

Máquinas con vectores de Soporte - SVM

Misael López Ramírez

Conceptos matemáticos

- Las SVM son clasificadores derivados de la teoría de aprendizaje estadístico postulada por Vapnik y Chervonenkis.
- Las SVM fueron presentadas en 1992 y adquirieron fama cuando dieron resultados muy superiores a las redes neuronales en el reconocimiento de letra manuscrita, usando como entrada pixeles.
- Pretenden predecir a partir de lo ya conocido.

Conceptos matemáticos

Hay **l** observaciones y cada una consiste en un par de datos:

un vector  $x_i \in R^n, i = 1, \dots, l$

una etiqueta  $y_i \in \{+1, -1\}$

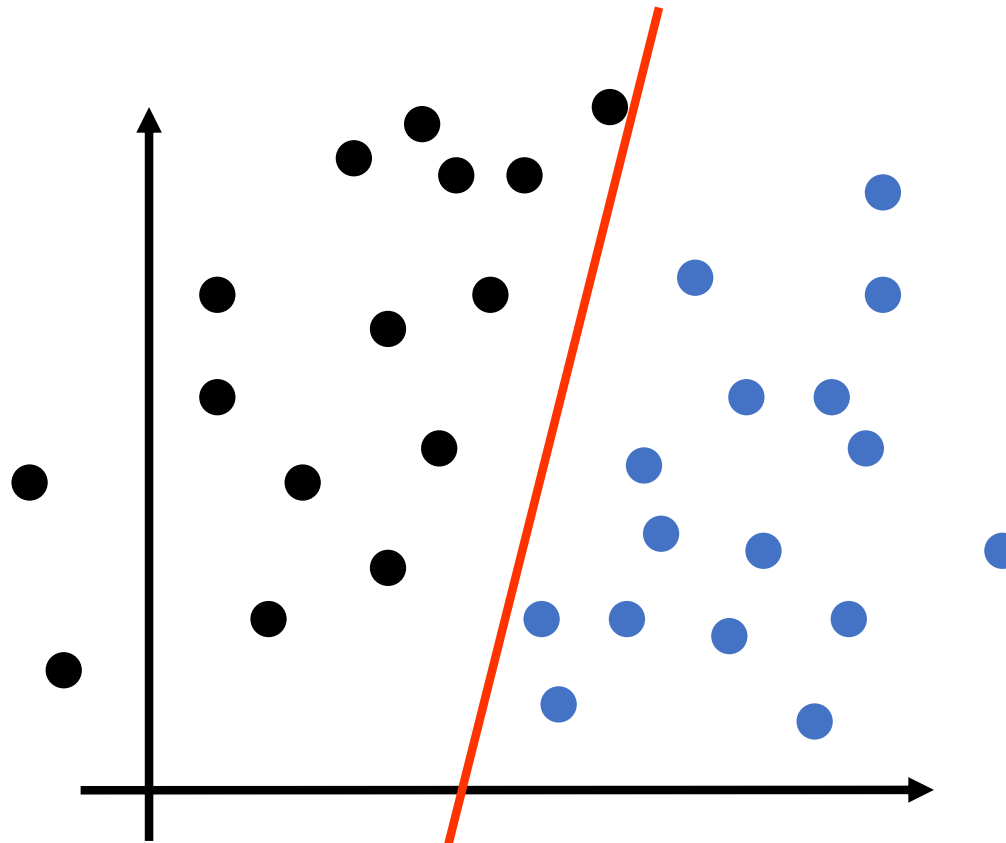
Supóngase que se tiene un **hiperplano** que separa las muestras positivas (+1) de las negativas (-1). Los puntos \mathbf{x}_i que están en el hiperplano satisfacen $\mathbf{w} \cdot \mathbf{x} + b = 0$.

 Ecuación de la línea recta

Idea inicial de separación

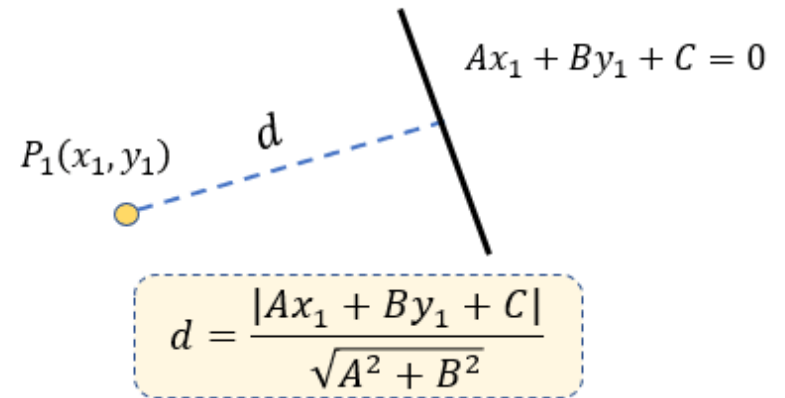
● +1

● -1

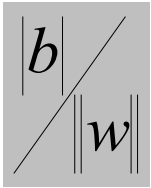


$$w \cdot x + b = 0$$

Conceptos matemáticos



w es normal al hiperplano.



es la distancia perpendicular del hiperplano al origen.



es la norma euclídea de **w**

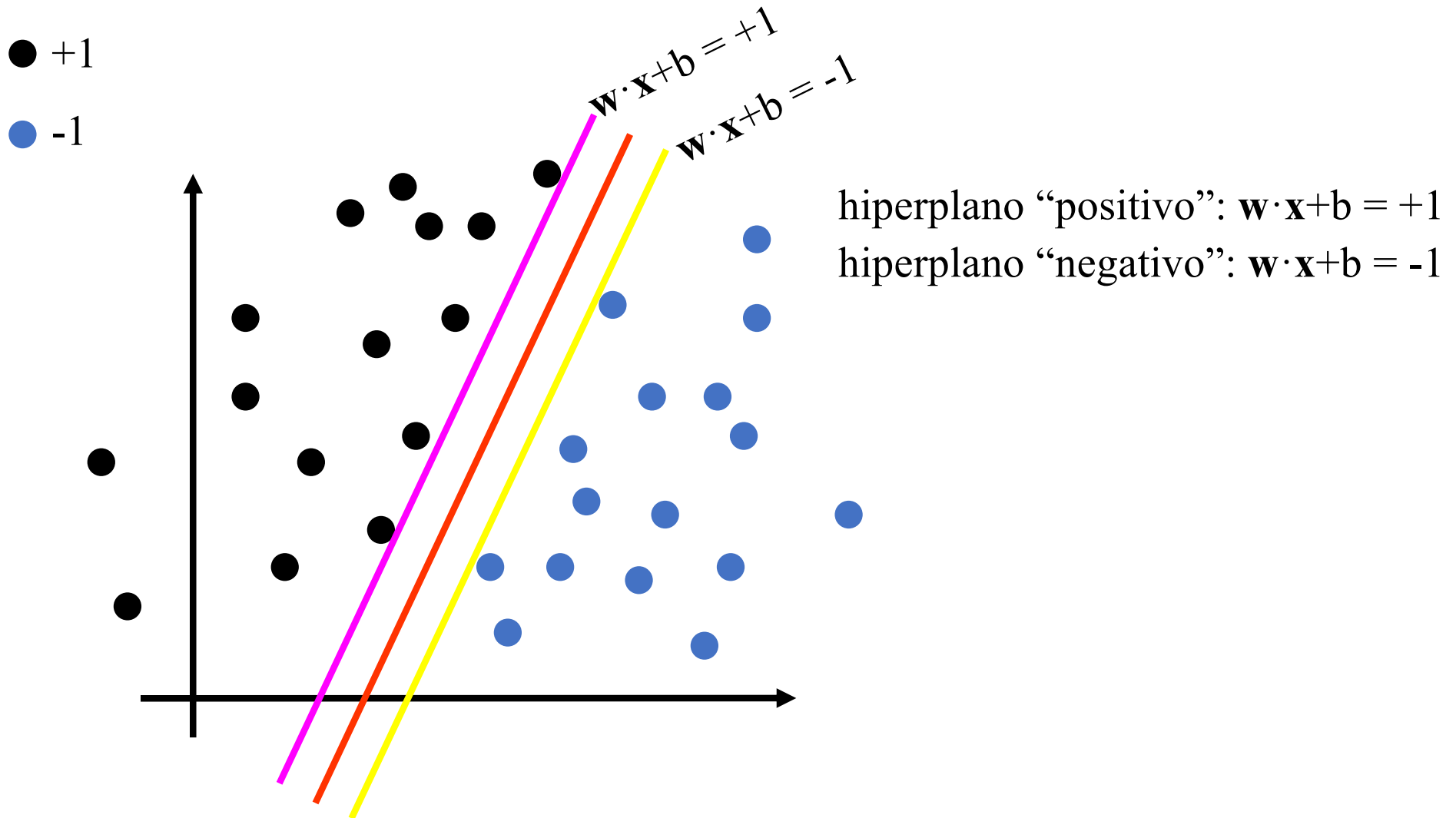
Lo que se quiere es separar los puntos de acuerdo al valor de su etiqueta y_i en dos hiperplanos diferentes:

