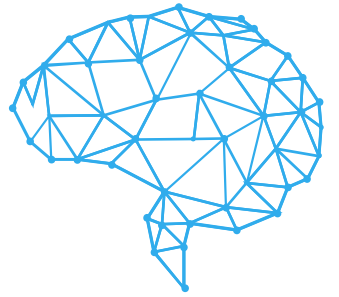
$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j$$

Se **observa** que se este **kernel corresponde** a la siguiente **transformación** ϕ :

$$\phi = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_1 x_n \\ \vdots \\ x_n x_n \end{bmatrix} \in \mathbb{R}^{n^2}$$

M Á Q U I N A D E S V

K E R N E L S



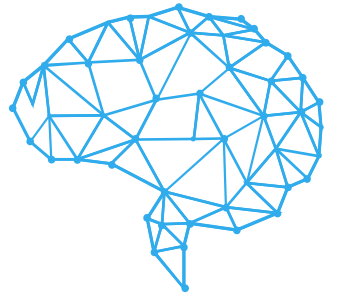
Es importante ver que, aunque calcular ϕ necesita n^2 **operaciones**, para el caso del **kernel solo** se necesitan $2n + 1$ **operaciones**

$$\phi = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_1 x_n \\ \vdots \\ x_n x_n \end{bmatrix} \in \mathbb{R}^{n^2}$$

$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j = (x^T z)^2$$

M Á Q U I N A D E S V

K E R N E L S



EJEMPLO 2:

De manera general, se puede obtener un **kernel** para poder **mapear polinomios** de **cualquier orden d** :

$$K(x, z) = (x^T z + c)^d$$

Veamos el **ejemplo** con $d = 2$ y $n = 3$:

$$K(x, z) = (x^T z + c)^2$$

$$K(x, z) = (x^T z)^2 + 2cx^T z + c^2$$

M Á Q U I N A D E S V

K

E

R

N

E

L

S



EJEMPLO 2:

$$K(x, z) = (x^T z)^2 + 2cx^T z + c^2$$

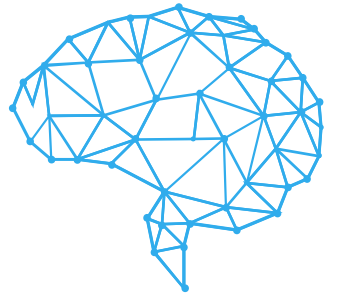
$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j + \sum_{i=1}^n 2cx_i z_i + c^2$$

$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j + \sum_{i=1}^n 2cx_i z_i + c^2$$

$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j + \sum_{i=1}^n (\sqrt{2cx_i})(\sqrt{2cz_i}) + c^2$$

M Á Q U I N A D E S V

K E R N E L S



EJEMPLO 2:

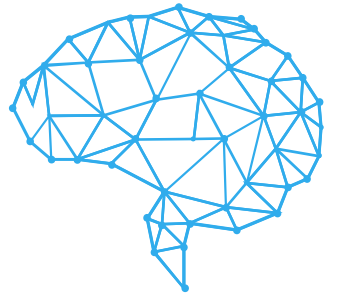
$$K(x, z) = \sum_{i,j=1}^n x_i x_j z_i z_j + \sum_{i=1}^n (\sqrt{2c} x_i)(\sqrt{2c} z_i) + c^2$$

Observamos la transformación ϕ .

$$\phi(x) =$$

$$\begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ \sqrt{2c} x_3 \\ c \end{bmatrix}$$

M Á Q U I N A D E S V
K E R N E L S



Recordando que para el **caso general**:

$$K(x, z) = (x^T z + c)^d$$

La **matriz** ϕ requeriría de $O(n^d)$ **cálculos** para computarse, **pero el kernel solo necesita** $O(n)$ **cálculos.**

Por lo tanto, se ha **encontrado** una **forma** de trabajar con **vectores** en un **espacio** de **infinitas dimensiones** sin tener que **representarlos explícitamente.**

M Á Q U I N A D E S V
K E R N E L S

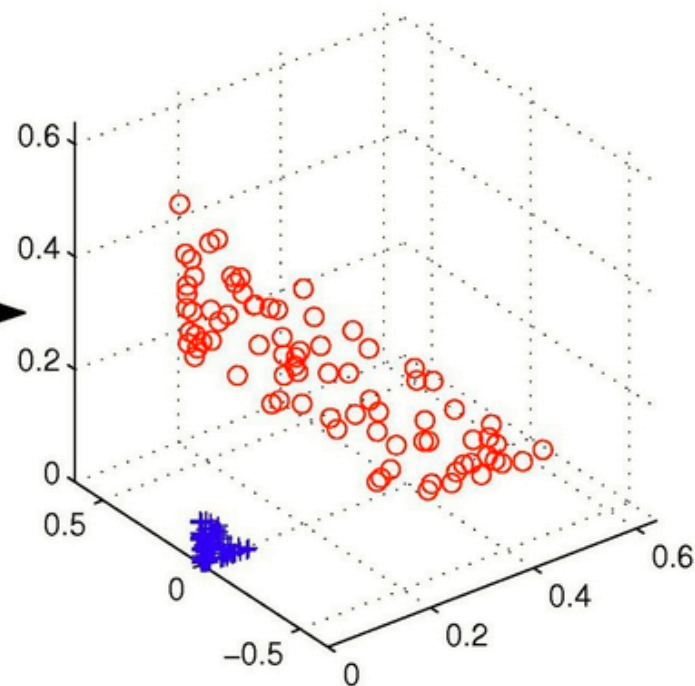
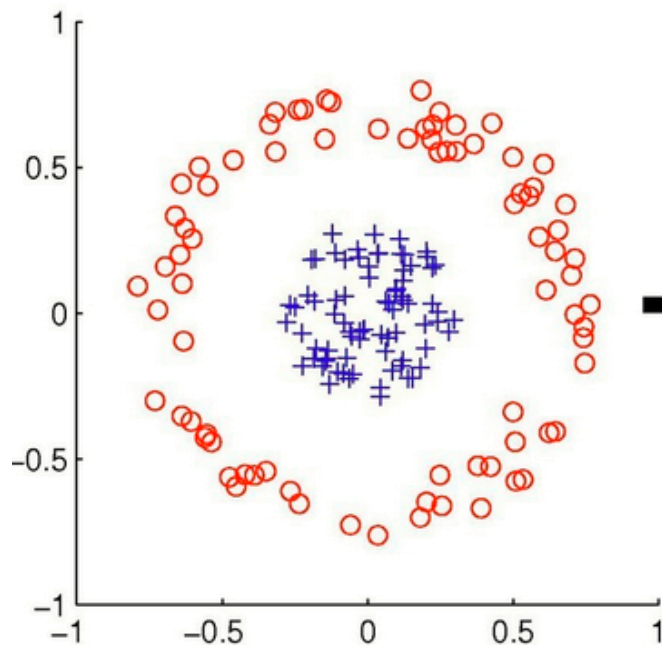


¿PORQUÉ NECESITAMOS EL KERNEL?

¿PORQUÉ TRANSFORMAR LOS VECTORES?

