

MAQUINAS DE VECTORES DE SOPORTE

Modelos Lineales

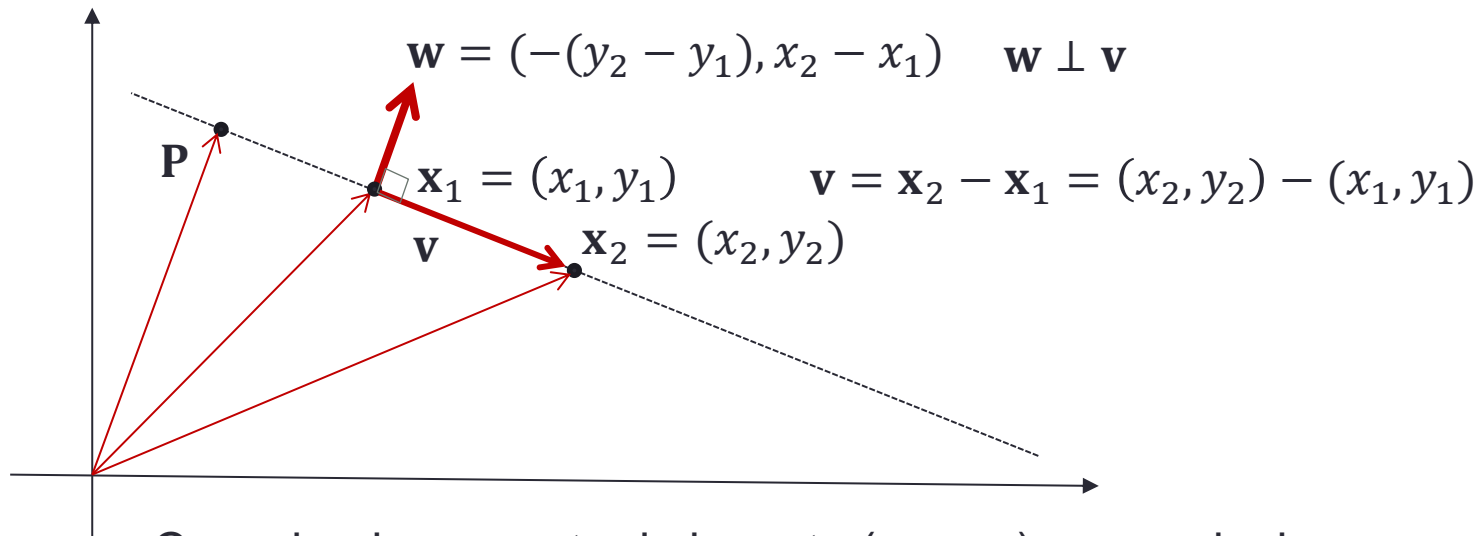
Dr. Jorge Hermosillo Valadez (jhermosillo@uaem.mx)

Depto. de Computación

CInC – IICBA, UAEM

Agosto - 2017

Geometría básica en el plano 2D



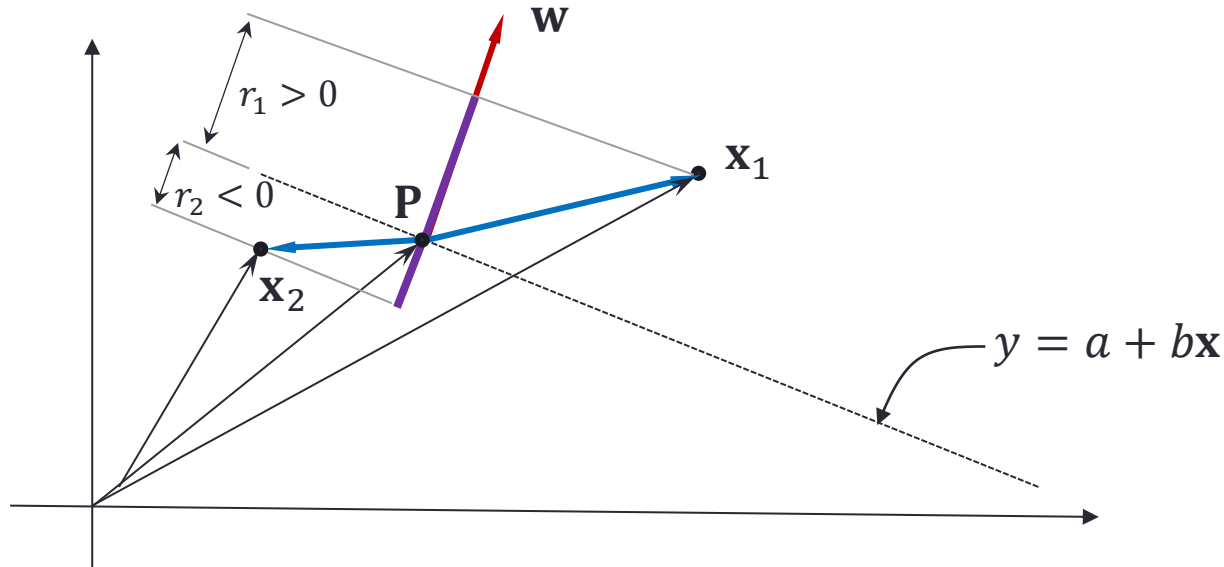
Conociendo un punto de la recta (e.g. \mathbf{x}_1) y \mathbf{w} cualquier punto $\mathbf{P} = (x_P, y_P)$ que pertenezca a la recta satisface:

$$\mathbf{w} \cdot (\mathbf{P} - \mathbf{x}_1) = 0 \rightarrow (y_1 - y_2, x_2 - x_1) \cdot (x_P - x_1, y_P - y_1) = 0$$

$$(y_1 - y_2)(x_P - x_1) + (x_2 - x_1)(y_P - y_1) = 0$$

$$Ax_P + By_P + C = 0$$

Producto punto sobre vector perpendicular a una recta



Para todo punto $\mathbf{x} = (x, y)$ que no está sobre la recta: $y = a + bx$

$$\langle \mathbf{w}, \mathbf{x} - \mathbf{P} \rangle = w_x(x - x_p) + w_y(y - y_p) = r \neq 0$$

$$\langle \mathbf{w}, \mathbf{x}_1 - \mathbf{P} \rangle = r_1 > 0$$

$$\langle \mathbf{w}, \mathbf{x}_2 - \mathbf{P} \rangle = r_2 < 0$$

donde $\langle \square, \square \rangle$ es el producto punto

Resumen del algoritmo Perceptron

- Datos de entrada:
 - $\mathbf{X} := \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subset \mathcal{X}(\text{datos})$, con $\mathbf{x}_i := (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$
 - $\mathbf{y} := \{y_1, y_2, \dots, y_N\} \subset \{\pm 1\}$ (etiquetas)
 - En coordenadas homogéneas:
$$\mathbf{X} := \{(\mathbf{x}_1, 1), (\mathbf{x}_2, 1), \dots, (\mathbf{x}_N, 1)\} \subset \mathbb{R}^N \times \mathbb{R}^{d+1}$$
- Parámetros:
 - $\mathbf{w} := (w_1, w_2, \dots, w_d) \in \mathbb{R}^d$ (vector de pesos)
 - $u \in \mathbb{R}$ (umbral)
- Función (modelo lineal) en coord. homogéneas:
 - $y_i = \text{sign}(\langle \mathbf{w}, \mathbf{x}_i \rangle)$
 - $\mathbf{x}_i := (x_1, x_2, \dots, x_d, 1) \in \mathbb{R}^{d+1}$,
 - $\mathbf{w} := (w_1, w_2, \dots, w_d, w_0) \in \mathbb{R}^{d+1}$ donde $w_0 = u$

Algoritmo de Aprendizaje Perceptron

Entrada: Datos etiquetados de entrenamiento \mathbf{X} en coordenadas homogéneas

Salida: Vector de pesos \mathbf{w} que define al clasificador

$\mathbf{w} = \mathbf{0}$ *#Otras inicializaciones son posibles*

converge = **Falso**

mientras *converge* == **Falso** :

converge = **Verdadero**

para i en $|\mathbf{X}|$:

si $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \leq 0$ **entonces**: *#xi mal clasificado*

$\mathbf{w} = \mathbf{w} + y_i \eta \mathbf{x}_i$ *# $0 < \eta \leq 1$ es la tasa de aprendizaje*

converge = **Falso**

fin

fin

fin

Algoritmo de Aprendizaje Perceptron

Nota que el perceptrón básico aprende los pesos de los atributos

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = y_1 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_1 + y_2 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_2 + \cdots + y_N \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_N \Rightarrow \mathbf{w} = \sum_{i=1}^N y_i \mathbf{x}_i$$

- Luego del entrenamiento cada ejemplo estuvo mal clasificado 0 o α veces (en función del número de épocas; una época es una pasada completa de todos los datos).

$$\mathbf{w} = \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j$$

Perceptron DUAL

- El vector de pesos es una combinación lineal de puntos de entrenamiento:

$$\mathbf{w} = \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j$$

- En este caso la salida del clasificador $y_i = \text{sign}(\langle \mathbf{w}, \mathbf{x}_i \rangle)$ se escribe:

$$y_i = \text{sign} \left(\sum_{j=1}^N \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

- Es decir que podemos aprender los pesos α de las instancias \mathbf{x}_i , en lugar de los pesos de los atributos.
- Por lo tanto la única información necesaria es la matriz n-por-n $\mathbf{G} = \mathbf{X}\mathbf{X}^T$, llamada *matriz Gram*, que contiene todos los productos punto.