$\mathbf{w} \cdot \mathbf{x}_i + b \geq +1$ para $y_i = +1$. (hiperplano “positivo”)

$\mathbf{w} \cdot \mathbf{x} + b \leq -1$ para $y_i = -1$ (hiperplano “negativo”)

Simplificando: $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0$

Idea inicial de separación



Conceptos matemáticos

Sea d_+ (d_-) la distancia más corta entre el hiperplano positivo (negativo) y el punto positivo (negativo) más cercano.

Sea el “margen” la distancia entre los hiperplanos “positivo” y “negativo”. El margen es igual a: $\frac{2}{\|W\|}$

La idea es encontrar un hiperplano con el máximo “margen”. Esto es un problema de optimización:

maximizar: $\frac{2}{\|W\|}$ sujeito a : $y_i(W \cdot X_i + b) \geq +1$


Ó minimizar

Conceptos matemáticos

El problema su puede expresar así:

$$\text{minimizar: } \frac{1}{2} \|W\|^2 \quad \text{sujeto a : } y_i(W \cdot X_i + b) \geq 1$$

Pero el problema se puede transformar para que quede más fácil de manejar! Se usan multiplicadores de Lagrange (α_i).

$$L_P \equiv \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^1 \alpha_i y_i (\mathbf{w} \cdot \mathbf{x}_i + b) + \sum_{i=1}^1 \alpha_i$$


de observaciones ó muestras

Conceptos matemáticos

Multiplicadores de Lagrange

minimizar: $\frac{1}{2} \|W\|^2$ sujeto a : $y_i(W \cdot X_i + b) \geq 1$

Teorema de Lagrange:

$$\nabla f(x, y, z) = \lambda \nabla g(x, y, z)$$

Método de los multiplicadores de Lagrange:

$$\begin{array}{ll} \text{1) } \frac{df}{dx}(x, y, z) = \lambda \frac{dg}{dx}(x, y, z) & \text{3) } \frac{df}{dz}(x, y, z) = \lambda \frac{dg}{dz}(x, y, z) \\ \text{2) } \frac{df}{dy}(x, y, z) = \lambda \frac{dg}{dy}(x, y, z) & \text{4) } g(x, y, z) = c \end{array}$$

$$L_P = f(x, y, z) - \lambda g(x, y, z)$$

Función de lagrange

$$L_P \equiv \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^1 \alpha_i y_i (\mathbf{w} \cdot \mathbf{x}_i + b) + \sum_{i=1}^1 \alpha_i$$

Conceptos matemáticos

Haciendo que los gradientes de L_p respecto a \mathbf{w} y b sean cero, se obtienen las siguientes condiciones:

$$\mathbf{w} = \sum_{i=1}^1 \alpha_i y_i \mathbf{x}_i \quad \sum_{i=1}^1 \alpha_i y_i = 0$$

Reemplazando en L_p se obtiene el problema dual:

$$L_P \equiv \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^1 \alpha_i y_i (\mathbf{w} \cdot \mathbf{x}_i + b) + \sum_{i=1}^1 \alpha_i$$

$$L_D = \sum_{i=1}^1 \alpha_i + \sum_{i=1, j=1}^1 \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

Conceptos matemáticos

$$\min f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}$$

La forma para optimizar es:

$$s.a. \mathbf{A} \mathbf{x} \leq \mathbf{b}$$

$$\mathbf{x} \geq 0$$

maximizar:

$$L_D = \sum_{i=1}^1 \alpha_i + \sum_{i=1, i=1}^1 \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

* L_D mayores a cero son nombrados vectores de Soporte (NV).

sujeto a :

$$\mathbf{w} = \sum_{i=1}^1 \alpha_i y_i \mathbf{x}_i \quad \sum_{i=1}^1 \alpha_i y_i = 0$$

El problema se reduce a encontrar el máximo de L_D con las restricciones anteriores

*PROGRAMACIÓN CUADRÁTICA

*Se puede resolver con el método del punto interior

Conceptos matemáticos

La forma para optimizar es:

maximizar:

$$L_D = \sum_{i=1}^1 \alpha_i + \sum_{i=1, j=1}^1 \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

* L_D mayores a cero son nombrados vectores de Soporte (NV).

sujeto a :

$$\mathbf{w} = \sum_{i=1}^1 \alpha_i y_i \mathbf{x}_i \quad \sum_{i=1}^1 \alpha_i y_i = 0$$

Descripción

Solver para funciones objetivas cuadráticas con restricciones lineales.

encuentra un mínimo para un problema especificado por quadprog

$$\min_x \frac{1}{2} x^T H x + f^T x \text{ such that } \begin{cases} A \cdot x \leq b, \\ Aeq \cdot x = beq, \\ lb \leq x \leq ub. \end{cases}$$

, y son matrices, y,,, y son vectores. $H A Aeq f b beq lb ub x$

Conceptos matemáticos

Al encontrar el vector L_D máximo calculamos los valores de pendiente W y b con las formulas optimos:

$$W = \sum_{i=1}^L \alpha_i y_i x_i \quad b = \frac{1}{NV} \sum_{i=1}^{NV} (y_i - W * x_i)$$

*NV numero de vectores de soporte.

Sustituimos W y b en la ecuación de la línea recta de la forma

$$W \cdot x + b = -1$$

$$W \cdot x + b = +1$$

*Por lo tanto clasifica entre dos clases únicamente con el signo de tal forma nos queda como:

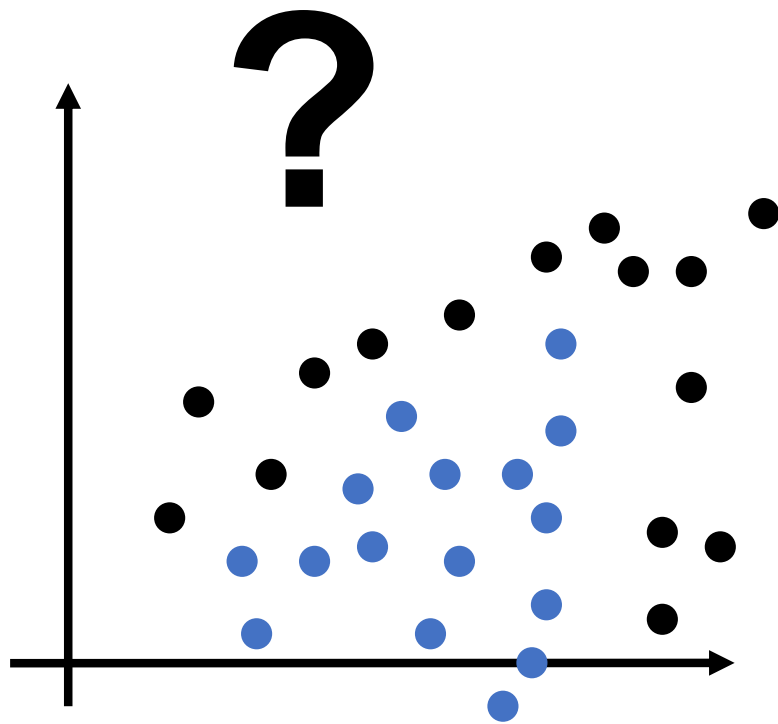
$$f(x) = \text{sign}(W \cdot x + b)$$

Conceptos matemáticos

Cuando los datos no se pueden separar linealmente se hace un cambio de espacio mediante una función que transforme los datos de manera que se puedan separar linealmente. Tal función se llama ***Kernel***.

