Proyecto 2 Optimización en ML – SVM

MM3025 Métodos Numéricos 2









¿Qué es SVM?





SVM

Support Vector Machines – SVM es un algoritmo de aprendizaje supervisado que se utiliza en problemas de clasificación y regresión.



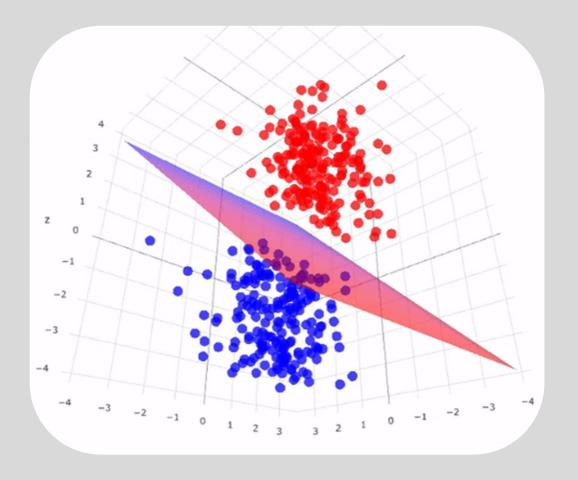
SVM

Cada objeto que se desea clasificar se representa como un punto en un espacio n-dimensional.

El objetivo del algoritmo SVM es encontrar un hiperplano que separe de la mejor forma posible dos clases diferentes de puntos.

El hiperplano es una superficie de dimensión n-1.

Nota: En 2 dimensiones es una recta y en 3 dimensiones un plano.

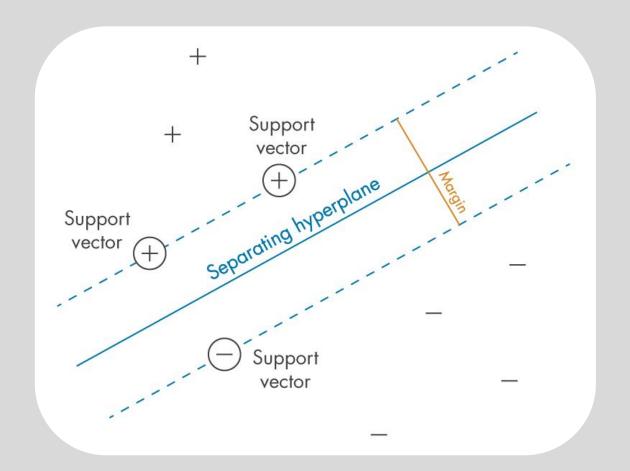




SVM

El margen es la distancia entre el hiperplano y el dato más cercano de cada clase, llamados vectores de soporte.

El algoritmo busca maximizar el margen, por lo que la optimización en SVM es esencial.



Para encontrar el hiperplano, SVM requiere un conjunto de entrenamiento que consiste de n puntos $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n)$, en donde los \mathbf{x}_i son vectores n-dimensionales & $y_i = \pm 1$ indica la clase a la que pertenece \mathbf{x}_i .

Un hiperplano se puede describir como el conjunto de puntos \mathbf{x} que satisfacen: $\mathbf{w}^T\mathbf{x} - b = 0$,

en donde w es un vector normal al hiperplano.

Recordemos que la forma normal de la ecuación de una recta en \mathbb{R}^2 es:

$$\mathbf{n}^{T}(\mathbf{x} - \mathbf{p}) = \mathbf{n}^{T}\mathbf{x} - \mathbf{n}^{T}\mathbf{p} = 0$$

Para encontrar el hiperplano, SVM requiere un conjunto de entrenamiento que consiste de n puntos $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n)$, en donde los \mathbf{x}_i son vectores n-dimensionales & $y_i = \pm 1$ indica la clase a la que pertenece \mathbf{x}_i .

Un hiperplano se puede describir como el conjunto de puntos \mathbf{x} que satisfacen: $\mathbf{w}^T\mathbf{x} - b = 0$,

en donde w es un vector normal al hiperplano.

El problema consiste en encontrar dos hiperplanos paralelos que separen las dos clases de datos, de forma que la distancia entre ellos sea la mayor posible.