Algoritmo del Perceptron DUAL

$\alpha := (\alpha_1, \alpha_2, \dots, \alpha_N) = 0$

converge = *Falso*

mientras *converge* == *Falso* :

converge = *Verdadero*

para *i* **en** $|X|$:

 toma (\mathbf{x}_i, y_i) de los datos

si $y_i \sum_{j=1}^{|X|} \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \leq 0$ **entonces**: #xi mal clasificado

$\alpha_i = \alpha_i + 1$

converge = *Falso*

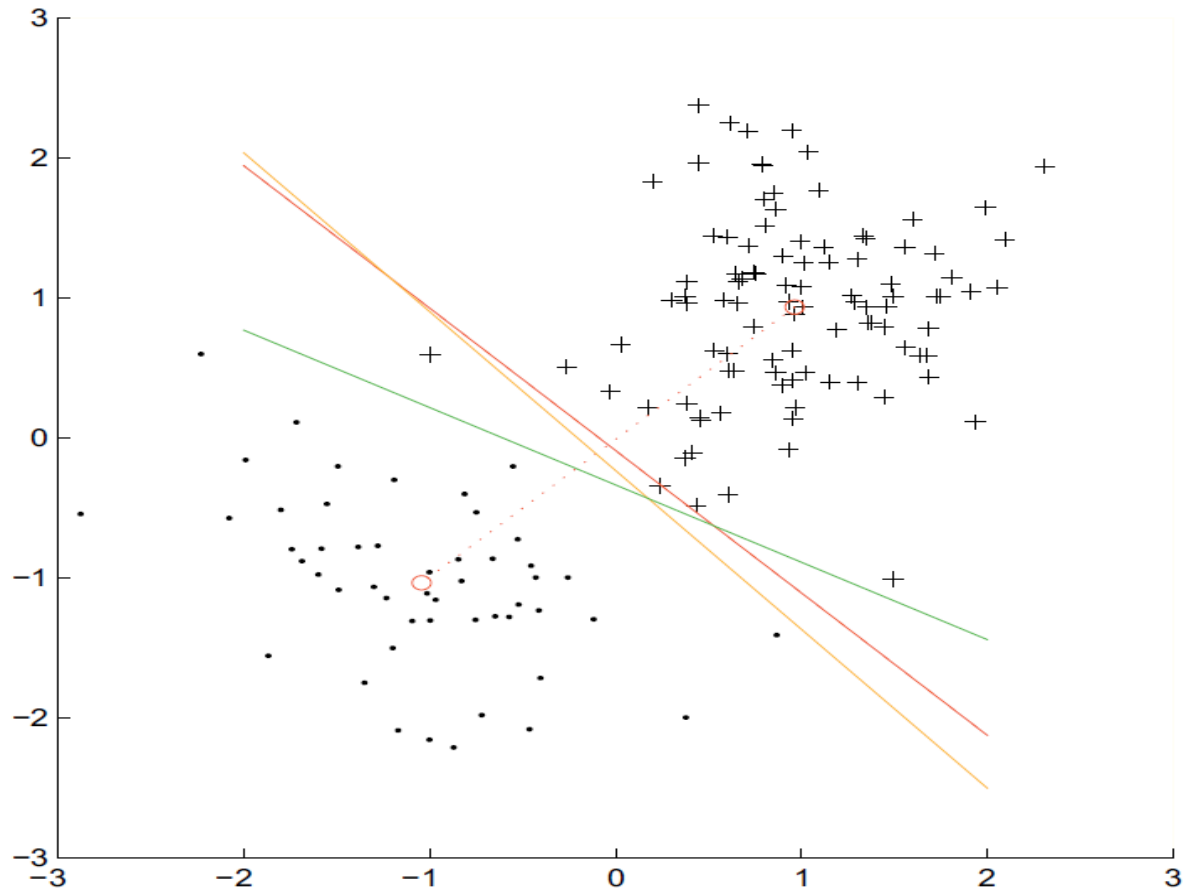
fin

fin

fin

¿Cuál es mejor?

