

# SVM (clasificación-regresión)

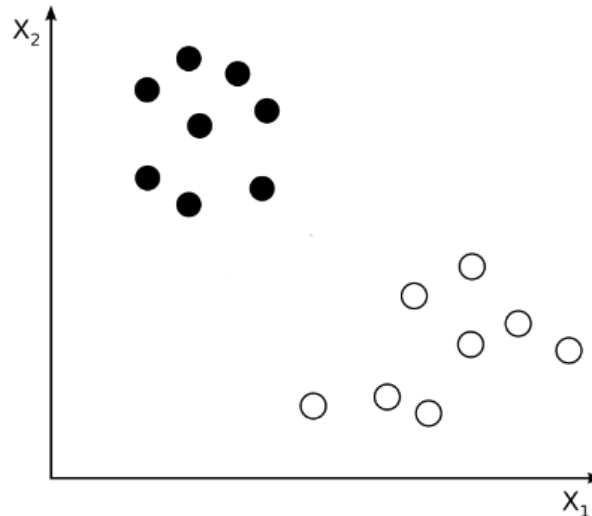
PhD(e). Jonnatan Arias Garcia – [jonnatan.arias@utp.edu.co](mailto:jonnatan.arias@utp.edu.co) –  
[jariasg@uniquindio.edu.co](mailto:jariasg@uniquindio.edu.co)

PhD. David Cardenas peña - [dcardenasp@utp.edu.co](mailto:dcardenasp@utp.edu.co)

PhD. Hernán Felipe Garcia - [hernanf.garcia@udea.edu.co](mailto:hernanf.garcia@udea.edu.co)

# Support Vector Machines

- (*Máquinas de Vectores de Soporte*) son un conjunto de versátiles y potentes algoritmos desarrollados en los laboratorios AT&T, útiles tanto en escenarios de clasificación, clúster y regresión.
- La idea detrás de SVM Partimos de un conjunto de puntos en un plano que pertenecen a dos clases distintas:



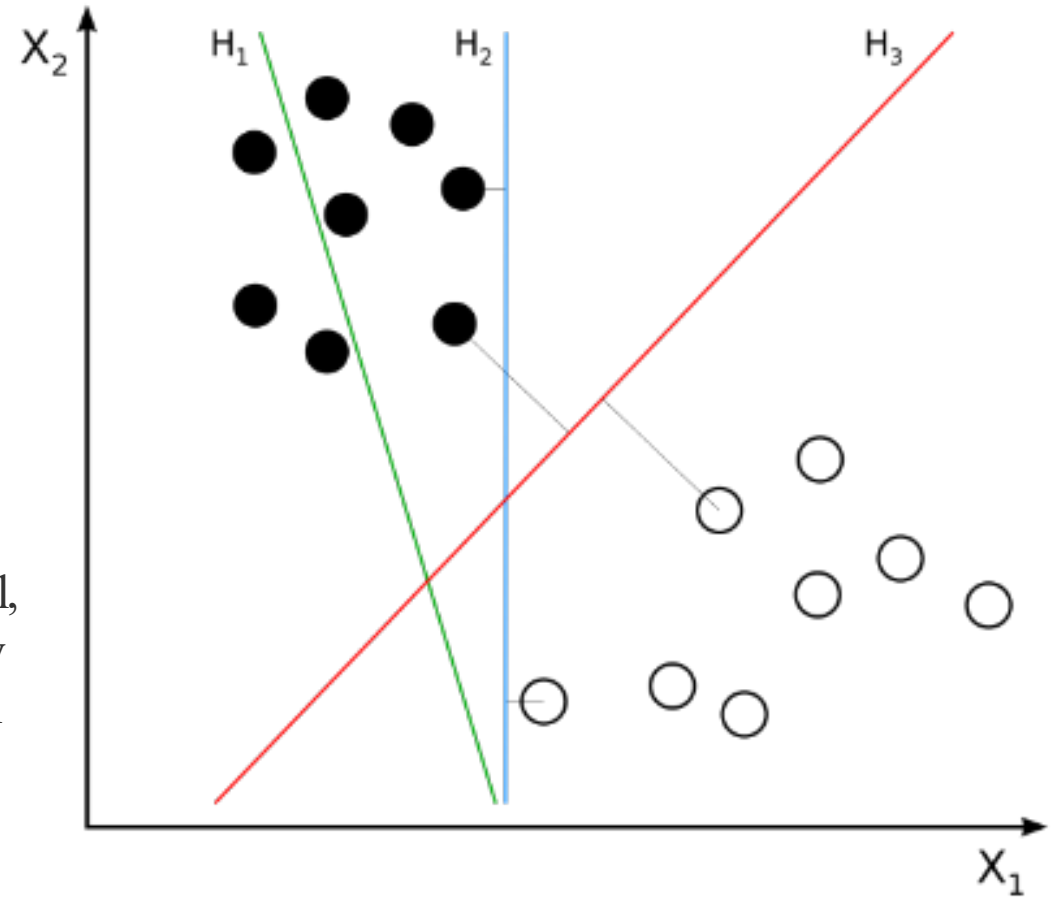
# Support Vector Machines

El objetivo es encontrar una recta (un hiperplano, en general) que permita separar ambos bloques de puntos. En el caso mostrado, existe un número infinito de rectas candidatas a resolver el problema:

La recta  $H_1$  no divide correctamente los dos bloques

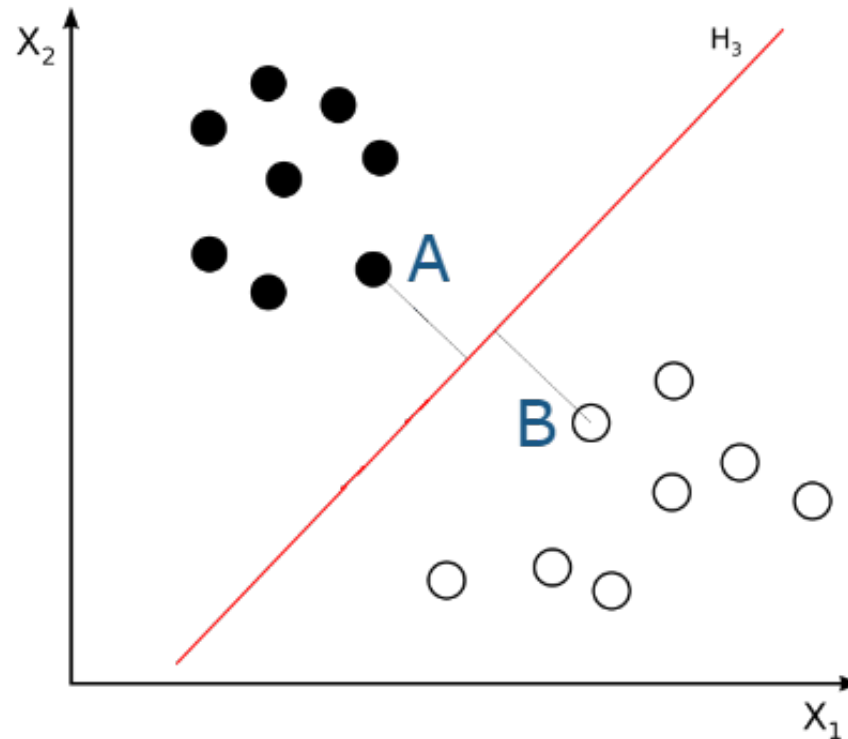
La recta  $H_2$  sí lo hace, pero su proximidad a los puntos negros hace que difícilmente pueda generalizarse el resultado (puede existir en el conjunto de puntos sobre los que realizar la predicción alguno que quede cerca de los puntos negros)

La recta  $H_3$ , desde cierto punto de vista, es la recta ideal, pues maximiza la distancia mínima a los puntos negros y blancos, optimizando la capacidad de generalización del algoritmo:



# Support Vector Machines

Los puntos A y B, los más próximos a la recta  $H_3$ , son los que determinan la posición de la recta. Si esta recta existe (si los puntos son linealmente separables) se denomina *hiperplano de máximo margen* (*maximum-margin hyperplane*), y los puntos A y B se denominan *vectores de soporte* (*support vectors*). La suma de las distancias que separan los puntos A y B del hiperplano de máximo margen se denomina *margen*.