M Á Q U I N A D E S V

K E R N E L S

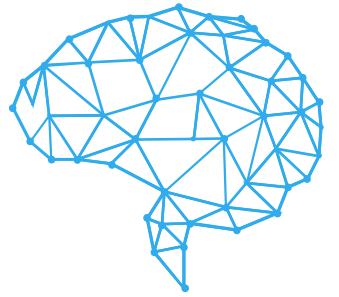


$$\phi = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2$$

$$\phi = \begin{bmatrix} x_1 x_1 \\ \sqrt{2} x_1 x_2 \\ x_2 x_2 \end{bmatrix} \in \mathbb{R}^3$$

M Á Q U I N A D E S V

K E R N E L S



Ahora veamos un **kernel** muy importante: el **kernel Gaussiano** o **kernel de función de base radial**:

$$K(x, z) = \phi(x)^T \phi(z) = e^{-\left(\frac{\|x-z\|^2}{2\sigma^2}\right)}$$

Se puede **entender**, de **manera intuitiva**, que:

- Sí x y z se encuentran **muy cercanos**, el **kernel** $K(x, z)$ será **1**.
- Sí x y z se encuentran **muy lejanos**, el **kernel** $K(x, z)$ será **0**.

Este **kernel corresponde** a un **mapeado** $\phi(x) \in \mathbb{R}^\infty$

M Á Q U I N A D E S V

K E R N E L S



Vamos a **demostrar** que $K(x, z) = \phi(x)^T \phi(z) = e^{-\left(\frac{\|x-z\|^2}{2\sigma^2}\right)}$ nos da un **mapeado** $\phi(x) \in \mathbb{R}^\infty$.

Veremos el **ejemplo** más sencillo $x, z \in \mathbb{R}$ y $\sigma^2 = \frac{1}{2}$.

$$K(x, z) = \phi(x)^T \phi(z) = e^{-\left(\frac{(x-z)^2}{2\sigma^2}\right)} = e^{-(x-z)^2}$$

$$K(x, z) = e^{-x^2+2xz-z^2} = e^{-x^2} e^{-z^2} e^{2xz}$$

Recordando que $e^k = \sum_{n=0}^{\infty} \frac{k^n}{n!}$:

$$K(x, z) = e^{-x^2+2xz-z^2} = e^{-x^2} e^{-z^2} \sum_{n=0}^{\infty} \frac{(2xz)^n}{n!}$$

M Á Q U I N A D E S V

K E R N E L S



$$K(x, z) = e^{-x^2} e^{-z^2} \sum_{n=0}^{\infty} \frac{2^n x^n z^n}{n!}$$

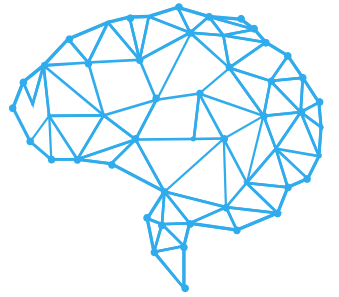
$$K(x, z) = \sum_{n=0}^{\infty} \frac{2^n \left[e^{-x^2} x^n \right] \left[e^{-z^2} z^n \right]}{n!} = \phi(x)^T \phi(z)$$

Por lo tanto:

$$\phi(x) = \begin{bmatrix} e^{-x^2} \\ \frac{\sqrt{2^1}}{1!} e^{-x^2} x \\ \frac{\sqrt{2^2}}{2!} e^{-x^2} x^2 \\ \frac{\sqrt{2^3}}{3!} e^{-x^2} x^3 \\ \vdots \end{bmatrix}$$

M Á Q U I N A D E S V

K E R N E L S



Hasta el momento se han **planteado kernels**, pero **no** se ha **establecido** cuando un **kernel es válido** o **no**. Es decir, que $\forall x, z \exists K(x, z) = \phi(x)^T \phi(z)$.

Para esto, se **deben cumplir** las **dos condiciones** de Mercer:

1. **Simetría del kernel:**

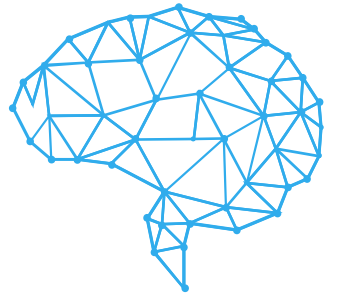
$$K(x, z) = K(z, x)$$

2. La **matriz K** es **positiva semi definida**:

$$K = \begin{bmatrix} K(x^{(1)}, x^{(1)}) & K(x^{(1)}, x^{(2)}) & \dots & K(x^{(1)}, x^{(m)}) \\ K(x^{(2)}, x^{(1)}) & K(x^{(2)}, x^{(2)}) & \dots & K(x^{(2)}, x^{(m)}) \\ \vdots & \vdots & \ddots & \vdots \\ K(x^{(m)}, x^{(1)}) & K(x^{(m)}, x^{(2)}) & \dots & K(x^{(m)}, x^{(m)}) \end{bmatrix}$$

M Á Q U I N A D E S V

K E R N E L S



Teorema de Mercer:

Establezcamos que $K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$. Entonces, para que K sea un **kernel válido** de Mercer es **necesario y suficiente** que **para cualquier** $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, ($m < \infty$), la correspondiente **matriz de kernel** sea simétrica positiva semi definida.

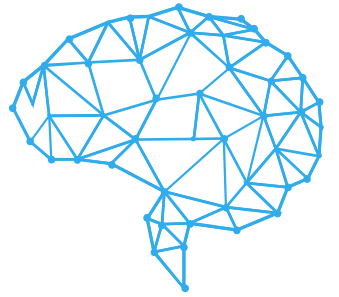
Aplicaciones

Realmente, los **kernels no solo** pueden aplicarse **para SVM**, pero para **cualquier algoritmo** que pueda ser **representado** en **términos** de **productos punto**.

- Regresión lineal.
- Regresión logística.
- Perceptrón.
- Procesos gaussianos.
- PCA.