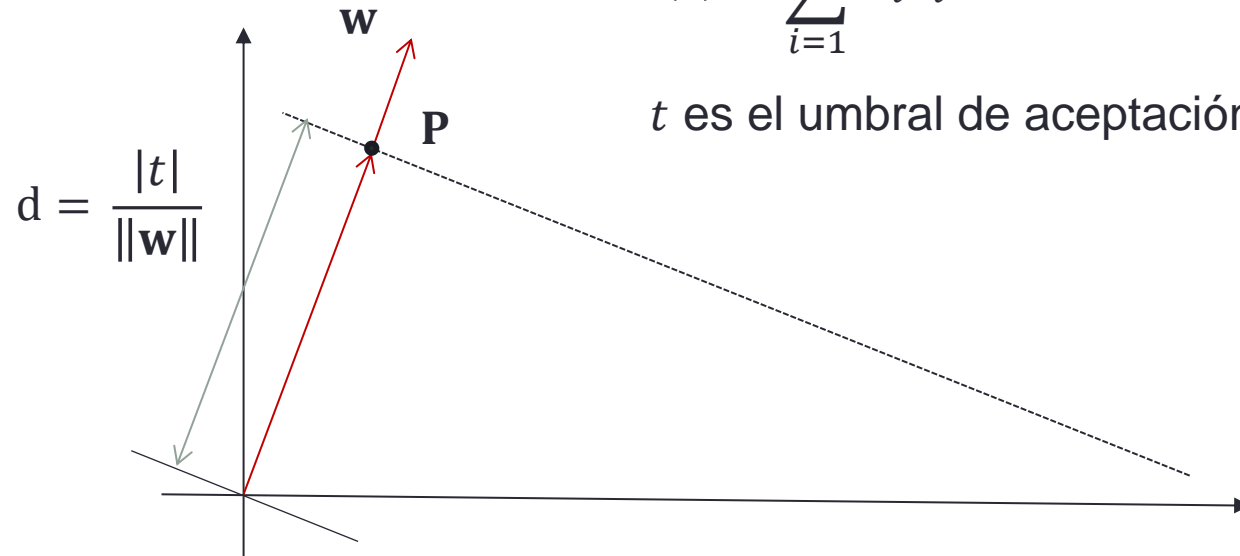
Tres distintos clasificadores lineales: ¿cómo estandarizar la ubicación del hyperplano?