Estos hiperplanos pueden describirse mediante las ecuaciones:

$$\mathbf{w}^T \mathbf{x} - b = 1$$

todo lo que esté en o por encima de este límite pertenece a la clase con la etiqueta 1

$$\mathbf{w}^T \mathbf{x} - b = -1$$

todo lo que esté en o por debajo de este límite pertenece a la clase con la etiqueta -1

Geométricamente, la distancia entre los hiperplanos es $\frac{2}{\|\mathbf{w}\|}$, así que para maximizar dicha distancia, se requiere minimizar $\|\mathbf{w}\|$.

Además, se tienen las siguientes restricciones:

$$\mathbf{w}^T \mathbf{x} - b \ge 1$$
, si $y_i = 1 \& \mathbf{w}^T \mathbf{x} - b \le -1$, si $y_i = -1$

Esto último se puede reescribir como:

$$y_i(\mathbf{w}^T\mathbf{x} - b) \ge 1, \forall i = 1, 2, ..., n$$



En conclusión, el problema de optimización en forma estándar es:

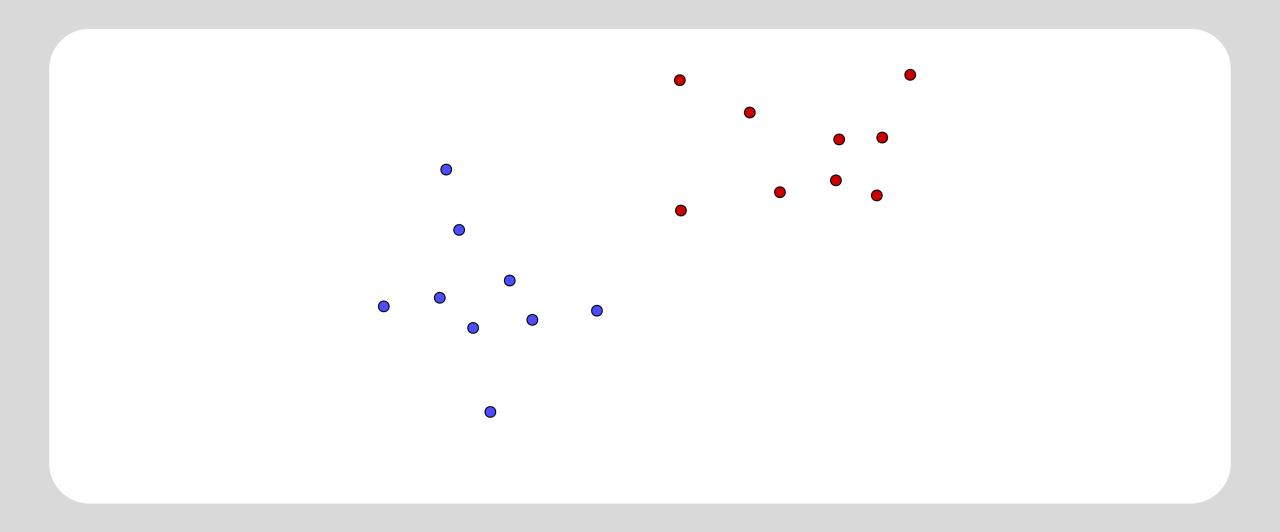
$$\min_{\mathbf{W},b} \frac{1}{2} \|\mathbf{w}\|^2 = \min_{\mathbf{W},b} \frac{1}{2} \mathbf{w}^T \mathbf{w} = \min_{\mathbf{W},b} \frac{1}{2} \sum_{i=1}^n w_i^2$$

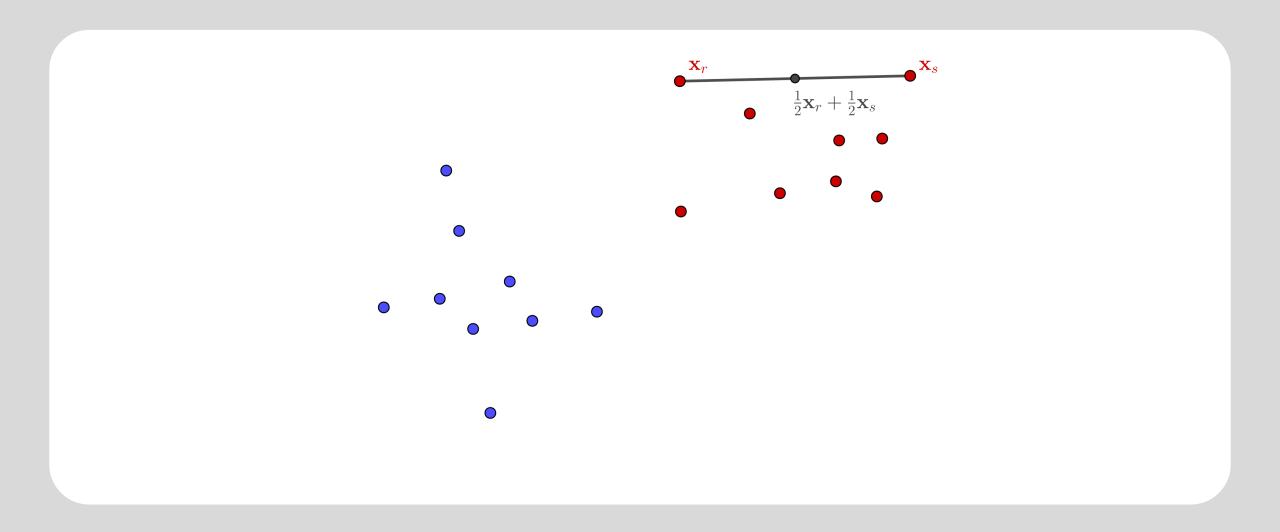
sujeto a la restricción:

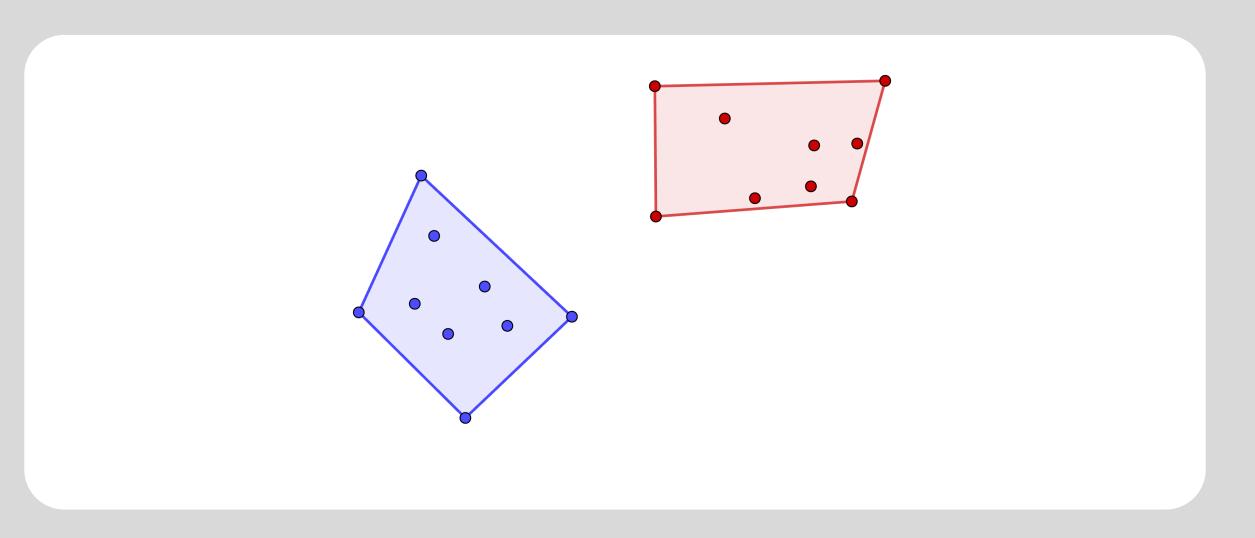
$$y_i(\mathbf{w}^T\mathbf{x} - b) \ge 1, \forall i = 1, 2, \dots, n$$

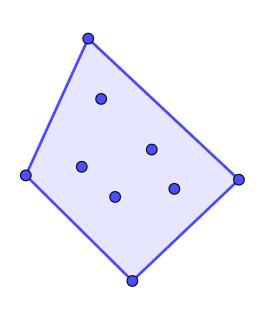
A este tipo de problemas, se les llama programación cuadrática.