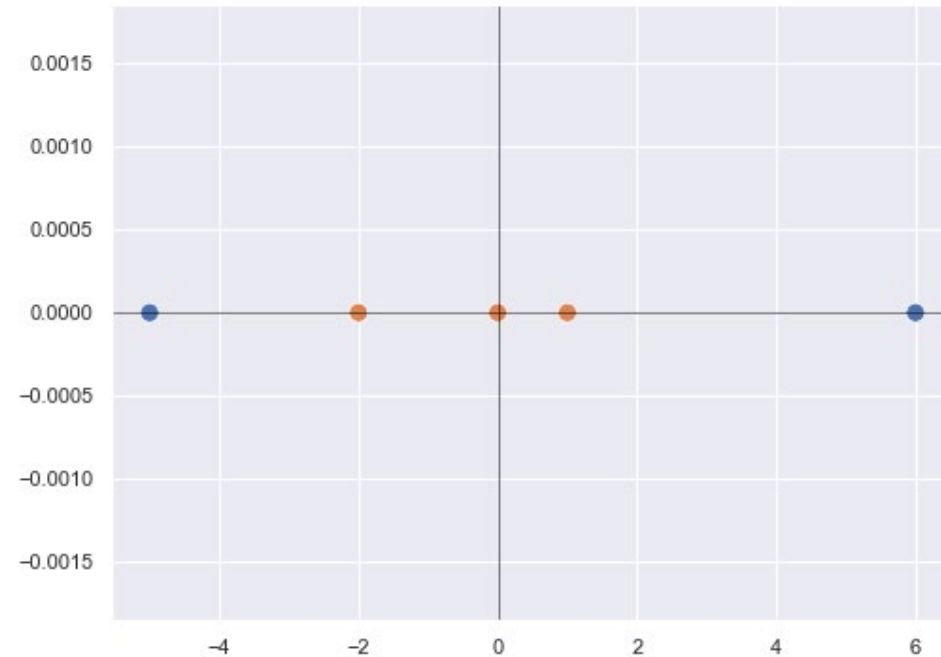


# El kernel Trick

No siempre van a ser linealmente separables.

Ejemplo, Existe un plano que separe puntos azules de naranjas?

Solución?



# El kernel Trick

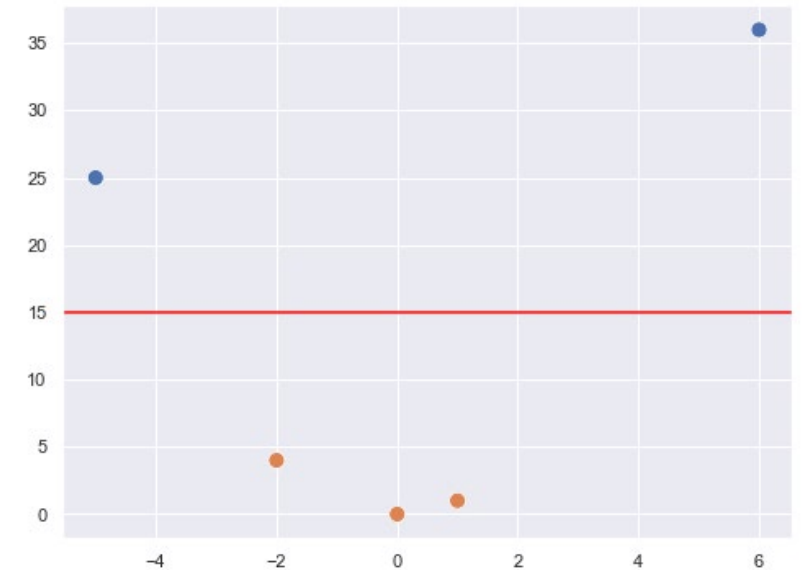
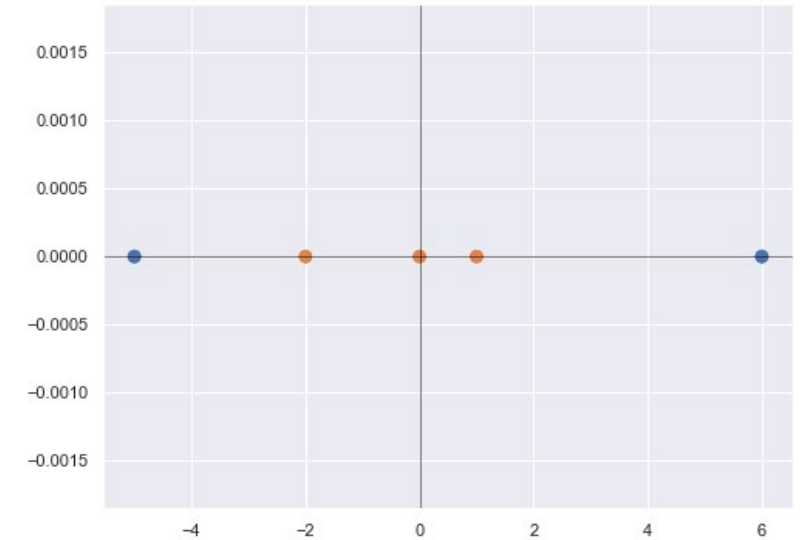
Solución?

Transformar los puntos a un espacio de mayor dimensionalidad donde sean separables.