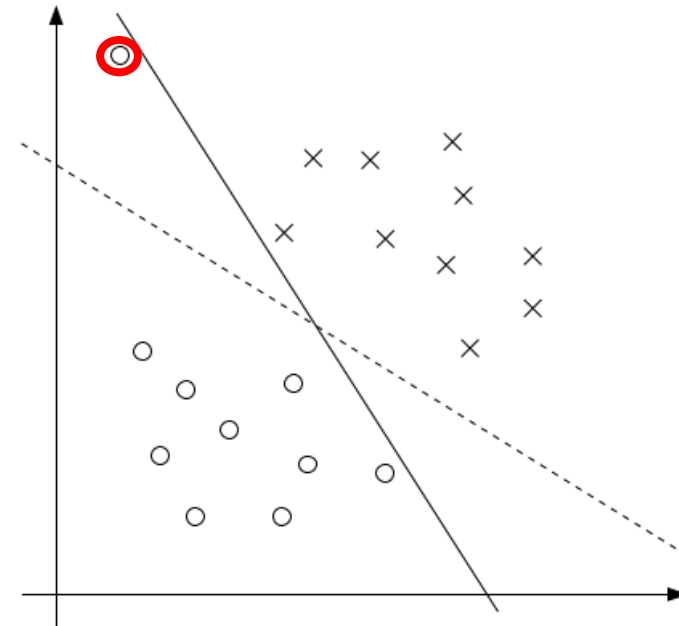
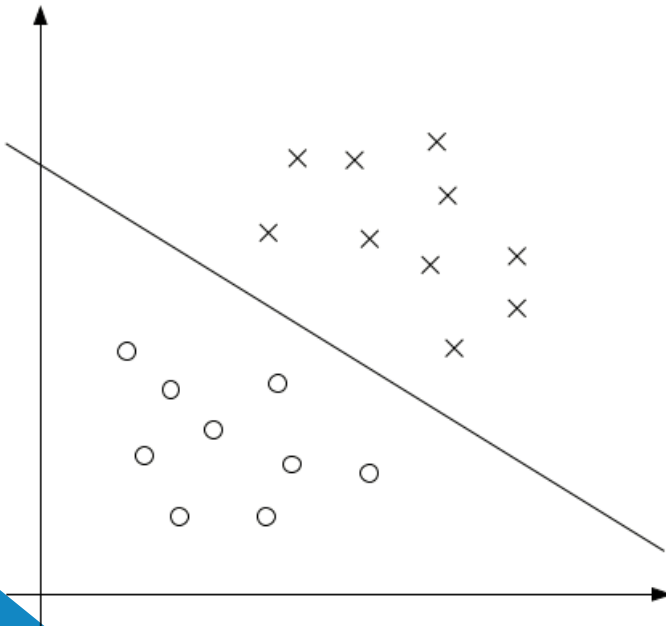
M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



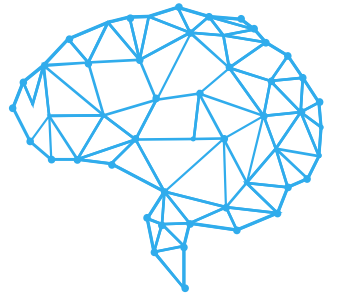
Aún cuando los **kernels** pueden **trasladar** el **problema** a un **espacio** donde los datos sean **separables linealmente**, **no existe** una **garantía** para que esto **siempre ocurra**.

Además, existe el **problema** de los **datos atípicos** que pueden llegar a **cambiar** el **hiperplano** **drásticamente** dejando **poco margen**.



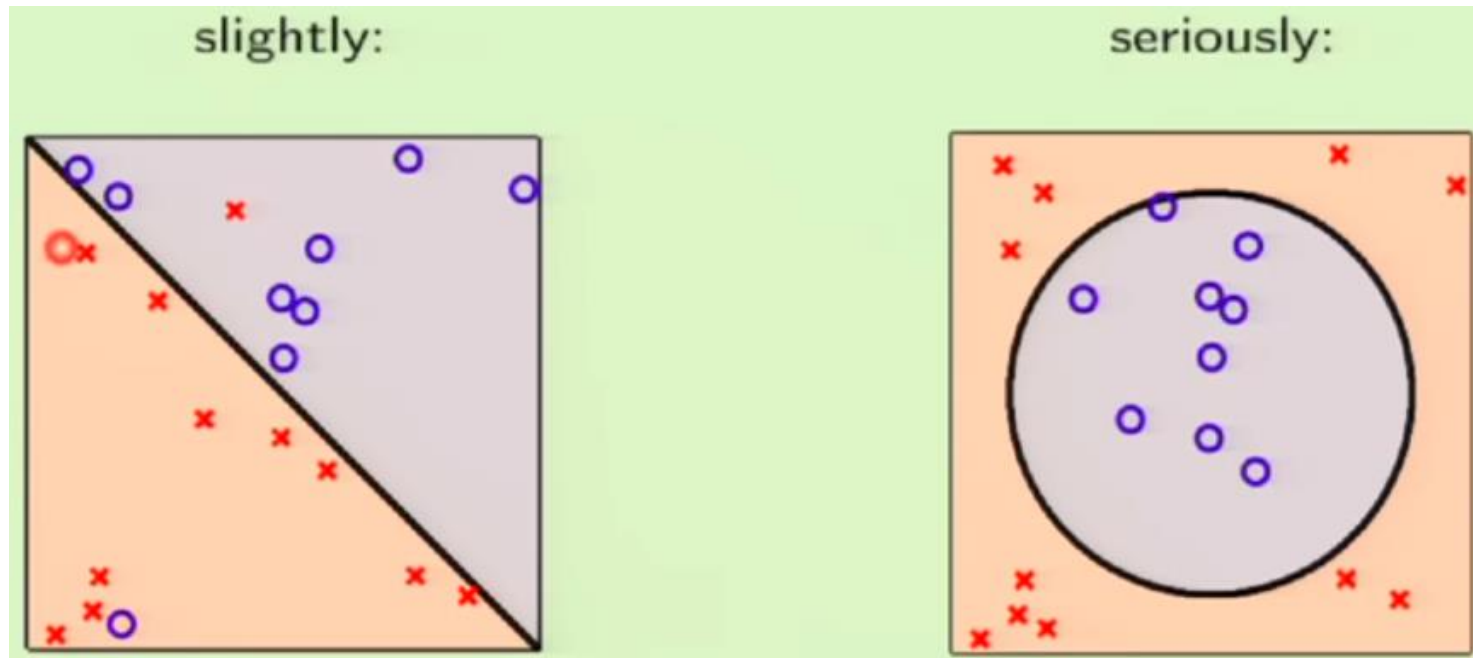
M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



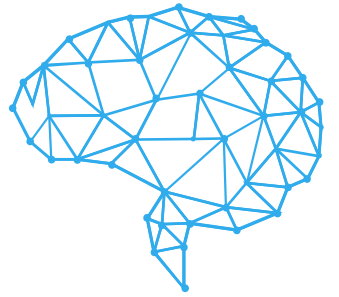
De manera **general**, puede haber **dos tipos** de no linealidad:

- **Poca no linealidad** por datos atípicos (**solución: regularización**).
- **No linealidad inherente** a los datos (**solución: kernels**).