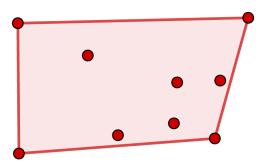






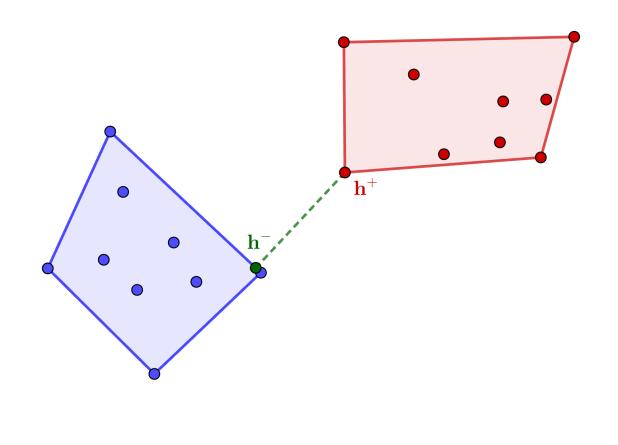




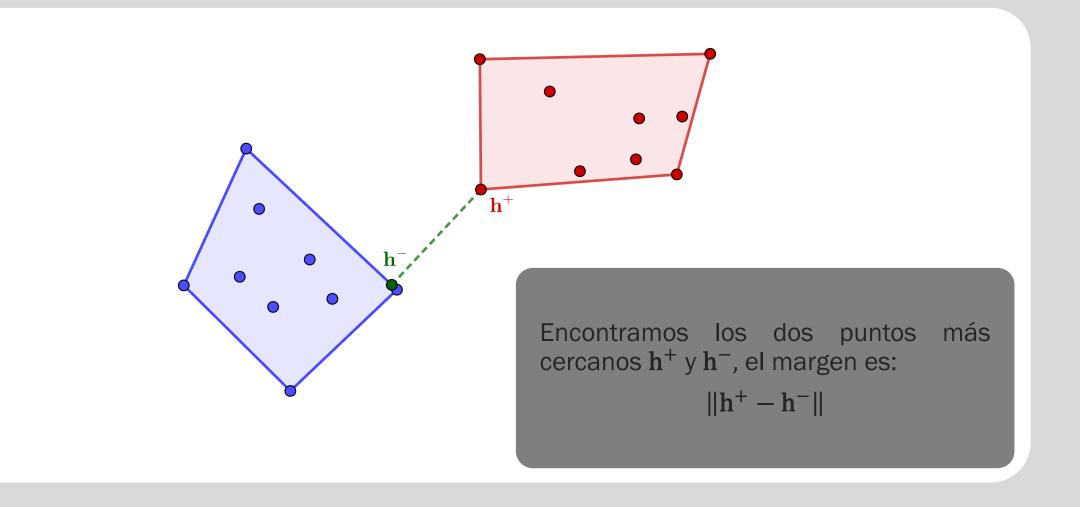


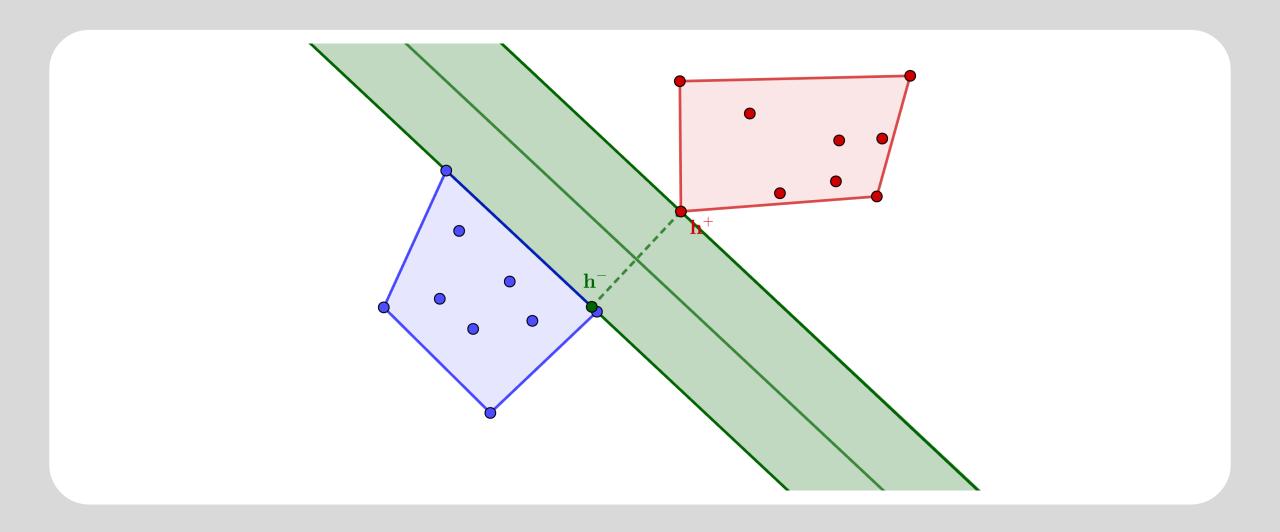
Todo punto **h** dentro del polígono es una combinación lineal de los puntos en las esquinas:

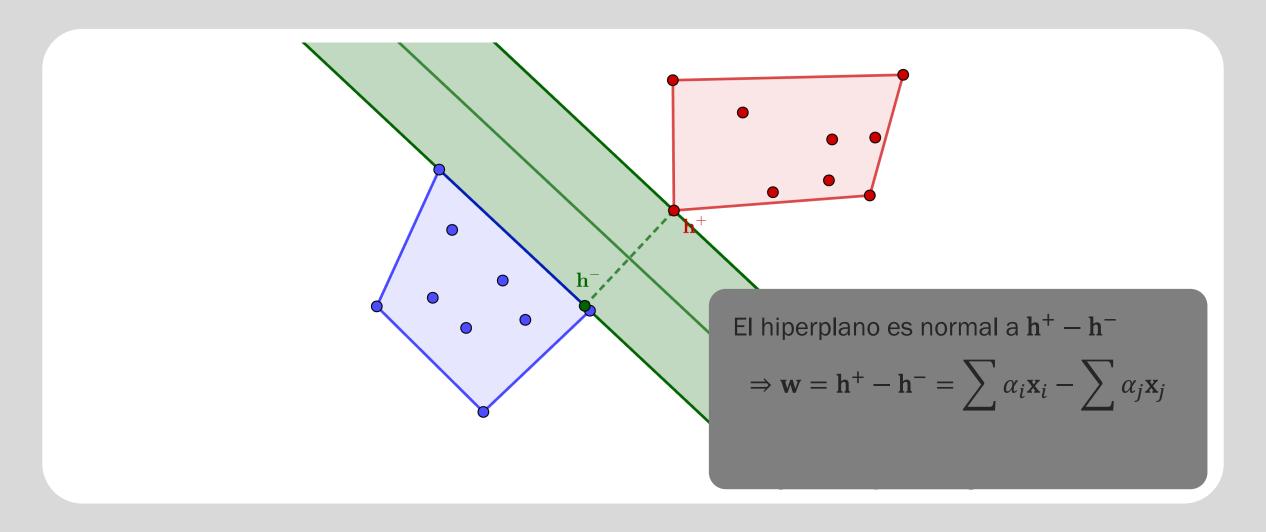
$$\mathbf{h} = \sum \alpha_i \mathbf{x}_i$$

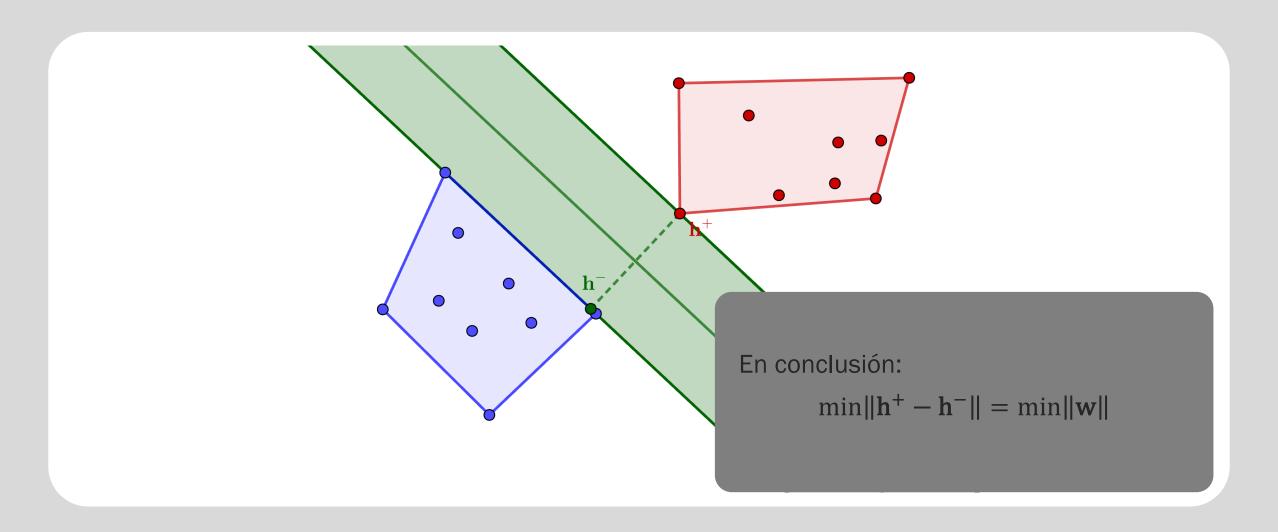














Algoritmo de optimización





Algunos algoritmos para programación cuadrática

Punto interior:

Optimiza iterativamente dentro del conjunto factible, evitando los bordes de las restricciones.