Por ejemplo, pasar los punto (dim=1) a un plano (dim=2)

$X$  a  $XX$

$-5, -2, 0, 1, 6 \rightarrow 25, 4, 0, 1, 36$



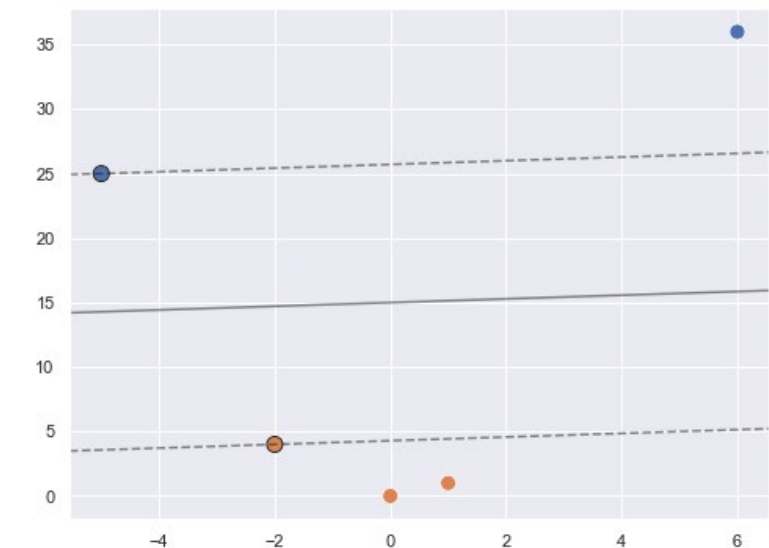
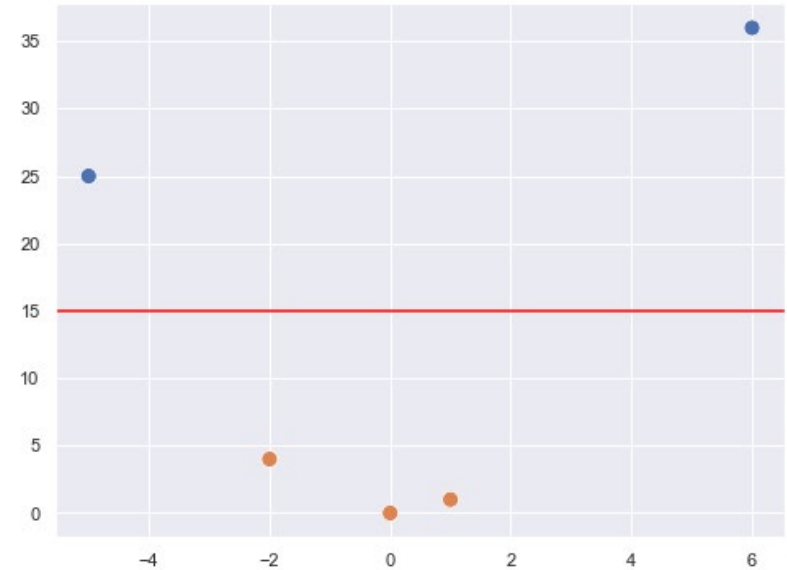
# El kernel Trick

En el nuevo espacio, es posible crear una recta que separe ambos grupos de puntos, por el ejemplo la recta (hay un infinito número de ellas, de hecho).

**SVM** buscaría la recta que maximizase la distancia a ambos grupos de puntos en este nuevo espacio.

Podemos distinguir el hiperplano de máximo margen ?  
Y los vectores de soporte ?

Esta transformación de las muestras desde el espacio original hasta un espacio de mayor dimensionalidad es llamada "**kernel trick** "



# Márgenes duros y blandos

- Si, tras aplicar el Kernel Trick los puntos que queremos clasificar son linealmente separables. Sin embargo, puede ocurrir que las clases no sean linealmente separables, no será posible encontrar el hiperplano en cuestión.
- “*Margen duro*” hace referencia al escenario en el que no se permiten errores en el entrenamiento: si se encuentra un hiperplano de máximo margen es porque clasifica correctamente todas las muestras. Si las clases no son linealmente separables, resulta más práctico permitir ciertos errores en la clasificación a cambio de poder seguir encontrando el hiperplano de máximo margen. Este segundo enfoque es el que denominamos “*de margen blando*”.
- **Margen Blando** -determinado por la función de coste del algoritmo- está regulado con el parámetro  $C$ :
  - Un valor mayor de  $C$  implica un coste mayor derivado de las muestras mal clasificadas.
  - Un valor menor implica que las muestras mal clasificadas van a suponer un coste menor, por lo que se tiende a un escenario en el que se permite un mayor número de errores.