También hay métodos para separar los datos (x_i, y_i) directamente aún no siendo separables linealmente, mediante ***funciones polinómicas*** y otro tipo de funciones, las ***Funciones de Base Radial*** (RBF).

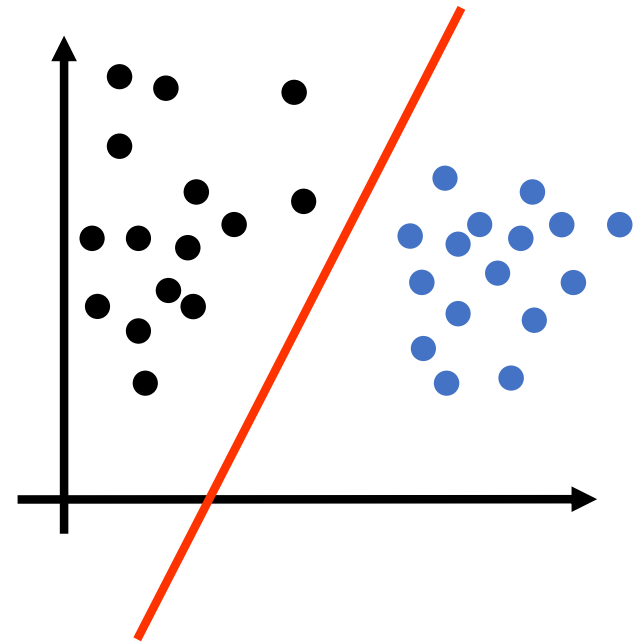
Conceptos matemáticos



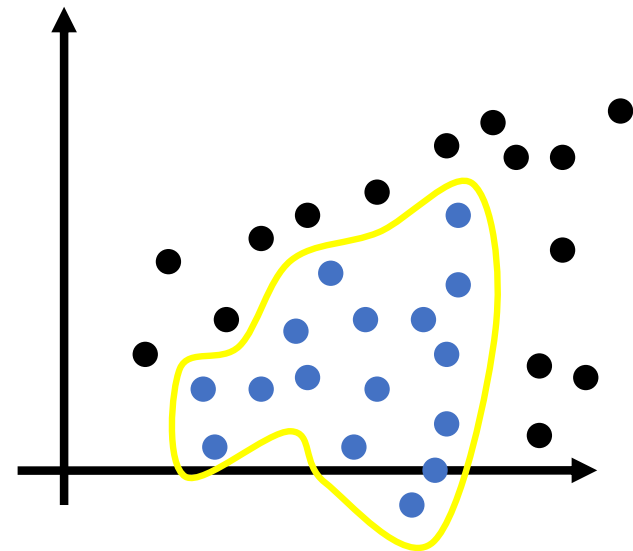
● +1

● -1

Kernel

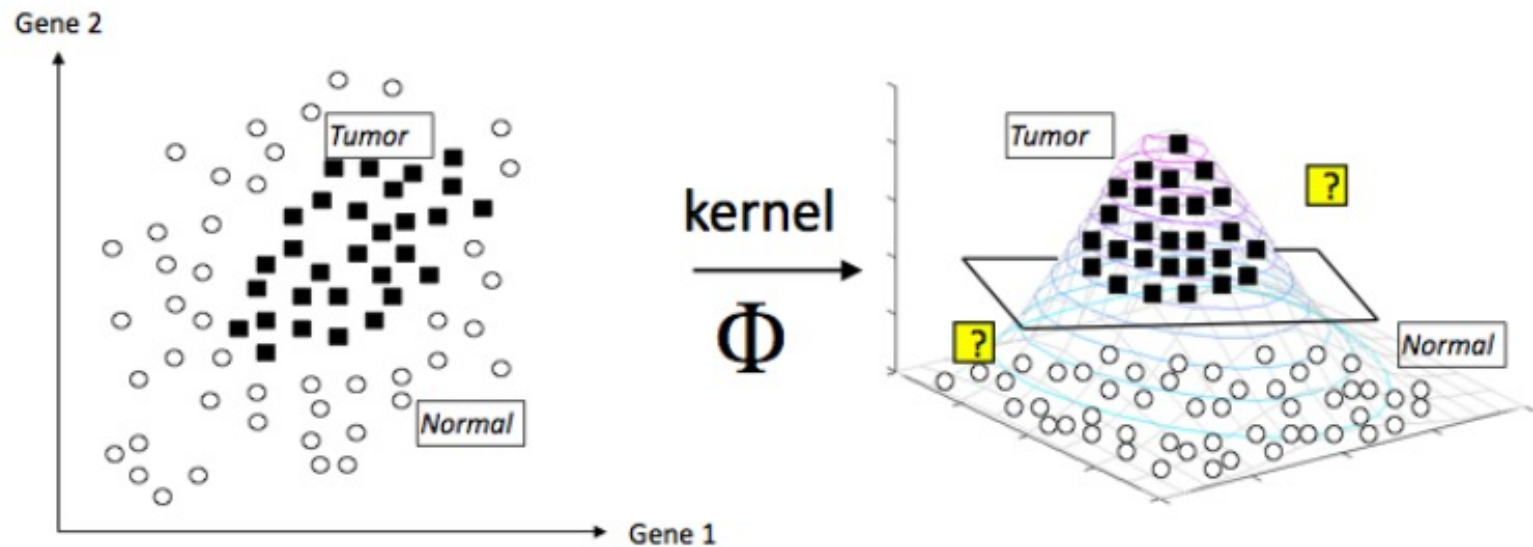


RBF



Conceptos matemáticos

La forma del kernel nos puede ayudar a obtener una mejor discriminación entre clases como lo muestra la siguiente figura:



Conceptos matemáticos

Algunos problemas con las SVM:

Overtraining: se han aprendido muy bien los datos de entrenamiento pero no se pueden clasificar bien ejemplos no vistos antes. Ej.: un botánico que conoce mucho.

La porción n de los datos no conocidos que será mal calificada, está limitada por:

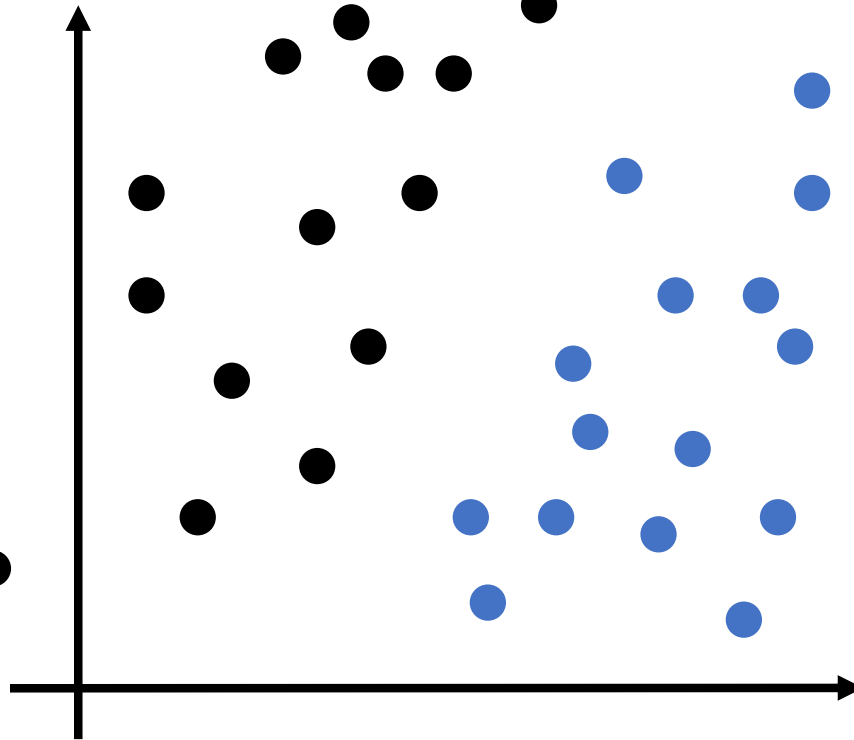
$$n = \frac{\text{No. vectores de soporte}}{\text{No. de ejemplos de entrenamiento}}$$

Se aplica el principio de Ockham.