M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L₁

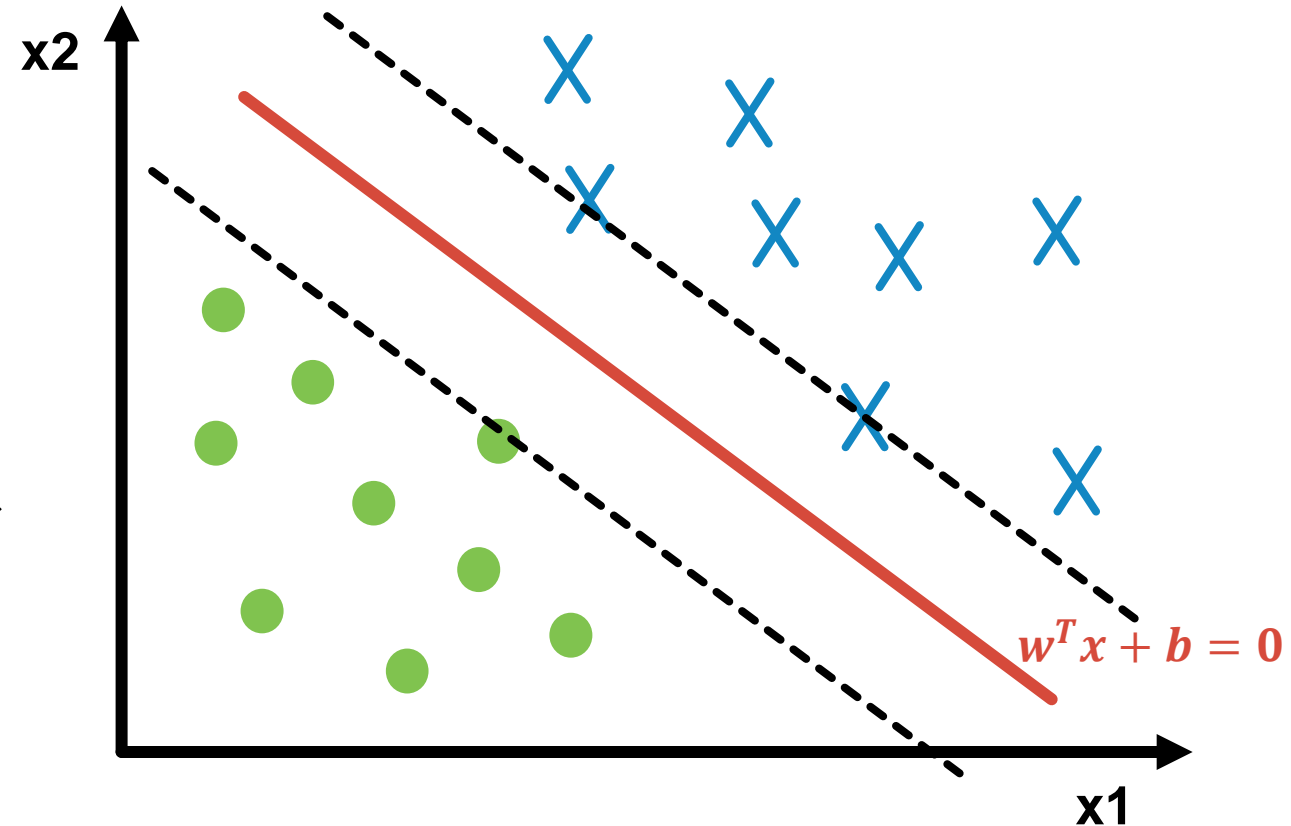


Vamos a **resolver** el **problema** de los **datos atípicos**. Para esto, **recordemos** cual era el **problema original**:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

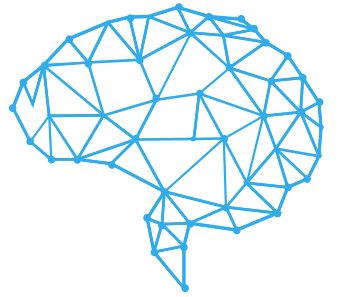
Sujeto a:

$$y^{(i)}(w^T x^{(i)} + b) \geq 1; i = 1, \dots, m$$

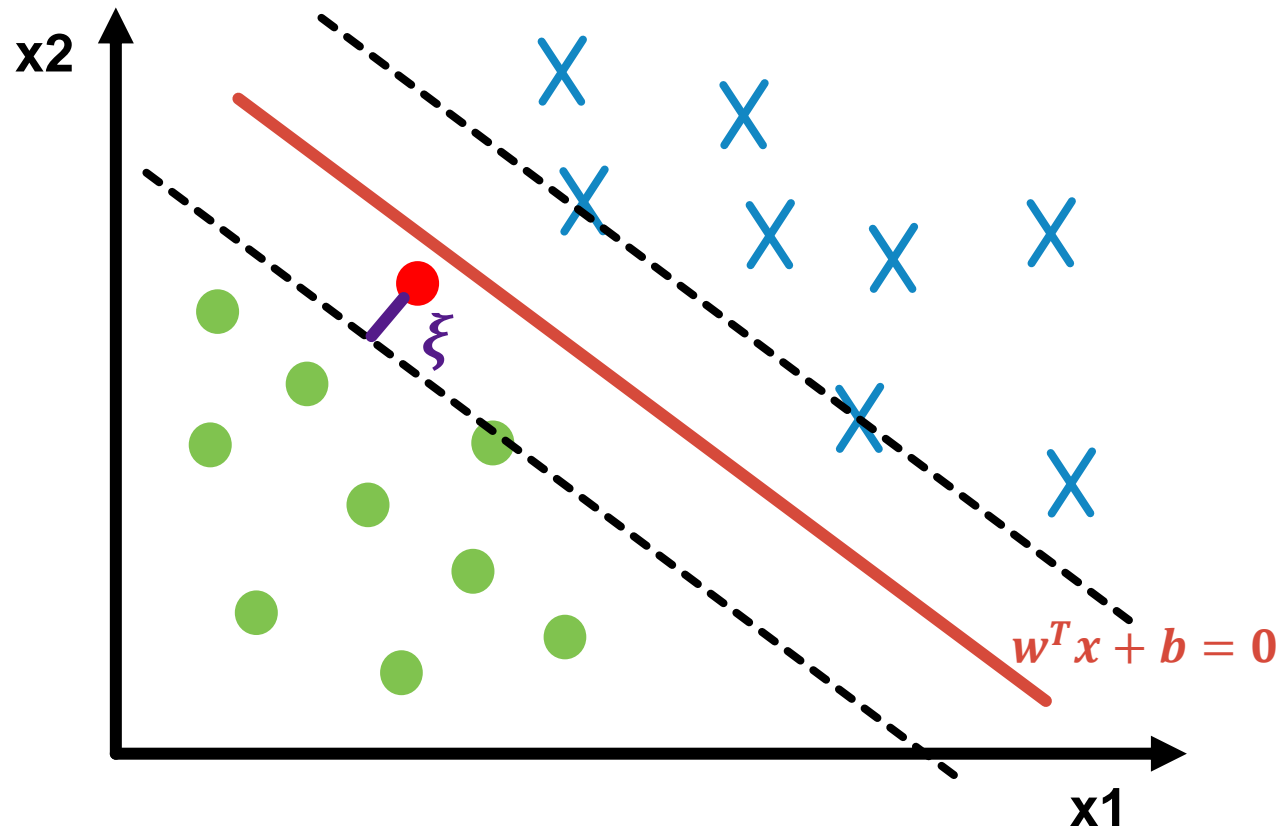


M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1

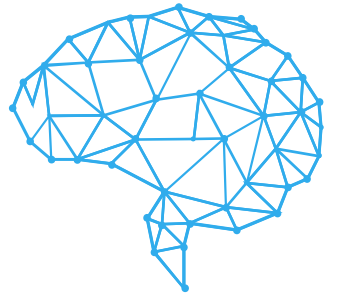


Ahora consideraremos **un error** como una **violación del margen** por una **cantidad $\xi \geq 0$** :



M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L₁



Consideramos el **total de violación generada** por **todos** los **datos** como (**regularización L1**):

$$error\ total = \sum_{i=1}^m \xi_i$$

El problema se **convierte** en:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