# Kernels

Función de transformación que permite aumentar dimensionalidad y dar posiblemente un plus en la separación de datos.

Un kernel viene determinado por una matriz cuadrada de dimensión igual al número de muestras siendo analizadas, y el contenido de esta matriz es la resultante de aplicar la función de mapeo al conjunto de entrenamiento.

Linear

La función aplicada es:

$$K(x_i, x_j) = x_i^T \cdot x_j$$

# Kernels

## Linear

La función aplicada es:

$$K(x_i, x_j) = x_i^T \cdot x_j$$

## Polynomial

En el kernel polinómico la función aplicada es:

$$K(x_i, x_j) = (\gamma x_i^T \cdot x_j + r)^d$$

...donde  $\gamma$  viene dada por el parámetro **gamma** de la función SVC,  $r$  por **coef0** y  $d$  por el parámetro **degree**.

# Kernels

## Radial Basis Function

La función aplicada es:

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

$\|X_i - X_j\|$  representa una matriz cuadrada con la distancia euclídea entre cada par de puntos de X. Esta matriz puede conseguirse con la función [sklearn.metrics.pairwise.euclidean\\_distances](#). Al igual que en el caso del kernel polinómico, y viene dado por el parámetro **gamma** de SVC.

## Sigmoid

En este último caso la función aplicada es:

$$K(x_i, x_j) = \tanh(\gamma x_i^T \cdot x_j + r)$$

...donde  $\gamma$  y  $r$  vienen dados por los parámetros **gamma** y **coef0**, respectivamente.

# Matemática del método

# Máxima margen

Modelos

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + b$$

Donde la decisión es

$$t^* = \text{sign}(y(\mathbf{x})), \mathbf{w}, b$$

La margen es el valor  $y(\mathbf{x}_n)$

La distancia de un punto a un (hiper)plano:

$$d(\mathbf{x}_n, \mathbf{w}) = \frac{t_n y_n}{||\mathbf{w}||} = \frac{t_n (\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + b)}{||\mathbf{w}||}$$

donde  $t_n$  se incluye para indicar si la muestra está bien clasificada. De los modelos lineales sabemos que una muestra bien clasificada cumple  $y(\mathbf{x}_n)t_n > 0$  porque los signos coinciden.

- Para que **el error sea mínimo**, **la margen debe ser máxima**, es decir, el punto **más cercano** a la frontera debe estar **lo más lejos posible**. El punto más cercano a la frontera (sin importar la clase) estará a distancia:

$$\min_n d(x_n, w) = \min_n \frac{t_n(w^T \varphi(x_n) + b)}{\|w\|}$$

- Como queremos que esta distancia sea lo más grande posible, una *protofunción* de costo para maximizar la margen es:

$$\max_{w, b} \min_n d(x_n, w) = \max_{w, b} \left\{ \min_n \frac{t_n(w^T \varphi(x_n) + b)}{\|w\|} \right\}$$

$$\max_{w, b} \min_n d(x_n, w) = \max_{w, b} \left\{ \frac{1}{\|w\|} \min_n t_n(w^T \varphi(x_n) + b) \right\}$$

- Como escalar  $w, b$  no afecta la distancia,  $d(x_n, kw, kb) = d(x_n, w, b)$ , facilitamos la solución del problema imponiendo una restricción para aquellas muestras que estén sobre la margen:

$$t_n y_n = t_n(w^T \varphi(x_n) + b) = 1$$

- Con esta restricción, la *protofunción* se vuelve el **problema primal de optimización**:

$$\max_{\mathbf{w}, b} \min_n d(\mathbf{x}_n, \mathbf{w}) = \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n 1 \right\} = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} = \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}^2\|$$

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

$$\text{s. t. } t_n (\mathbf{w}^\top \varphi(\mathbf{x}_n) + b) \geq 1; n = 1, \dots, N$$

- Para optimizar con restricciones usamos el Lagrangiano:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1\}$$

donde  $\mathbf{a} \in \mathbb{R}^N$  es el vector de los  $N$  multiplicadores de Lagrange  $a_n$  (uno por cada restricción).

- Para optimizar con restricciones usamos el Lagrangiano:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1\}$$

donde  $\mathbf{a} \in \mathbb{R}^N$  es el vector de los  $N$  multiplicadores de Lagrange  $a_n$  (uno por cada restricción).