Interpretación geométrica

● +1

● -1



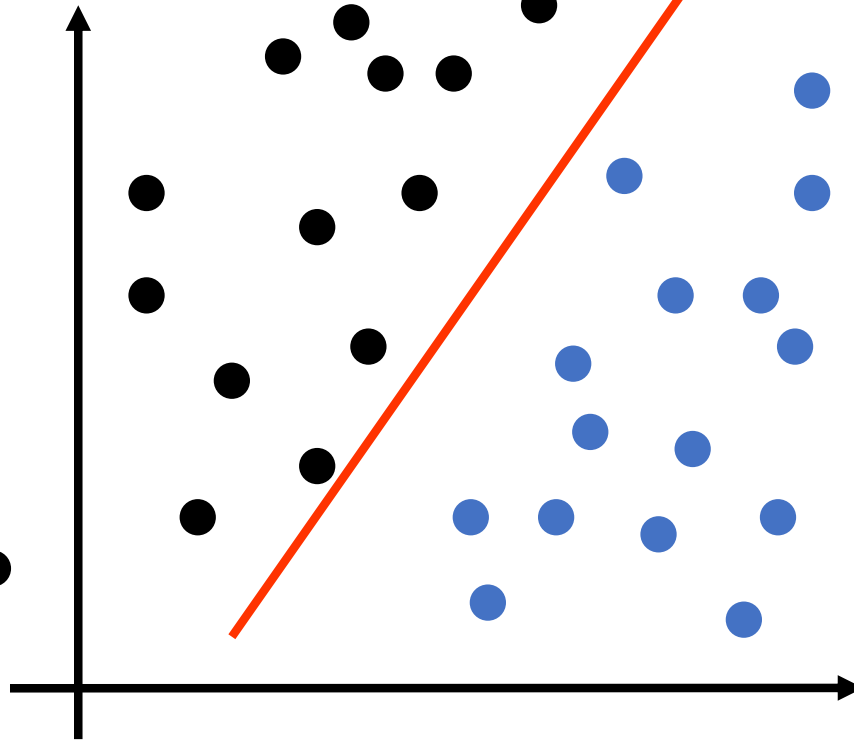
2 clases a clasificar

¿cómo clasificar estos datos?

Interpretación geométrica

● +1

● -1



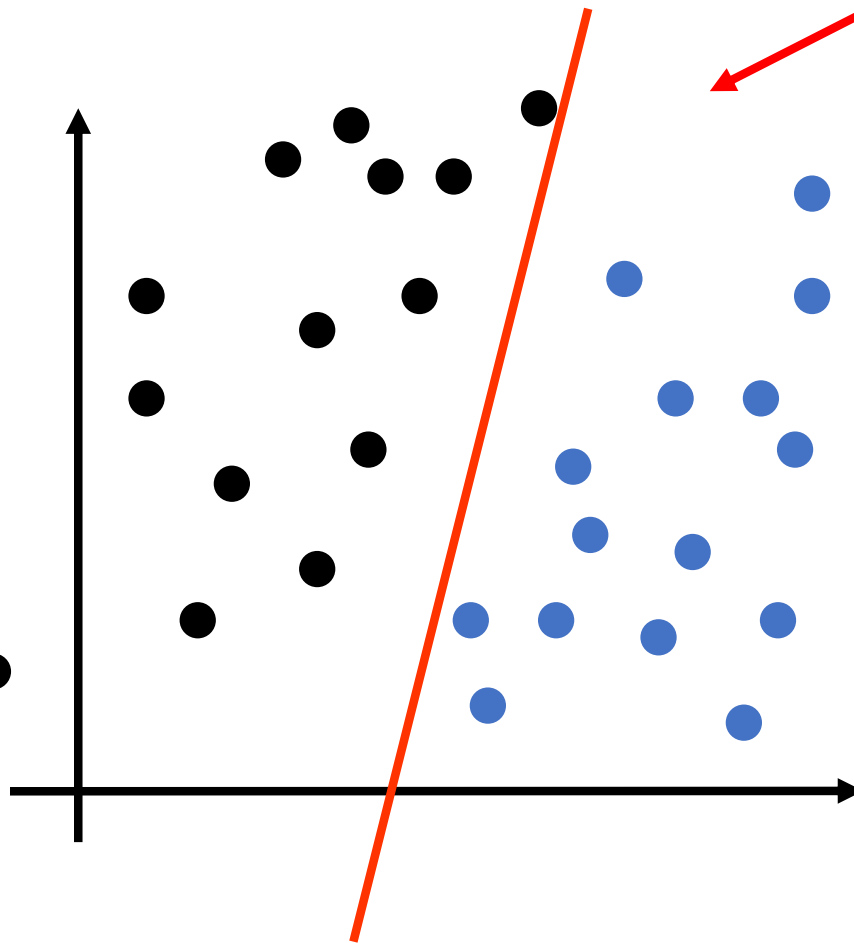
Linealmente separables

¿cómo clasificar estos datos?

Interpretación geométrica

● +1

● -1



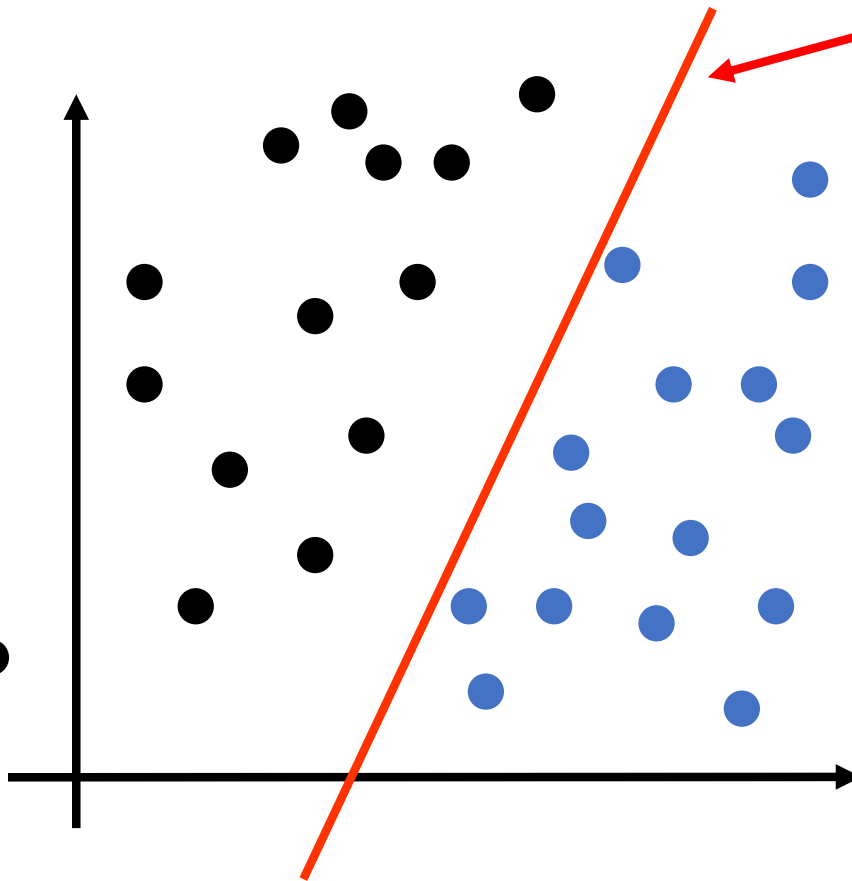
Otra forma de clasificar

¿cómo clasificar estos datos?