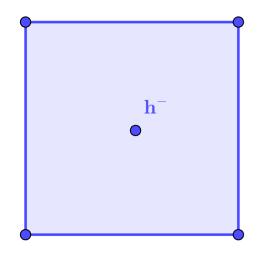
Método del Lagrangiano Aumentado:

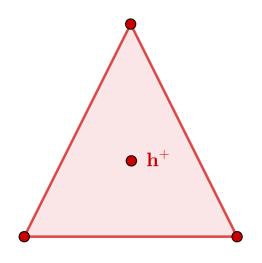
Utiliza multiplicadores de Lagrange, e introduce un término penalizador en la función objetivo para penalizar violaciones de las restricciones.

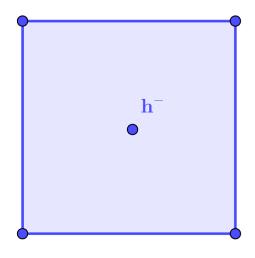


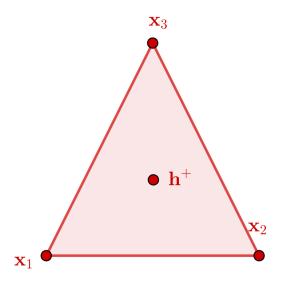
La optimización mínima secuencial – SMO es un algoritmo para resolver el problema de programación cuadrática que surge durante el entrenamiento de SVM.





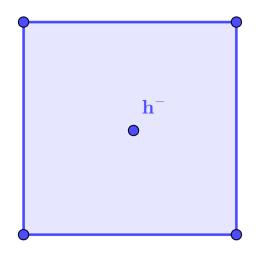


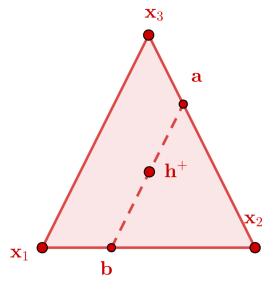




Se comienza con una elección de $\alpha_i=\frac{1}{3}$: $\mathbf{h}^+=\frac{1}{3}\mathbf{x}_1+\frac{1}{3}\mathbf{x}_2+\frac{1}{3}\mathbf{x}_3$

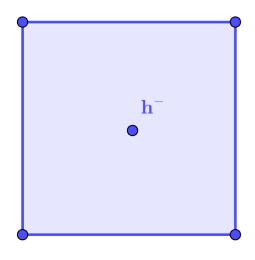


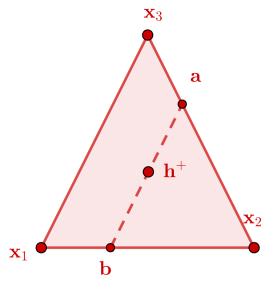




Se escogen dos α , por ejemplo, α_1 y α_3 , y se fija α_2 .

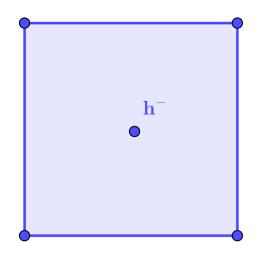


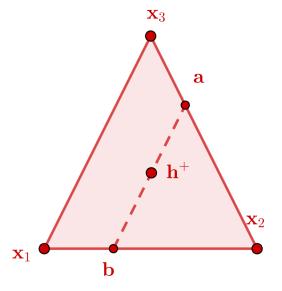




Entonces todos los puntos sobre el segmento que une **a** y **b**:

$$\mathbf{x} = \mathbf{b} + t(\mathbf{a} - \mathbf{b})$$





Se calcula el t que minimiza:

$$\|\mathbf{b} + t(\mathbf{a} - \mathbf{b}) - \mathbf{h}^-\|$$



Como min $\|\mathbf{x}\|=\min\frac{1}{2}\|\mathbf{x}\|^2$, se hace $\Delta=\frac{1}{2}\sum_i (b_i+t(a_i-b_i)-h_i^-)^2$, y entonces:

$$\frac{\partial \Delta}{\partial t} = \sum_{i} (b_i + t(a_i - b_i) - h_i^-)(a_i - b_i) = \sum_{i} (t(a_i - b_i)^2 + (b_i - h_i^-)(a_i - b_i))$$