- Tomamos derivada respecto al vector de pesos:

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} = \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} \mathbf{w}^\top \mathbf{w} - \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^N a_n \{t_n (\mathbf{w}^\top \boldsymbol{\varphi}_n + b) - 1\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}} \{a_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n + t_n a_n b - a_n\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n = 0$$

$$\boxed{\mathbf{w} = \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n} \quad \textcircled{\hat{L}}$$



- Tomamos derivada respecto al intercepto:

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} = \frac{\partial}{\partial b} \sum_{n=1}^N \mathbf{a}_n \{t_n (\mathbf{w}^\top \boldsymbol{\varphi}_n + b) - 1\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} = \sum_{n=1}^N \frac{\partial}{\partial b} \{ \mathbf{a}_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n + \mathbf{a}_n t_n b - \mathbf{a}_n \}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a})}{\partial b} = \sum_{n=1}^N \mathbf{a}_n t_n = 0 \quad (2)$$

- Reemplazamos estos dos resultados en el Lagrangiano:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \mathbf{a}_n \{t_n y(\mathbf{x}_n) - 1\}$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \mathbf{a}_n t_n y(\mathbf{x}_n) + \sum_{n=1}^N \mathbf{a}_n \quad (5)$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N \mathbf{a}_n t_n (\mathbf{w}^\top \boldsymbol{\varphi}_n + b) + \sum_{n=1}^N \mathbf{a}_n$$

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N \mathbf{a}_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n - \sum_{n=1}^N \mathbf{a}_n t_n b + \sum_{n=1}^N \mathbf{a}_n$$

- El tercer término se elimina usando la condición de *equilibrio*:

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} - \sum_{n=1}^N \mathbf{a}_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n + \sum_{n=1}^N \mathbf{a}_n$$

- Reemplazamos la ecuación de  $\mathbf{w}$  en términos de  $\mathbf{a}$ :

$$L(\mathbf{a}) = \frac{1}{2} \left( \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n \right)^T \left( \sum_{m=1}^N a_m t_m \boldsymbol{\varphi}_m \right) - \sum_{n=1}^N a_n t_n \left( \sum_{m=1}^N a_m t_m \boldsymbol{\varphi}_m \right)^T \boldsymbol{\varphi}_n + \sum_{n=1}^N a_n$$

$$L(\mathbf{a}) = \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n t_n t_m \boldsymbol{\varphi}_n^T \boldsymbol{\varphi}_m a_m - \sum_{n=1}^N \sum_{m=1}^N a_n t_n t_m \boldsymbol{\varphi}_n^T \boldsymbol{\varphi}_m a_m + \sum_{n=1}^N a_n$$

$$L(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n t_n t_m \boldsymbol{\varphi}_n^T \boldsymbol{\varphi}_m a_m + \sum_{n=1}^N a_n$$

- Hemos llegado al **problema dual**, restringido a las dos ecuaciones que reemplazamos en el primal:

$$L(\mathbf{a}) = -\frac{1}{2} \mathbf{a}^T (\mathbf{T} \circ \mathbf{K}) \mathbf{a} + \mathbf{1}_N^T \mathbf{a}$$

$$\text{s. t. } \mathbf{a} \geq \mathbf{0}$$

$$\mathbf{t}^T \mathbf{a} = 0$$

quadprog  
linear restrictions.

$$K_{ij} = \mathbf{x}_i^T \mathbf{x}_j ; K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$$

correspondiente a un problema de optimización cuadrática con restricciones de igualdad y desigualdad lineales donde:

- $\mathbf{K} \in \mathbb{R}^{N \times N}$  es una matriz con todos los productos internos entre parejas de muestras de la base de datos de entrenamiento en el RHKS  $K_{nm} = \boldsymbol{\varphi}_n^T \boldsymbol{\varphi}_m = k(\mathbf{x}_n, \mathbf{x}_m)$ , también conocida como *matriz Gram*.
- $\mathbf{T} \in \{+1, -1\}^{N \times N}$  es la matriz objetivo con elementos  $T_{nm} = t_n t_m$  comparando las etiquetas de las muestras.
- $\mathbf{T} \circ \mathbf{K}$  es el producto de Haddamard (elemento a elemento).  
↑
- $\mathbf{1}_N \in \mathbb{R}^N$  es un vector de unos.