Propiedades del Perceptron

$$h(\mathbf{x}) = \sum_{i=1}^d w_i x_i - t = \mathbf{w} \cdot \mathbf{x} - t$$

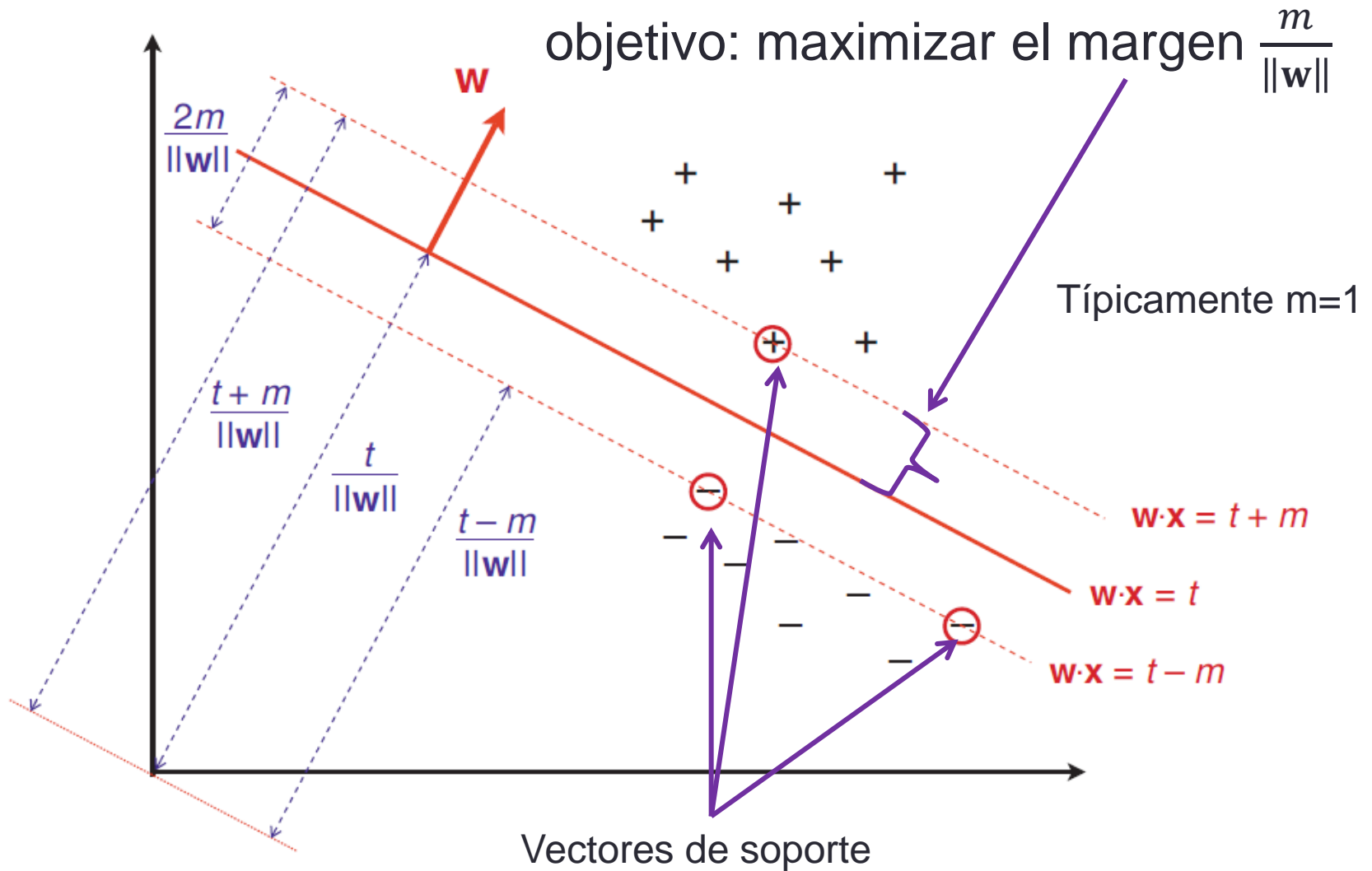
t es el umbral de aceptación



$$h(\mathbf{P}) = \mathbf{w} \cdot \mathbf{P} - t = 0$$

$$\|\mathbf{w}\| \|\mathbf{P}\| \cos \theta = |t|$$

Clasificador por vectores de soporte



Problema de optimización con restricciones