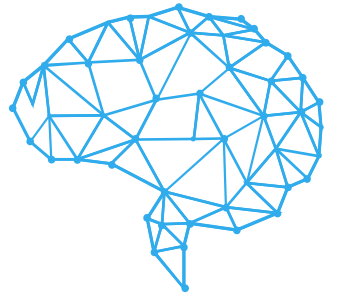
Sujeto a:

$$y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i; i = 1, \dots, m$$

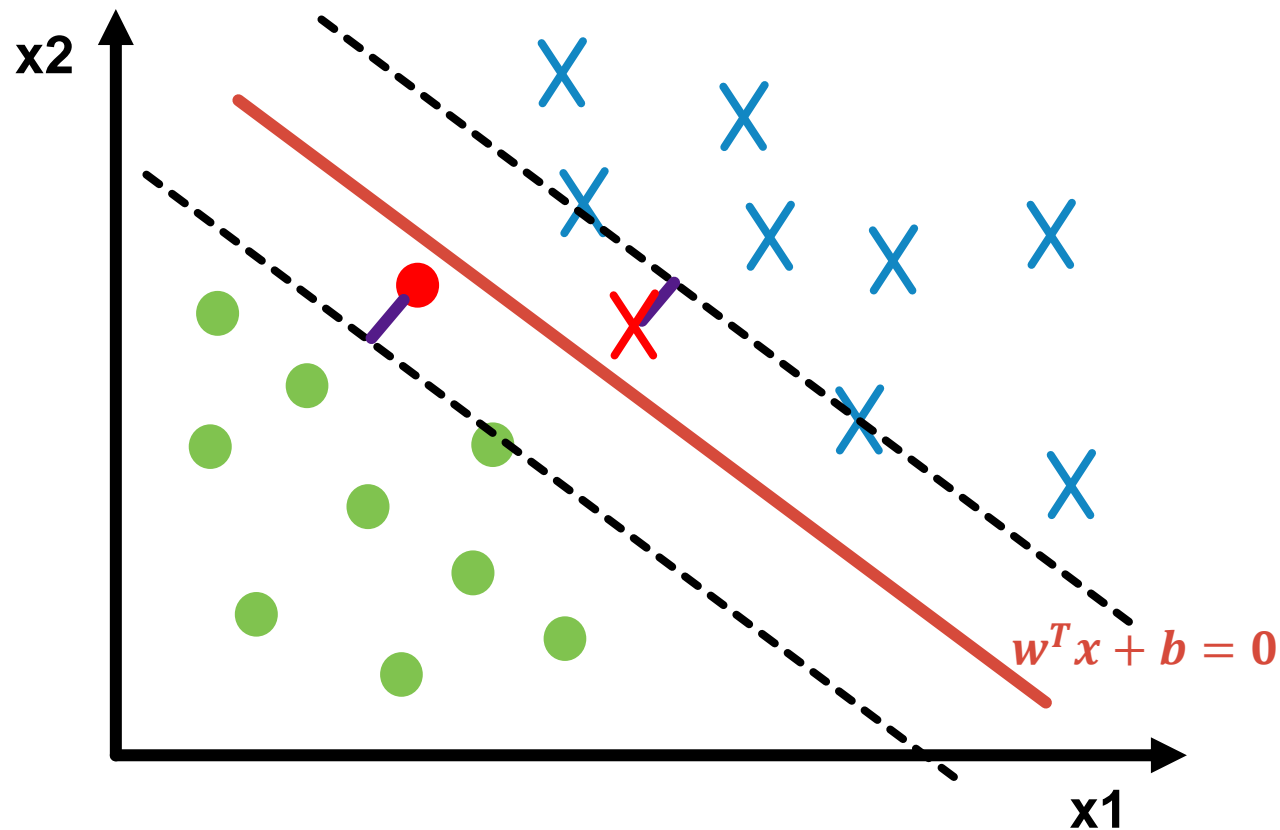
$$\xi_i > 0$$

M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1

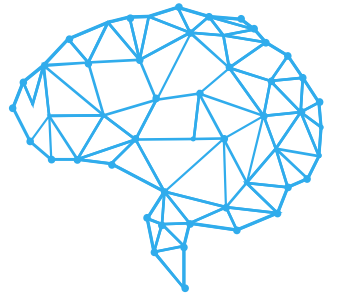


Es decir, ahora sí **se permiten datos** con **margen funcional < 1** que **violen el hiperplano**, solo que **cada vez** que **exista una violación** por un dato, **se penalizara la función objetivo** por una cantidad $C\xi_i$.



M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



Con esta **nueva función objetivo**, se construye el **Lagrangiano**:

$$L(\mathbf{w}, b, \xi, \alpha, r) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi_i] - \sum_{i=1}^m r_i \xi_i$$

donde los **multiplicadores de Lagrange** son α_i y r_i .

M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



Si se **resuelve** de nuevo el **problema dual**, **derivando** respecto a w , b y ξ_i se obtiene lo siguiente:

$$\frac{\partial L(w, b, \xi, \alpha, r)}{\partial b} = \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

$$\nabla_w L(w, b, \xi, \alpha, r) = w - \sum_{i=1}^m \alpha_i [y^{(i)}(x^{(i)})] = 0$$

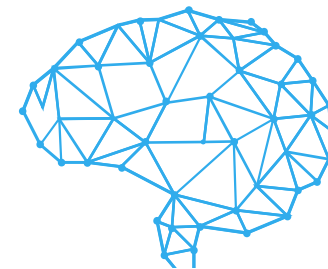
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial L(w, b, \xi, \alpha, r)}{\partial \xi_i} = C - \alpha_i - r_i = 0$$



M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



Sustituyendo w :

$$\min L(w, b, \xi, \alpha, r)$$

=

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m r_i \xi_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \alpha_i \left[y^{(i)} \left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b \right) - 1 \right]$$

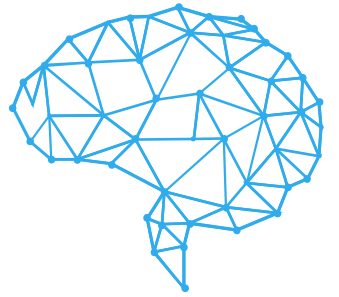
$$\min L(w, b, \xi, \alpha, r) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} + (C - r_i - \alpha_i) \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i \left[y^{(i)} \left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b \right) - 1 \right]$$

Sustituyendo $C - \alpha_i - r_i = 0$:

$$\min L(w, b, \xi, \alpha, r) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \alpha_i \left[y^{(i)} \left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b \right) - 1 \right]$$

M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



$$\min L(w, b, \xi, \alpha, r) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \alpha_i \left[y^{(i)} \left(\left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x^{(i)} + b \right) - 1 \right]$$

$$\min L(w, b, \xi, \alpha, r) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)} + \sum_{i=1}^m \alpha_i$$

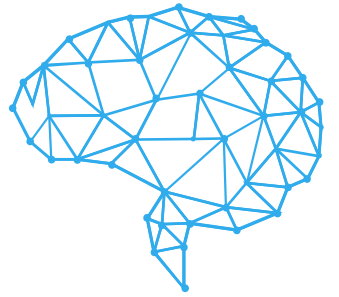
$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)} + \sum_{i=1}^m \alpha_i$$

Sustituyendo $\sum_{i=1}^m \alpha_i y^{(i)} = 0$:

$$\min L(w, b, \xi, \alpha, r) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y^{(i)} \alpha_j y^{(j)} x^{(i)T} x^{(j)} + \sum_{i=1}^m \alpha_i$$

M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



Se tiene el **mismo problema dual** pero **con regularización**:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

Lo **único** que **cambia** es la **restricción** para α_i .