- La interpretación gráfica del resultado es:

- Además, reemplazando la ecuación del vector que define el hiperplano en el modelo:

$$y(\mathbf{x}^*) = \mathbf{w}^\top \boldsymbol{\varphi}(\mathbf{x}^*) + b \quad \text{modelo primal}$$

$$y(\mathbf{x}^*) = \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n^\top \boldsymbol{\varphi}^* + b$$

$$y(\mathbf{x}^*) = \sum_{n=1}^N a_n t_n k(\mathbf{x}_n, \mathbf{x}^*) + b \quad \text{modelo dual}$$

- Según las condiciones de optimalidad:

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0 \quad \forall n$$

- Pero, sin que  $a_n$  y  $t_n y(\mathbf{x}_n) - 1$  sean 0 simultáneamente. Por lo tanto:

→ Si  $a_n = 0$ , entonces  $t_n y(\mathbf{x}_n) > 1$ , es decir,  $\mathbf{x}_n$  está fuera de la margen y NO suma para las predicciones.

→ Si  $a_n > 0$ , entonces  $t_n y(\mathbf{x}_n) = 1$ , es decir,  $\mathbf{x}_n$  está en la margen y SI suma para las predicciones.  $\mathbf{x}_n$  se vuelve un soporte para construir el vector  $\mathbf{w}$ , un vector de soporte para la máquina:

$$y(\mathbf{x}^*) = \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}^*) + b \quad \text{modelo dual final}$$

- Finalmente, el intercepto  $b$  se puede calcular sabiendo que todo vector de soporte vive sobre la margen y cumple  $y(\mathbf{x}_n)t_n = 1$ :

$$t_m \left( \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = 1 \quad \forall \mathbf{x}_m \in SV$$

$y_n = \begin{cases} 1 & \mathbf{x}_n \in SV, t_n = 1 \\ -1 & \mathbf{x}_n \in SV, t_n = -1 \end{cases}$

- Multiplicamos por  $\check{t}_m$  y sumamos para todos los  $SV$ :

$$\sum_{m \in SV} \left( \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}_m) + b \right) = \sum_{m \in SV} t_m$$

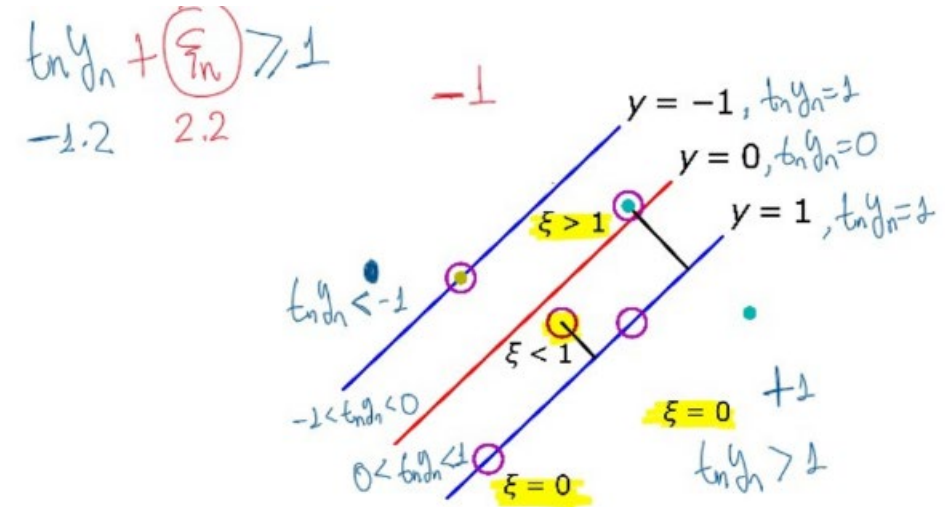
$$\sum_{m \in SV} \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}_m) + \sum_{m \in SV} b = \sum_{m \in SV} t_m$$

$$b N_{SV} = \sum_{m \in SV} \left\{ t_m - \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}_m) \right\}$$

$$b = \frac{1}{N_{SV}} \sum_{m \in SV} \left\{ t_m - \sum_{n \in SV} a_n t_n k(\mathbf{x}_n, \mathbf{x}_m) \right\}$$

# Margen Suave

- Si existe una separación lineal de las muestras en el RKHS, la SVM de margen dura la encontrará, así se vea no-lineal en el espacio de entrada.
- Sin embargo, si NO existe separación lineal en el RKHS, el problema de optimización NO tendrá solución.
- Además, si existe traslape en las distribuciones condicionales de clase en el espacio de entrada, en el RKHS se forzará la separación lineal, generando sobre-ajustes.
- Para evitar el sobre-ajuste, se permite que existan muestras dentro de la margen con una penalidad, haciendo que la **margen** sea **suave**.