Maximizar el margen es equivalente a minimizar \mathbf{w} o mejor aún $\frac{1}{2} \|\mathbf{w}\|^2$

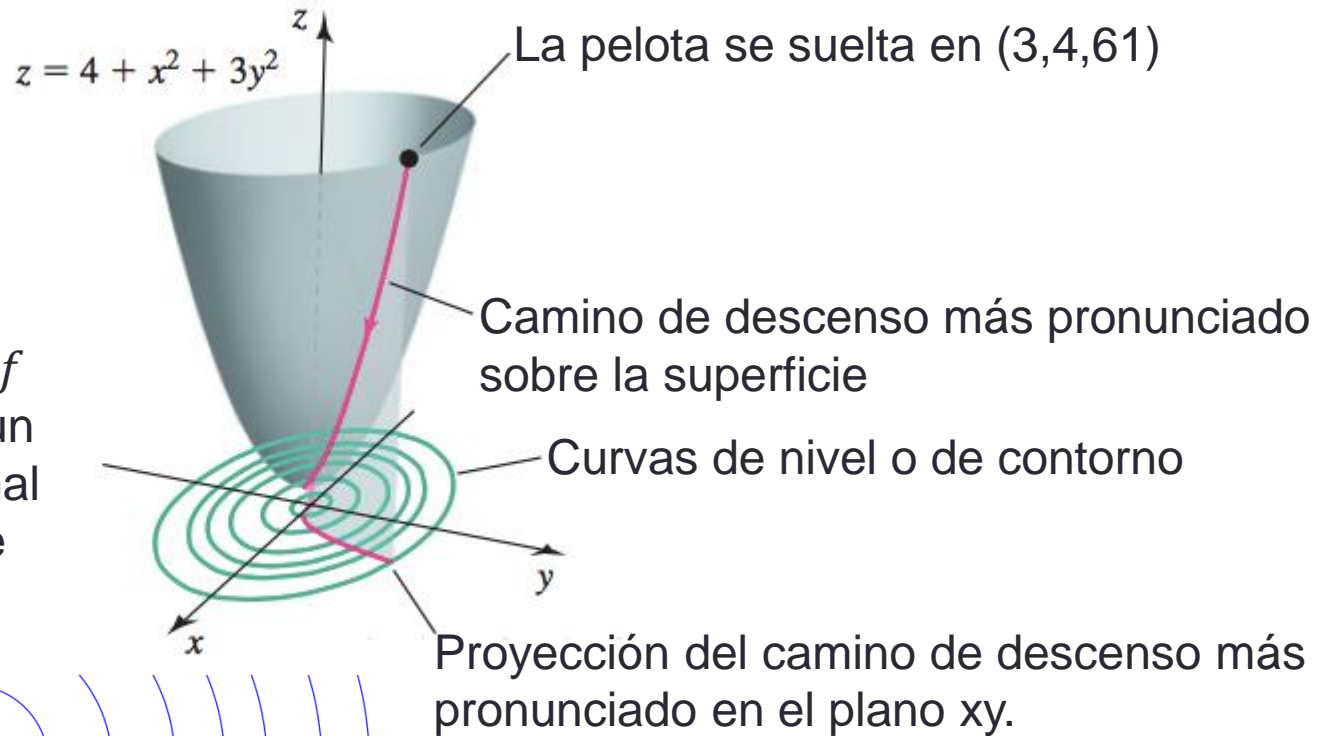
$$\mathbf{w}^*, t^* = \operatorname{argmin}_{\mathbf{w}, t} \frac{1}{2} \|\mathbf{w}\|^2 \text{ sujeto a } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \geq 1, \quad 1 \leq i \leq N$$