M Á Q U I N A D E S V

SVM CON REGULARIZACIÓN L_1



Las nuevas condiciones KKT complementarias serían:

$$\alpha_i = 0 \Rightarrow y^{(i)}(w^T x^{(i)} + b) \geq 1 \text{ A}$$

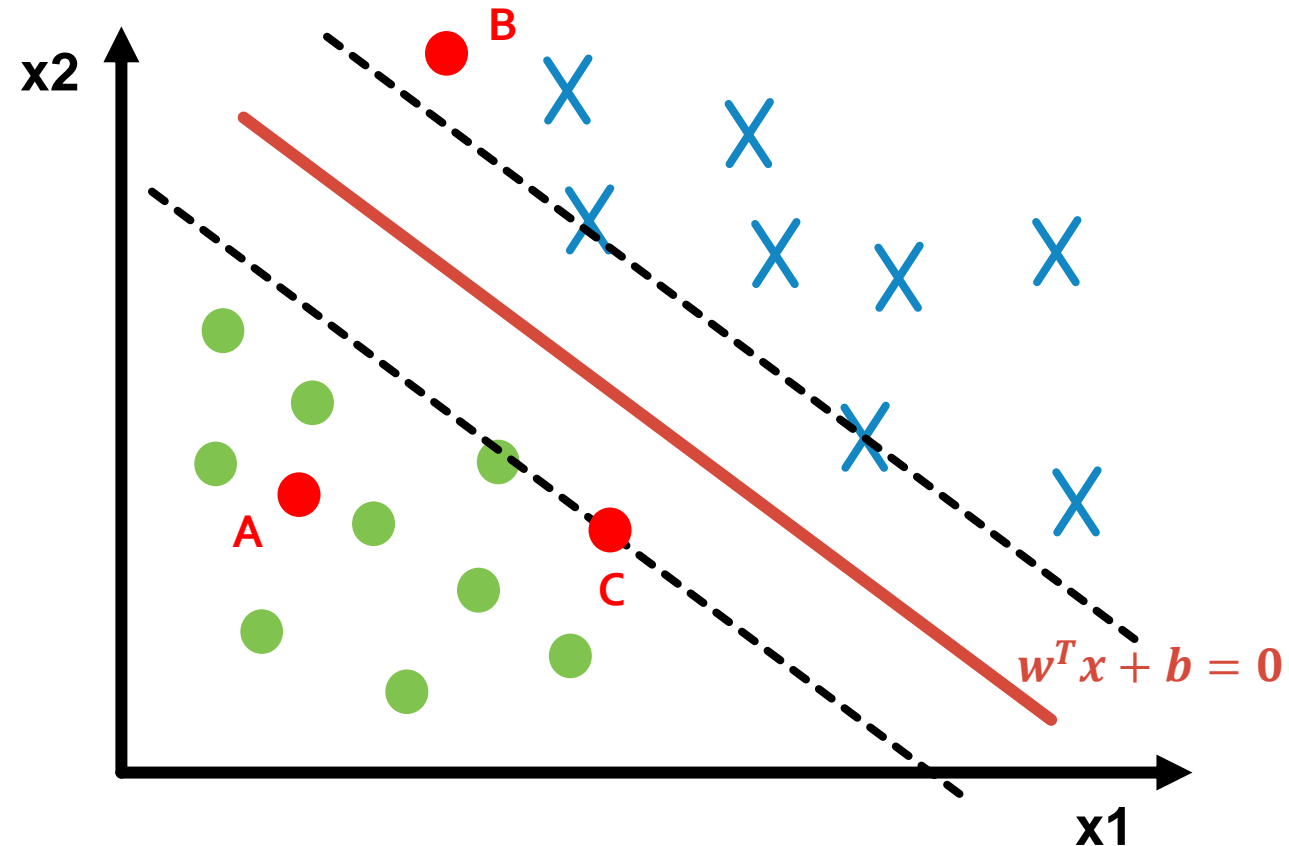
$$\alpha_i = C \Rightarrow y^{(i)}(w^T x^{(i)} + b) \leq 1 \text{ B}$$

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w^T x^{(i)} + b) = 1 \text{ C}$$

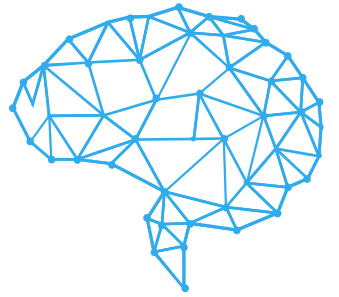
A = vectores **NO** de soporte

B = vectores de **soporte NO** marginales

C = vectores de **soporte marginales**):



M Á Q U I N A D E S V A L G O R I T M O S M O



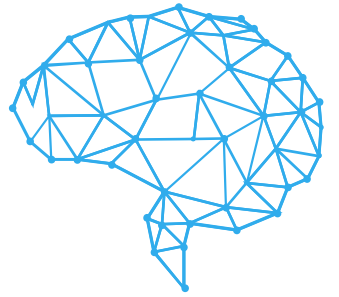
Hasta el momento **no hemos** discutido la **forma** en **como solucionar** el **problema dual**.

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

Para esto, se verá antes **otro algoritmo** denominado **ascenso de coordenadas** que permita **hacer eficiente** la **optimización**.

Esta es **otra** de las **razones** por las que se **desarrolló** la **forma dual** del **problema de optimización**.

M Á Q U I N A D E S V A L G O R I T M O S M O



ASCENSO DE COORDENADAS

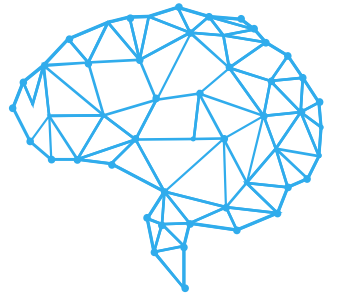
Supongamos que se **tiene** el siguiente **problema** de **optimización sin restricciones**:

$$\max_{\alpha} W(\alpha_1, \alpha_2, \dots, \alpha_m).$$

Se va a **maximizar** la **función** $W(\alpha)$ respecto a **una sola** α_i mientras que **los demás parámetros** se van a **dejar fijos**.

Se **repite** el **proceso** con **las demás** α hasta que **converja**.