Interpretación geométrica

● +1

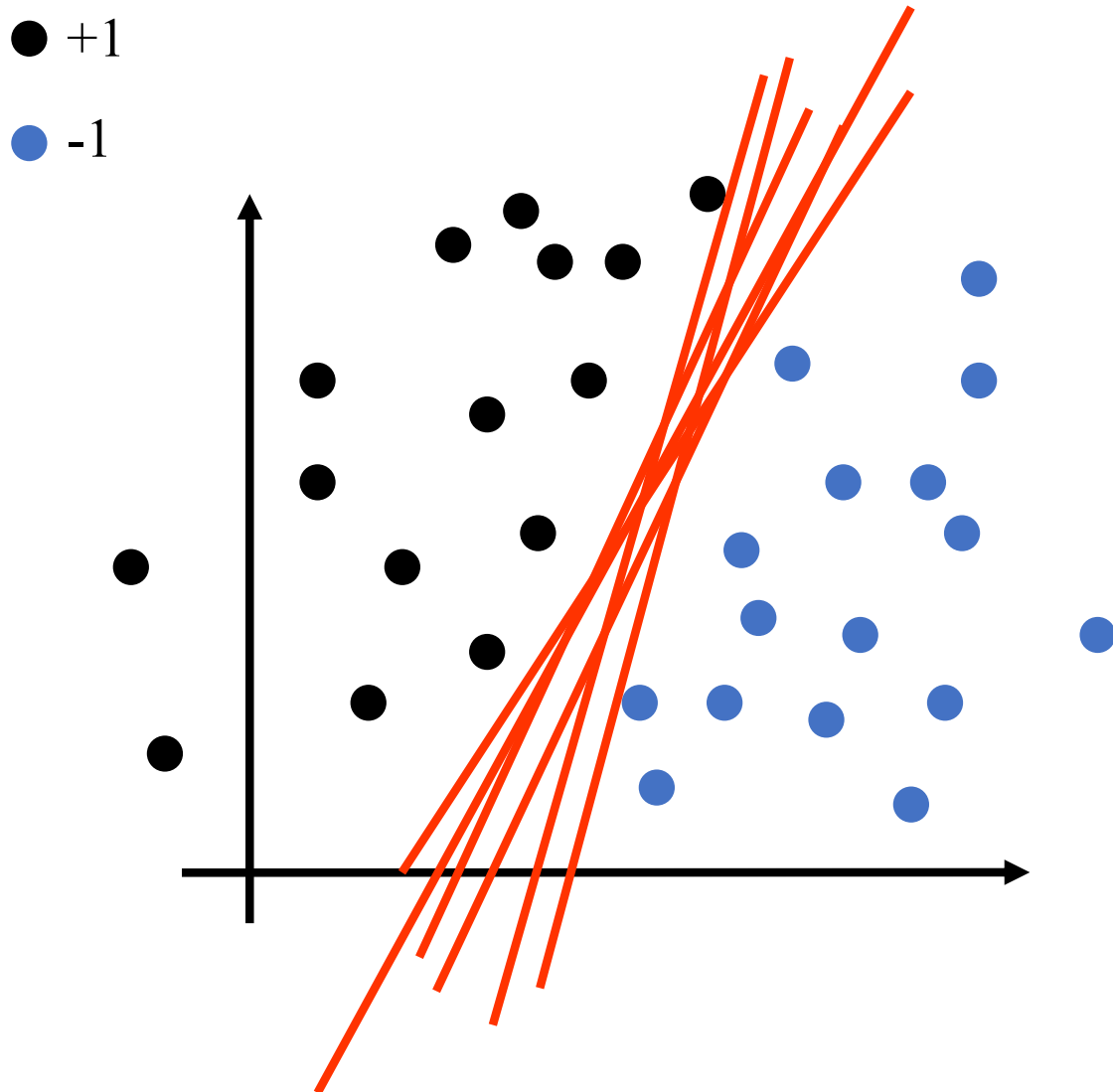
● -1



Podemos ver que puede tener un número infinito de rectas

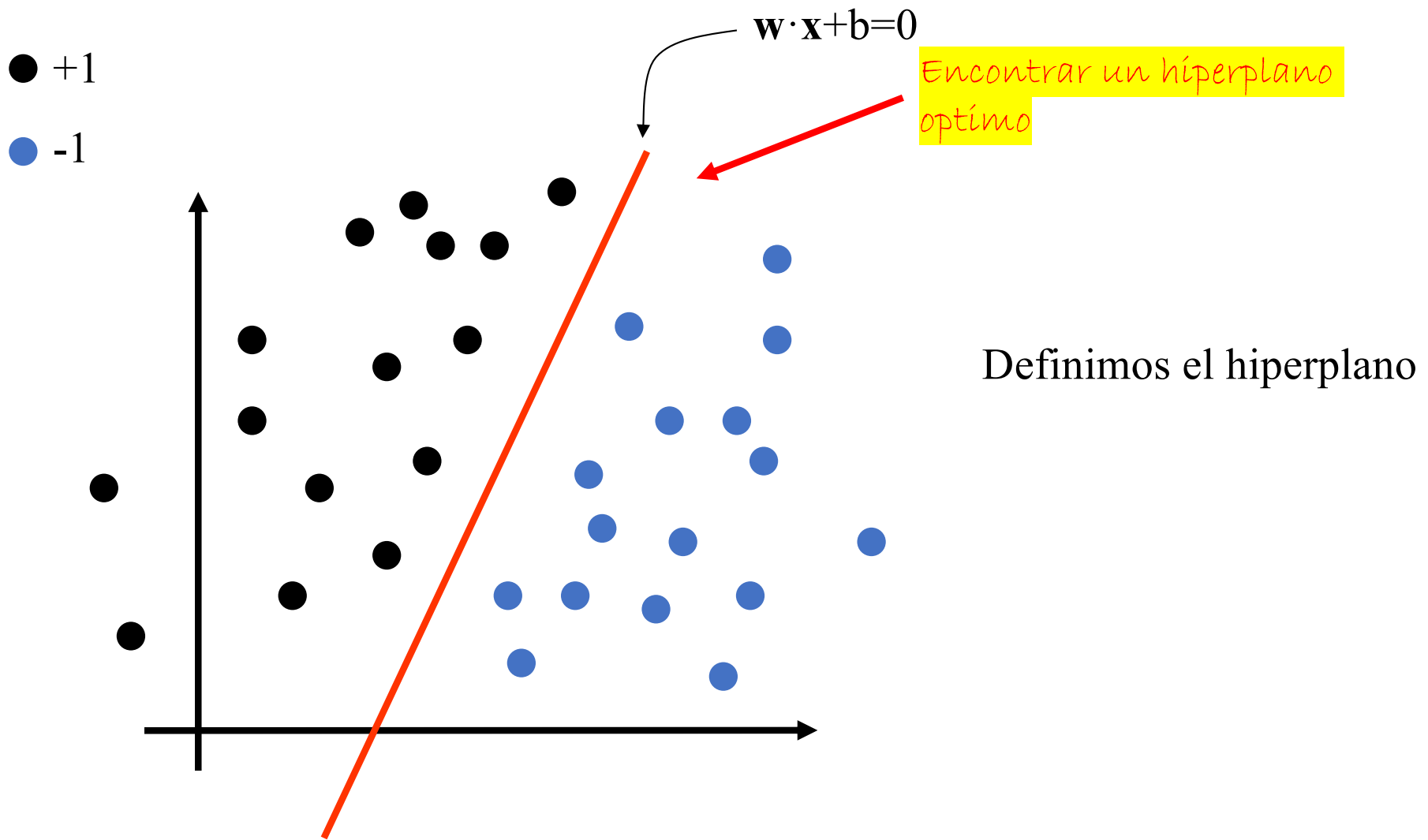
¿cómo clasificar estos datos?