Al hacer $\frac{\partial \Delta}{\partial t} = 0$ y despejar para t se obtiene:

$$t = \frac{(a - b)^{T}(h^{-} - b)}{\|a - b\|^{2}}$$

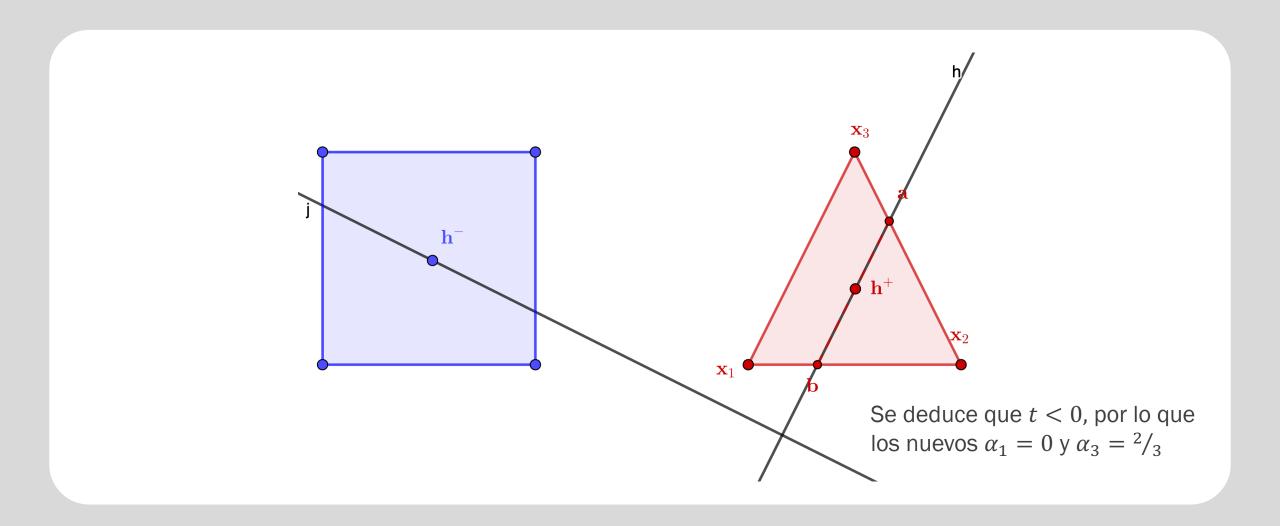


Se actualizan los α :

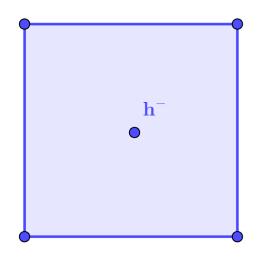
$$\alpha_1 = \begin{cases} 0, & t \le 0\\ \frac{2}{3}t, & t \in (0,1) & \alpha_3 = \frac{2}{3} - \alpha_1\\ \frac{2}{3}, & t \ge 1 \end{cases}$$

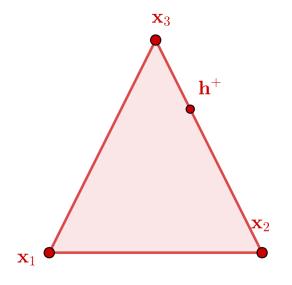
El algoritmo funciona de forma iterativa, actualizando únicamente dos variables a la vez, lo cual simplifica el problema considerablemente y reduce el tiempo de cálculo.











En conclusión:

$$h^+ = 0\mathbf{x}_1 + \frac{1}{3}\mathbf{x}_2 + \frac{2}{3}\mathbf{x}_3$$



El algoritmo SMO evita cálculos de matriz de alta dimensión, lo que reduce el costo computacional.

Debido a su enfoque en subproblemas de dos variables, el algoritmo puede implementarse sin depender de métodos complejos de optimización.

Su estructura iterativa lo hace ideal para conjuntos de datos grandes, donde los métodos tradicionales de optimización cuadrática serían demasiado costosos.