Para ello usaremos el metodo de Multiplicadores de Lagrange

Fuentes de consulta:

- <http://www-mtl.mit.edu/Courses/6.050/2004/unit9/wyatt.apr.7.pdf>
- <https://www.khanacademy.org/math/multivariable-calculus/applications-of-multivariable-derivatives/constrained-optimization/a/lagrange-multipliers-single-constraint>

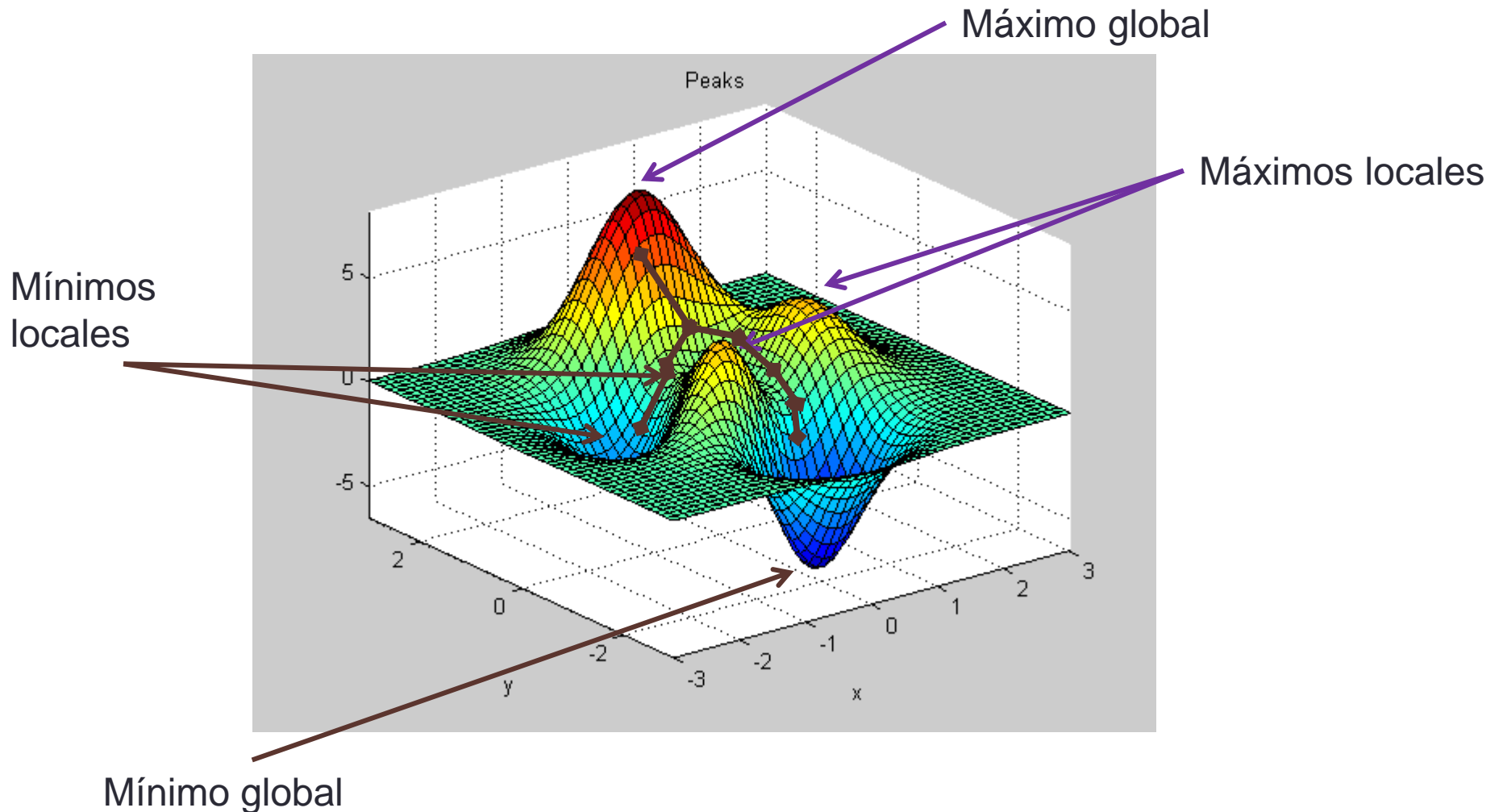
Descenso de gradiente



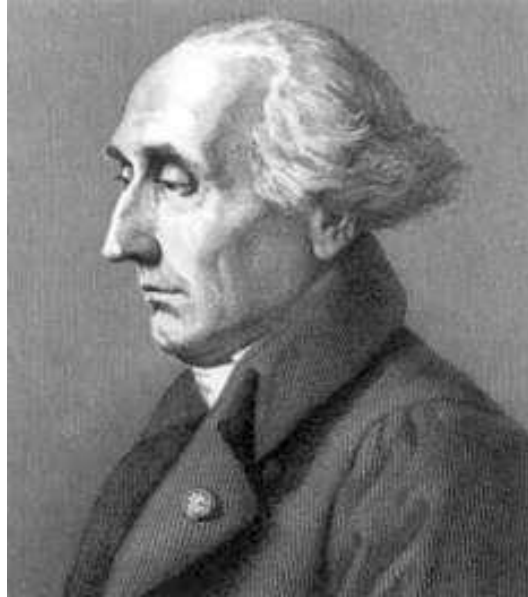
El gradiente ∇f de $f(y,x)$ es un vector ortogonal a las líneas de contorno



Mínimos/Máximos locales y globales



Método de Lagrange



Joseph Louis de Lagrange

Turín, 25 de enero de 1736-París, 10 de abril de 1813), fue un físico, matemático y astrónomo italiano naturalizado francés, que después de formarse en su Italia natal pasó la mayor parte de su vida en Prusia y Francia.

Método de Lagrange

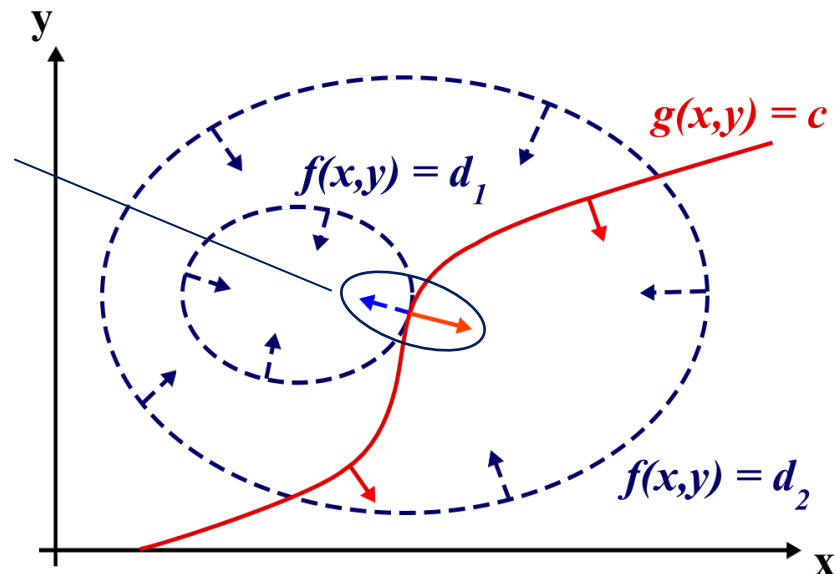
- La técnica del multiplicador de Lagrange permite encontrar el máximo o mínimo de una función multivariable $f(x_1, x_2, \dots, x_N)$ cuando hay alguna restricción en los valores de entrada del tipo:

$$g(x_1, x_2, \dots, x_N) = c \quad c \in \mathbb{R}$$

- La idea central es buscar puntos en los que las líneas de contorno de $f(x_1, x_2, \dots, x_N)$ y $g(x_1, x_2, \dots, x_N)$ son tangentes entre sí.

puntos donde los vectores de gradiente de ambas curvas son paralelos entre sí.