Interpretación geométrica



Cualquiera puede ser buena, ¿pero cuál es la mejor?

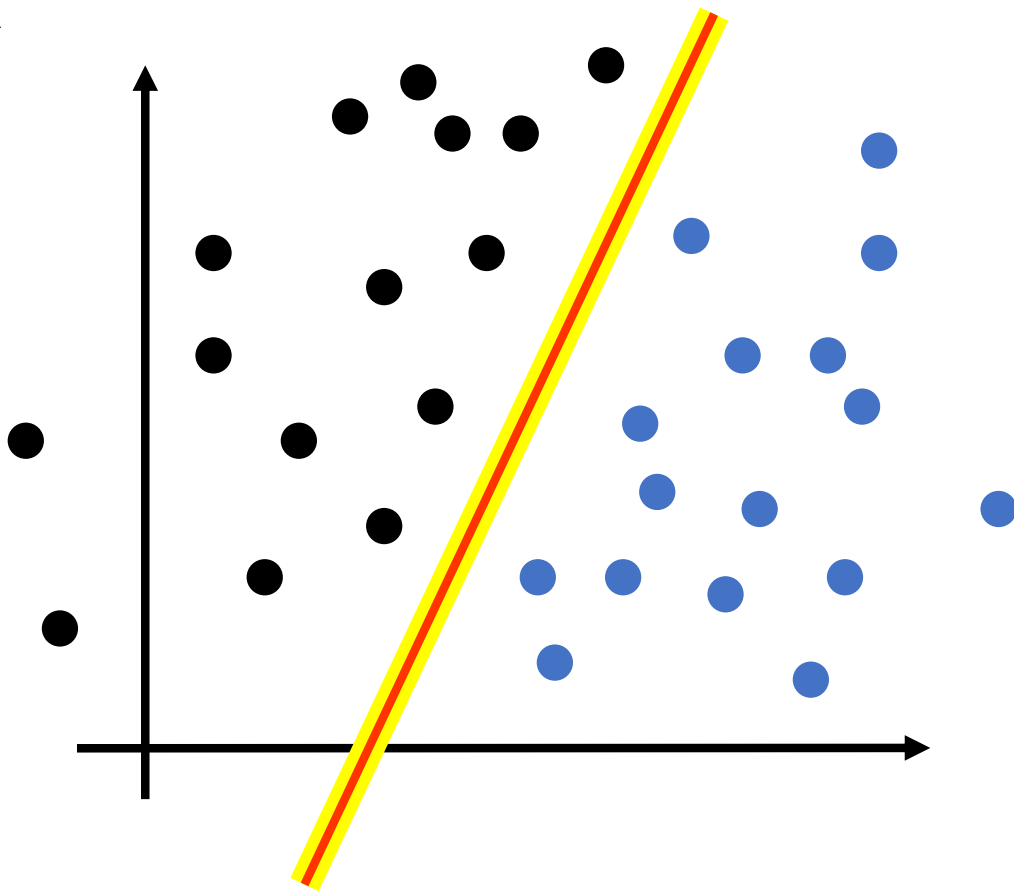
Interpretación geométrica



Interpretación geométrica

● +1

● -1

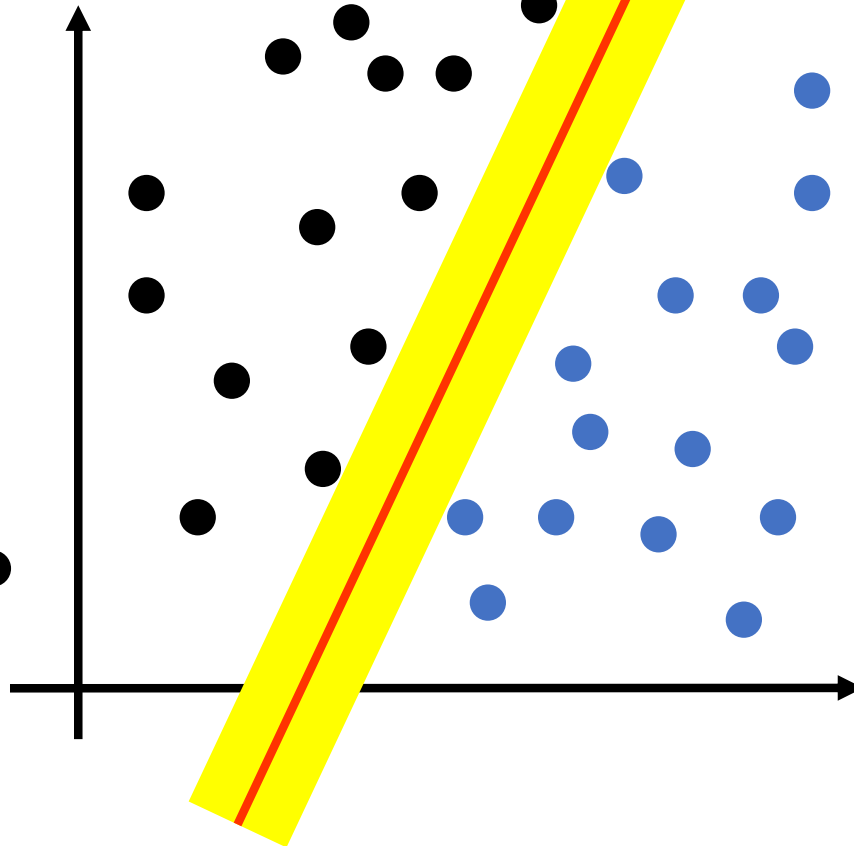


Definimos el margen

Interpretación geométrica

● +1

● -1



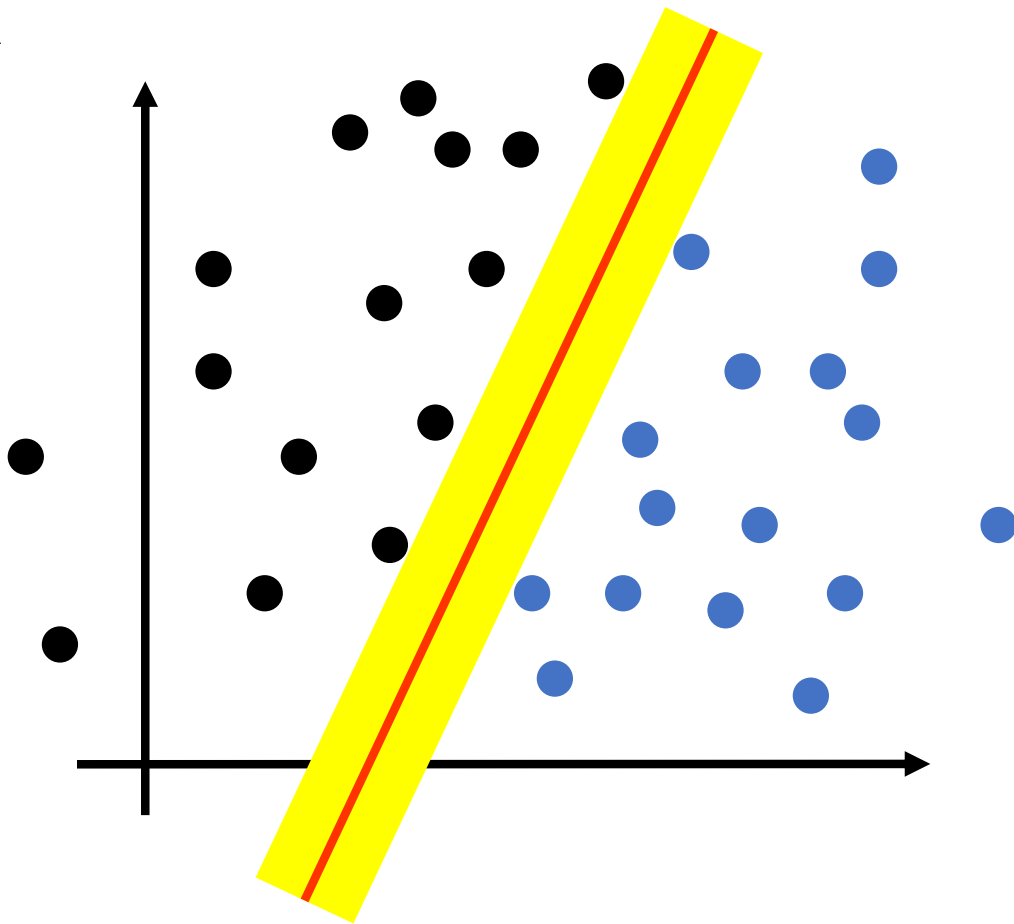
Maxima distancia entre
las muestras

La idea es maximizar el
margen.