- Agregamos una **variable de holgura** (slack variable),  $\xi_n \geq 0$ , para cada muestra tal que complete lo que le falta a la muestra para cumplir la restricción del problema primal de la SVM con margen dura:

$$t_n y(x_n) + \xi_n \geq 1$$

$$\xi_n \geq 0$$



- $\xi_n = 0$  Para muestras en el lado correcto de la margen.
- $0 < \xi_n < 1$  Para muestras entre la frontera y la margen de la clase correcta.
- $\xi_n = 1$  Para muestras sobre la frontera porque  $y(\mathbf{x}_n) = 0$ .
- $\xi_n > 1$  Para muestras entre la frontera y la margen de la clase incorrecta.

- Puesto que las penalizaciones deben ser las mínimas, el nuevo problema primal es:

$$\min_{\mathbf{w}, b} C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2$$

sub:  $t_n y_n + \xi_n \geq 1 \quad \forall n$   
 $\xi_n \geq 0 \quad \forall n$

donde  $C > 0$  actúa como parámetro de balance/trade-off.

- El nuevo Lagrangiano será:

$$L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu}) = C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n$$

- Tomamos derivada respecto al vector de pesos:

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial \mathbf{w}} = \mathbf{w} - \frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^N a_n \{t_n (\mathbf{w}^\top \boldsymbol{\varphi}_n + b) - 1 + \xi_n\}$$

$y_n = \mathbf{w}^\top \boldsymbol{\varphi}_n + b$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}} \{a_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n + t_n a_n b - a_n + a_n \xi_n\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n = 0$$

$$\mathbf{w} = \sum_{n=1}^N a_n t_n \boldsymbol{\varphi}_n \quad (1)$$

- Tomamos derivada respecto al intercepto:

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial b} = \frac{\partial}{\partial b} \sum_{n=1}^N a_n \{t_n (\mathbf{w}^\top \boldsymbol{\varphi}_n + b) - 1 + \xi_n\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial b} = \sum_{n=1}^N \frac{\partial}{\partial b} \{a_n t_n \mathbf{w}^\top \boldsymbol{\varphi}_n + a_n t_n b - a_n + a_n \xi_n\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\mu})}{\partial b} = \sum_{n=1}^N a_n t_n = 0 \quad (2)$$

- Tomamos derivada respecto a las holgas:

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \xi, \mu)}{\partial \xi_m}$$

$$= \frac{\partial}{\partial \xi_m} \left\{ C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n \right\}$$

$$\frac{\partial L(\mathbf{w}, b, \mathbf{a}, \xi, \mu)}{\partial \xi_m} = C - a_m - \mu_m = 0$$

$$a_m \leq C$$

- Reemplazamos estos tres resultados en el Lagrangiano:

$$L(\mathbf{w}, b, \mathbf{a}, \xi, \mu) = \sum_{n=1}^N C \xi_n - \sum_{n=1}^N \mu_n \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\}$$

$$L(\mathbf{w}, b, \mathbf{a}, \xi, \mu)$$

$$= \sum_{n=1}^N (C - \mu_n) \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n t_n y(\mathbf{x}_n) + \sum_{n=1}^N a_n - \sum_{n=1}^N a_n \xi_n$$

$$L(\mathbf{w}, b, \mathbf{a}, \xi, \mu) = \sum_{n=1}^N (C - \mu_n - a_n) \xi_n + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n t_n y(\mathbf{x}_n) + \sum_{n=1}^N a_n$$

$$L(\mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n t_n y(\mathbf{x}_n) + \sum_{n=1}^N a_n$$

$$d \{ C(\xi_1 + \xi_2) - a_1 \xi_1 - a_2 \xi_2 - \mu_1 \xi_1 - \mu_2 \xi_2 \}$$

$$= C - a_1 - \mu_1$$

- Volvemos al problema dual de la SVM con margen duro pero incluyendo una restricción más:

$$L(\mathbf{a}) = -\frac{1}{2} \mathbf{a}^T (T \circ K) \mathbf{a} + \mathbf{1}^T \mathbf{a}$$

$$\text{s.t. } 0 \leq a_n \leq C$$

$$\mathbf{t}^T \mathbf{a} = 0$$

② y ①  
quadprog

- Conclusión: Puesto que las muestras mal clasificadas tienen  $\xi_n > 1$ , la suma de las holguras es una cota superior al error de clasificación. Entonces:

- Si  $C$  aumenta, se le da más peso al término de holguras.
- Si  $C \rightarrow \infty$ , la máquina elimina por completo las holguras y vuelve a la solución de margen duro a costa de un modelo más complejo.
- Si  $C$  disminuye, las holguras no son tan importantes. Pesa más la simpleza del modelo a costa de muchos "préstamos".  $\xi_n \rightarrow 0$
- Por lo tanto,  $C$  se comporta como un parámetro de regularización inverso.