M Á Q U I N A D E S V A L G O R I T M O S M O



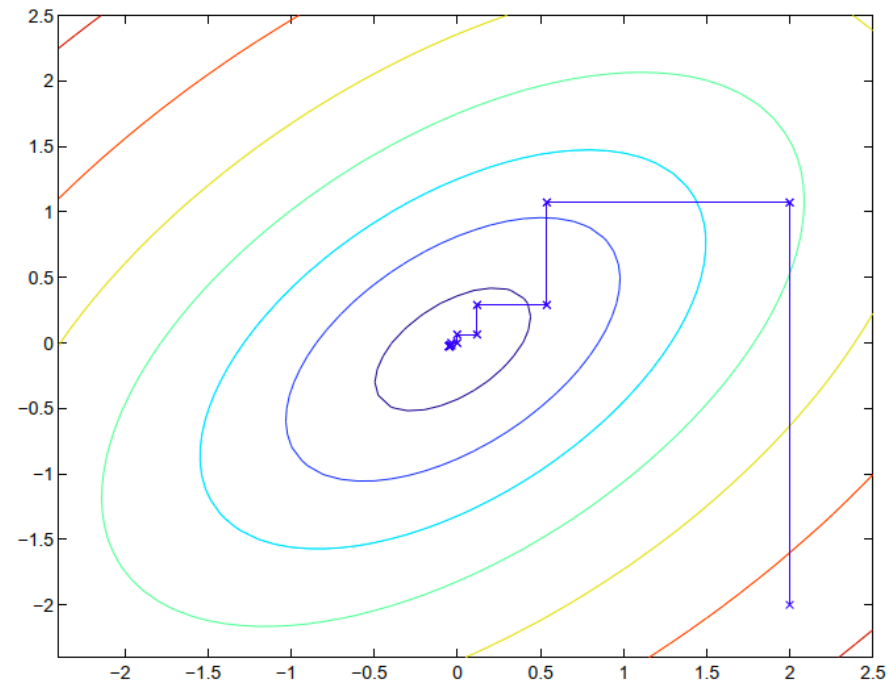
Loop until convergence: {

For $i = 1, \dots, m$, {

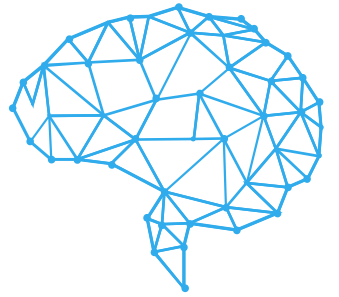
$$\alpha_i := \arg \max_{\hat{\alpha}_i} W(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_m).$$

}

}



M Á Q U I N A D E S V A L G O R I T M O S M O

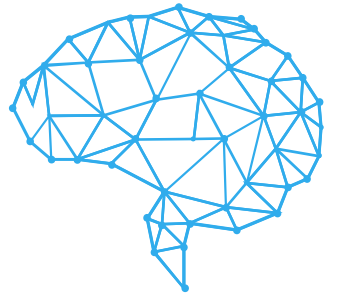


El **problema** del **algoritmo** de **ascenso** de **coordenadas** es que **no** toma en **consideración** **restricciones**. En el **problema dual** se tenía:

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0,$$

La idea del **algoritmo SMO (Sequential Minimal Optimization)** es que se **cambiarán dos** α_i y α_j **simultáneamente**.

M Á Q U I N A D E S V A L G O R I T M O S M O



Repeat till convergence {

1. Select some pair α_i and α_j to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
2. Reoptimize $W(\alpha)$ with respect to α_i and α_j , while holding all the other α_k 's ($k \neq i, j$) fixed.

}

M Á Q U I N A D E S V

A L G O R I T M O S M O



Se **seleccionan dos** α_1 y α_2 que **serán variables** mientras que **todas las demás** α serán tratadas como **constantes**. Sí esto es cierto entonces **se tiene** para la siguiente **restricción** lo siguiente:

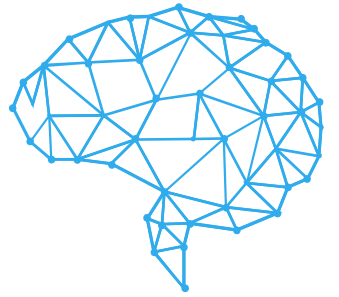
$$\sum_{i=1}^m \alpha_i y^{(i)} = 0,$$

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = - \sum_{i=3}^m \alpha_i y^{(i)}.$$

$$\alpha_1 y^{(1)} + \alpha_2 y^{(2)} = \zeta.$$

M Á Q U I N A D E S V

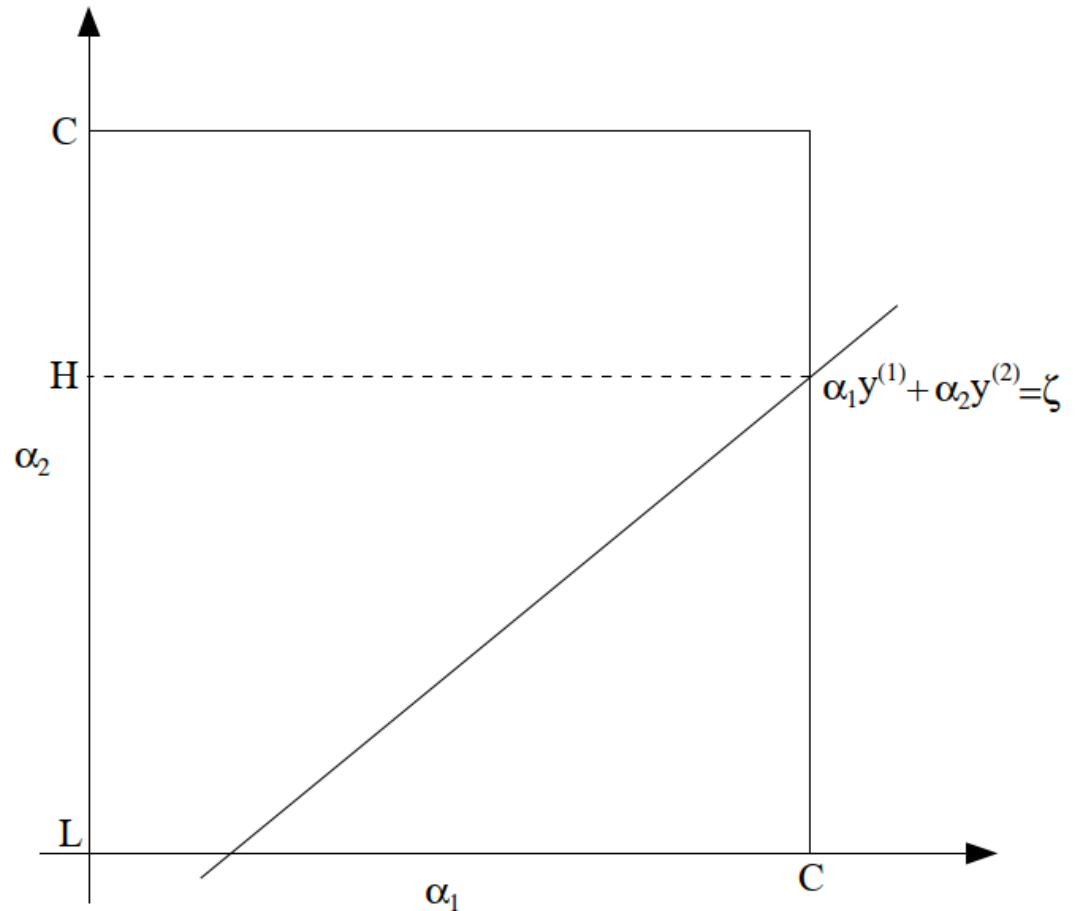
A L G O R I T M O S M O



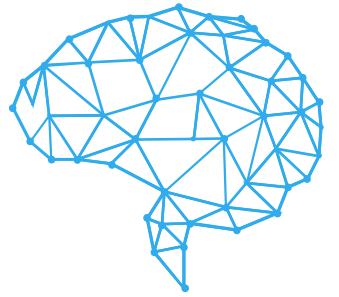
Gráficamente las dos restricciones se podría representar de la siguiente forma:

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0,$$

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$$



M Á Q U I N A D E S V A L G O R I T M O S M O



Podemos **representar** α_1 en **función** de α_2 , por lo que:

$$\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}.$$

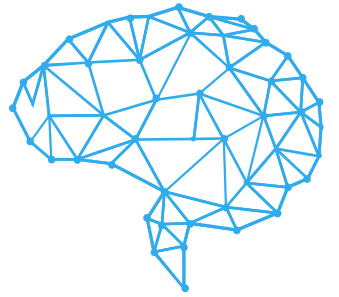
Lo que implica que **nuestro objetivo** se puede **representar como**:

$$W(\alpha_1, \alpha_2, \dots, \alpha_m) = W((\zeta - \alpha_2 y^{(2)}) y^{(1)}, \alpha_2, \dots, \alpha_m).$$

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

$$\max_{\alpha} W(\alpha) = (\xi - \alpha_2 y^{(2)}) y^{(1)} + \alpha_2 - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

M Á Q U I N A D E S V A L G O R I T M O S M O



Podemos **representar** α_1 en **función** de α_2 , por lo que:

$$\alpha_1 = (\zeta - \alpha_2 y^{(2)}) y^{(1)}.$$

Lo que implica que **nuestro objetivo** se puede **representar como**:

$$W(\alpha_1, \alpha_2, \dots, \alpha_m) = W((\zeta - \alpha_2 y^{(2)}) y^{(1)}, \alpha_2, \dots, \alpha_m).$$

$$W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

$$W(\alpha) = (\xi - \alpha_2 y^{(2)}) y^{(1)} + \alpha_2 - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

M Á Q U I N A D E S V A L G O R I T M O S M O



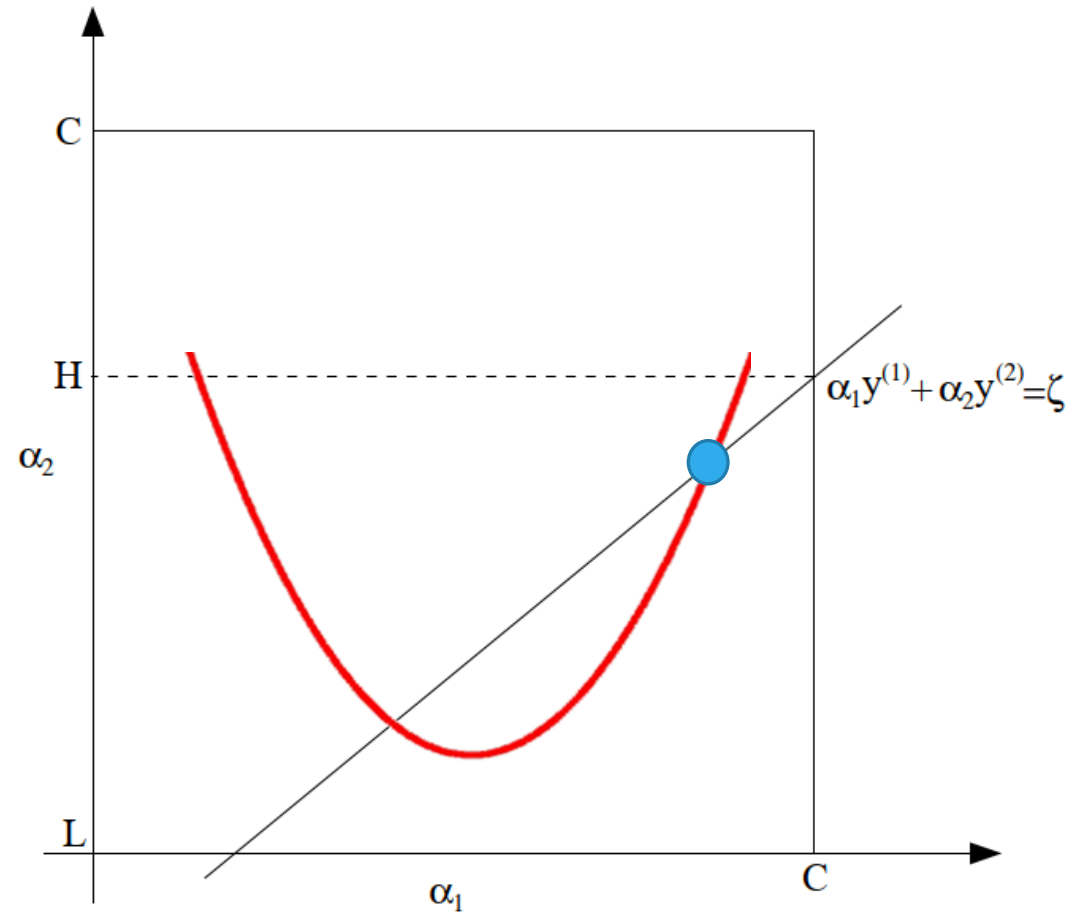
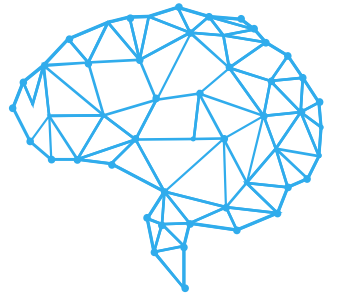
Vemos que $W(\alpha)$ es una **función cuadrática** en α_2

$$W(\alpha) = (\xi - \alpha_2 y^{(2)}) y^{(1)} + \alpha_2 - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

$$W(\alpha) = a\alpha_2^2 + b\alpha_2 + c$$

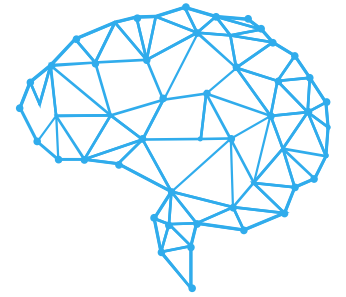
Es una **función fácil de optimizar**.

M Á Q U I N A D E S V A L G O R I T M O S M O



Se **sustituye** el **valor óptimo** de α_2 en la **recta** para obtener α_1

M Á Q U I N A D E S V A L G O R I T M O S M O



Vemos que $W(\alpha)$ es una **función cuadrática** en α_2

$$W(\alpha) = (\xi - \alpha_2 y^{(2)}) y^{(1)} + \alpha_2 - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

$$W(\alpha) = a\alpha_2^2 + b\alpha_2 + c$$

Es una **función fácil de optimizar**. Una vez optimizada, el **último paso** sería:



AI

K VECINOS

K

F

A

V

A

E

C

L

I

N

T

O

S

A





AI

ÁRBOLES DE DECISIÓN

ÁRBOLES DE DECISIÓN

F A L T A



A