Interpretación geométrica

● +1

● -1



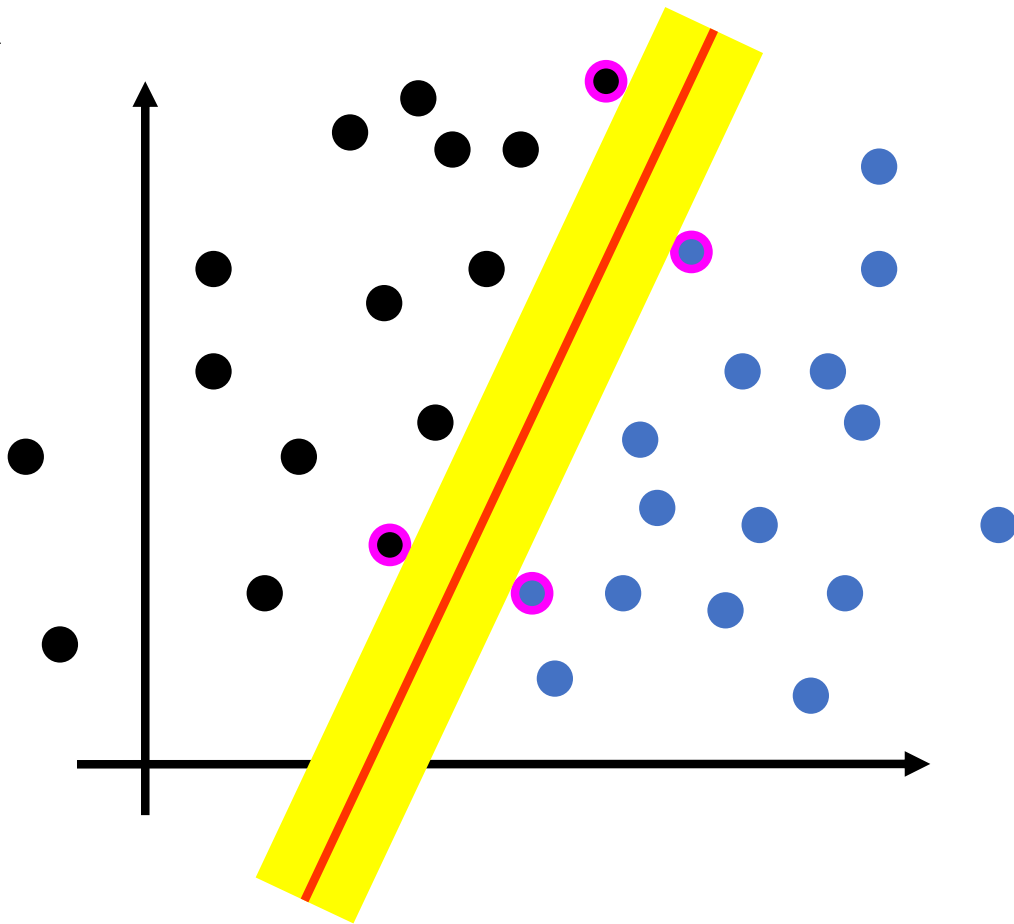
El hiperplano que tenga el mayor margen es el mejor clasificador de los datos.

Esta es la clase más simple de SVM, la LSVM.

Interpretación geométrica

● +1

● -1

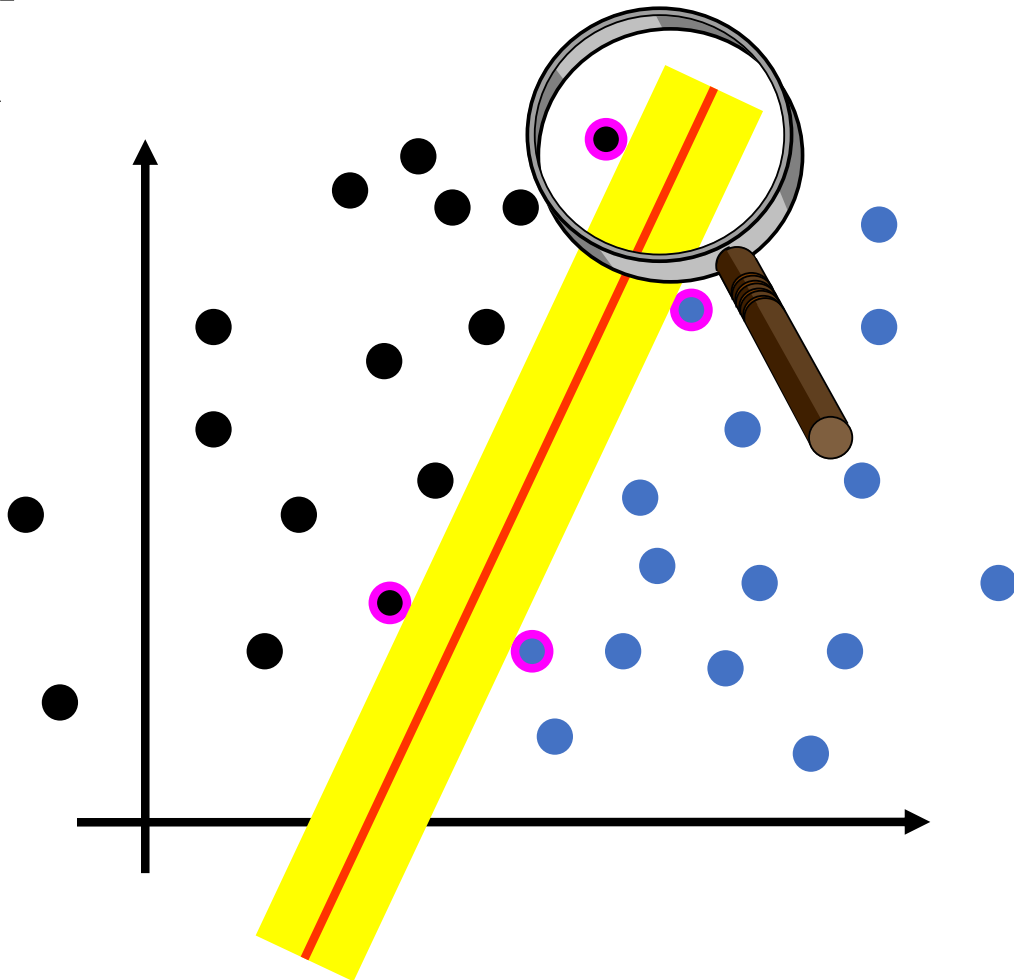


Los **vectores de soporte** son los puntos que tocan el límite del margen.