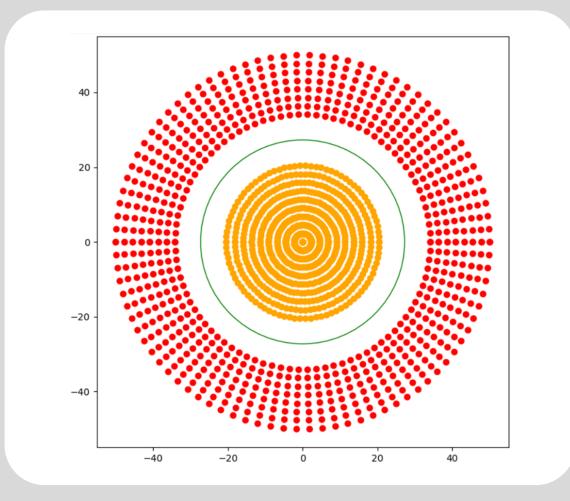




Cuando no es posible lograr la clasificación con un hiperplano lineal, los datos de entrada se transforman en un espacio de características de mayor dimensión, donde puede ser más fácil encontrar una frontera de decisión lineal que separe las clases.



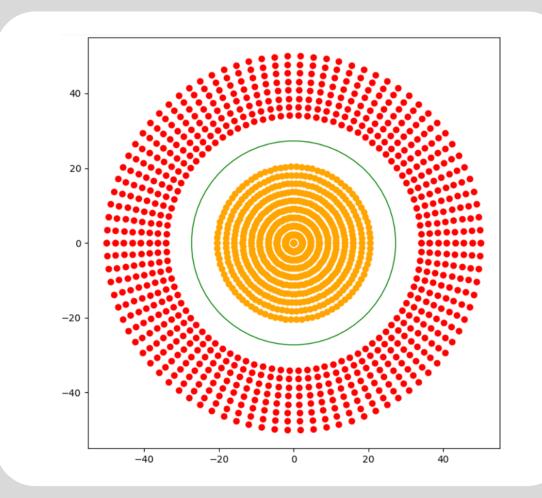
Suponga que se tienen dos características X e Y, y los datos no son linealmente clasificables, como los que se muestran en la figura.





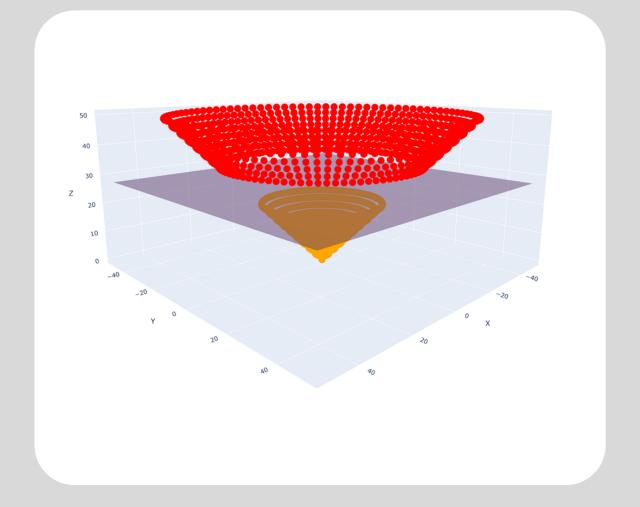
Lo que se tiene que hacer es añadir otra dimensión en la que los datos se vuelven linealmente separables.

La idea es usar una función kernel para transformar las características de los datos Z = f(X, Y).





Si $f(X,Y) = X^2 + Y^2$, entonces los datos de la forma (X,Y,Z) ahora son linealmente clasificables.



Tipos de Kernel

Las funciones de kernel asignan los datos a un espacio dimensional diferente, que suele ser superior, con el objetivo de que resulte más fácil separar las clases después de esta transformación. Entre los más conocidos están:

Lineal	$K(x_1, x_2) = x_1^T x_2$	Datos linealmente separables
RBF	$K(x_1, x_2) = \exp\left(-\frac{\ x_1 - x_2\ ^2}{2\sigma^2}\right)$	Datos no lineales con relaciones complejas.
Polinómica	$K(x_1, x_2) = (x_1^T x_2 + 1)^p$	Datos con relaciones polinómicas.
Sigmoide	$K(x_1, x_2) = \tanh(\beta_0 x_1^T x_2 + \beta_1)$	Datos que modelan comportamientos similares a redes neuronales o funciones sigmoides.



¿Preguntas?