Interpretación geométrica

● +1

● -1

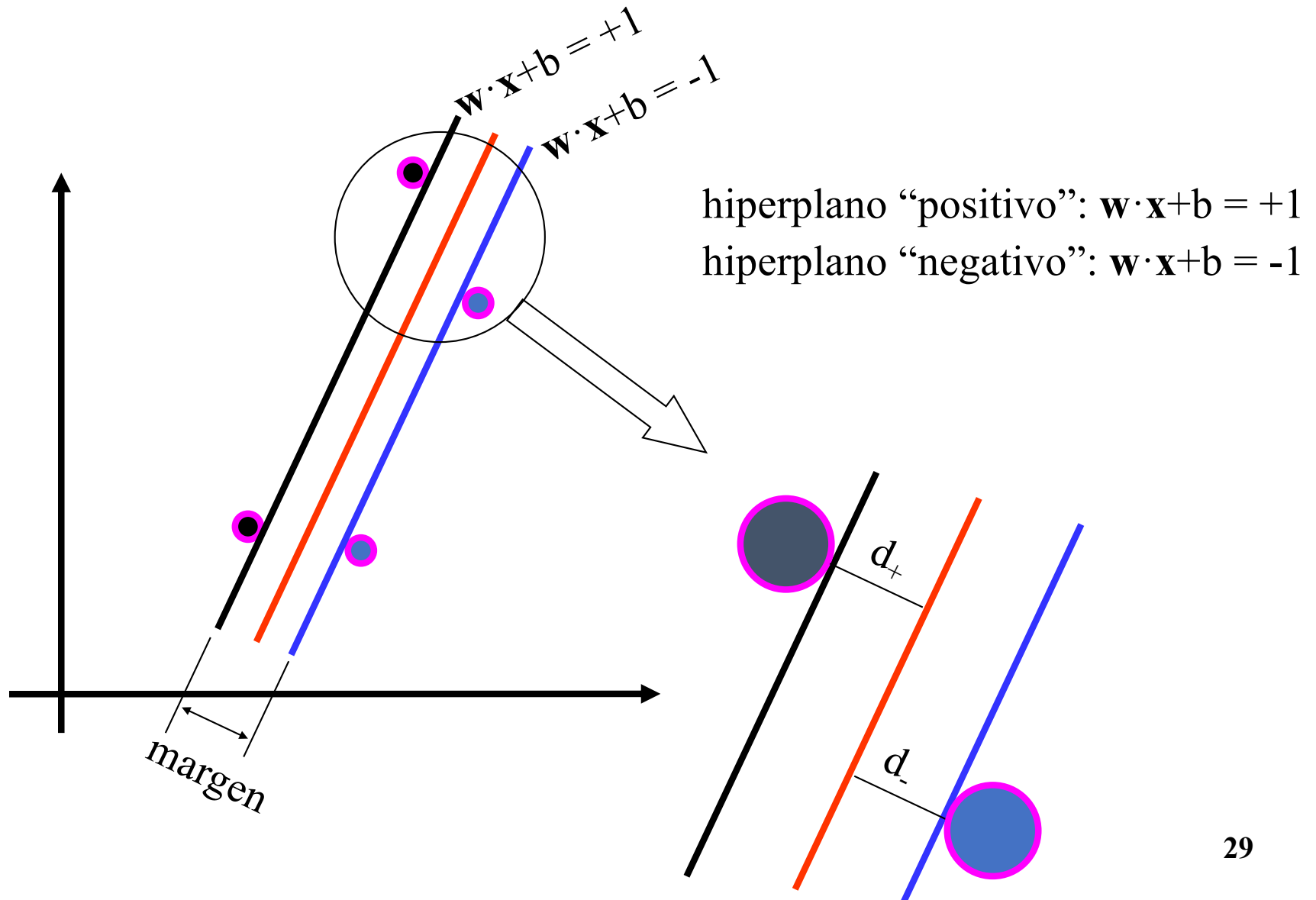


Veamos los hiperplanos
“positivo” y “negativo”

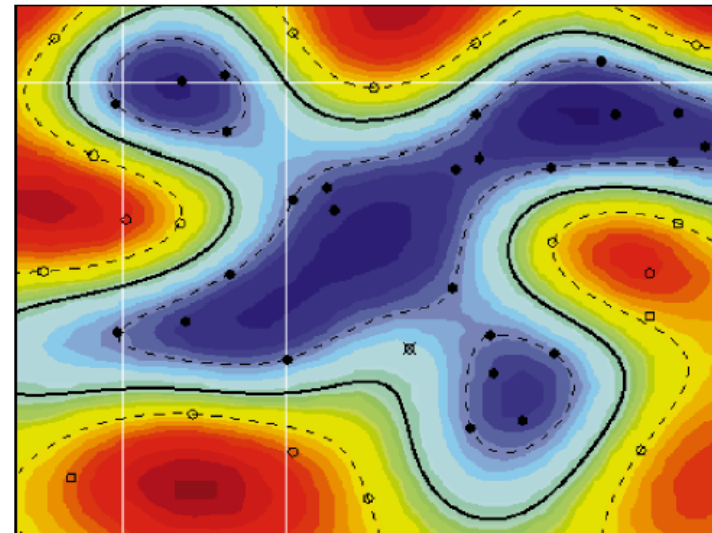
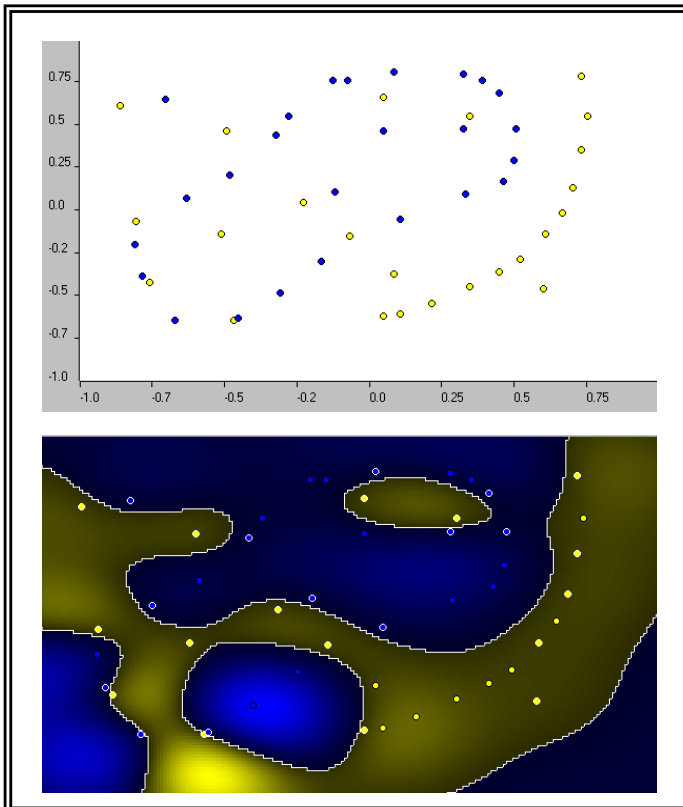
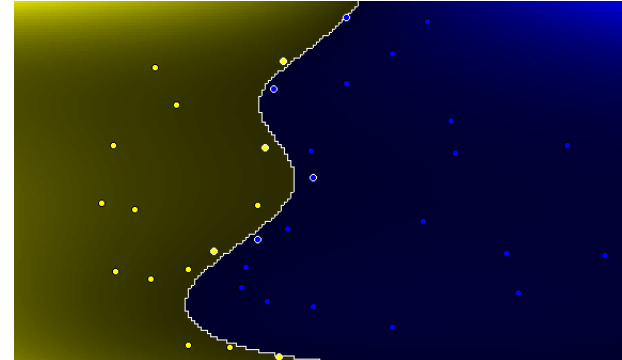
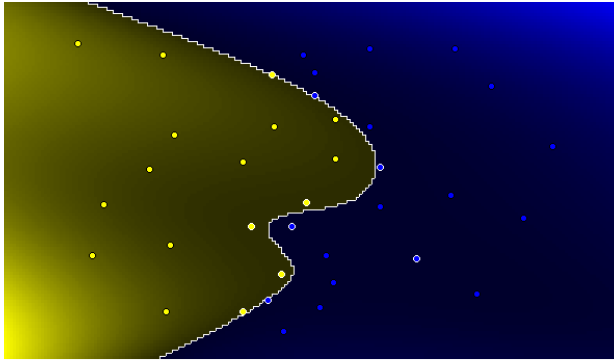
Interpretación geométrica

● +1

● -1



Separación polinómica y con RBF



Gracias!!!