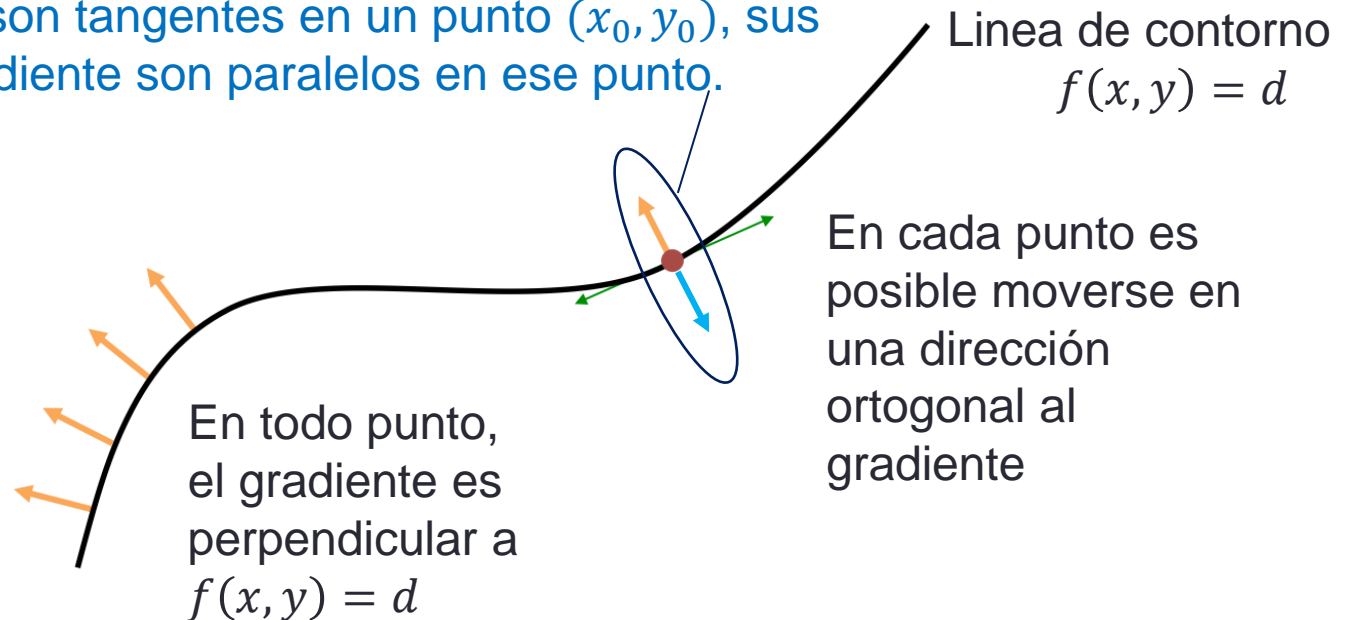
¡Veamos el gradiente!



Método de Lagrange

- Como hemos mencionado el gradiente ∇f de f evaluado en un punto (x_0, y_0) siempre da un vector perpendicular a la línea de contorno que pasa por ese punto.

Esto significa que cuando las líneas de contorno de dos funciones f y g son tangentes en un punto (x_0, y_0) , sus vectores de gradiente son paralelos en ese punto.



Multiplicador de Lagrange

- La condición de paralelismo se traduce en :

$$\nabla f(x_1, x_2, \dots, x_N) = \lambda_0 \nabla g(x_1, x_2, \dots, x_N) \quad (1)$$

- Para una función $f(x_1, x_2, \dots, x_N)$:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_N} \end{bmatrix}$$

- Notación práctica: $\nabla f = (\partial_{x_1} f, \partial_{x_2} f, \dots, \partial_{x_N} f)$

- En esta forma (1) se puede desglosar como:

$$\partial_{x_1} f = \lambda_0 \partial_{x_1} g$$

$$\partial_{x_2} f = \lambda_0 \partial_{x_2} g$$

$$\vdots$$

$$\partial_{x_N} f = \lambda_0 \partial_{x_N} g$$

Lagrangiano

- Lagrange propuso una nueva función especial (el Lagrangiano) que recoge las mismas variables de entrada que f y g , junto con λ que ahora se considera una variable en lugar de una constante:

$$\mathcal{L}(x_1, x_2, \dots, x_N, \lambda) = f(x_1, x_2, \dots, x_N) - \lambda(g(x_1, x_2, \dots, x_N) - c)$$

- Los puntos de tangencia se encuentran calculando:

$$\nabla \mathcal{L}(x_1, x_2, \dots, x_N, \lambda) = 0$$

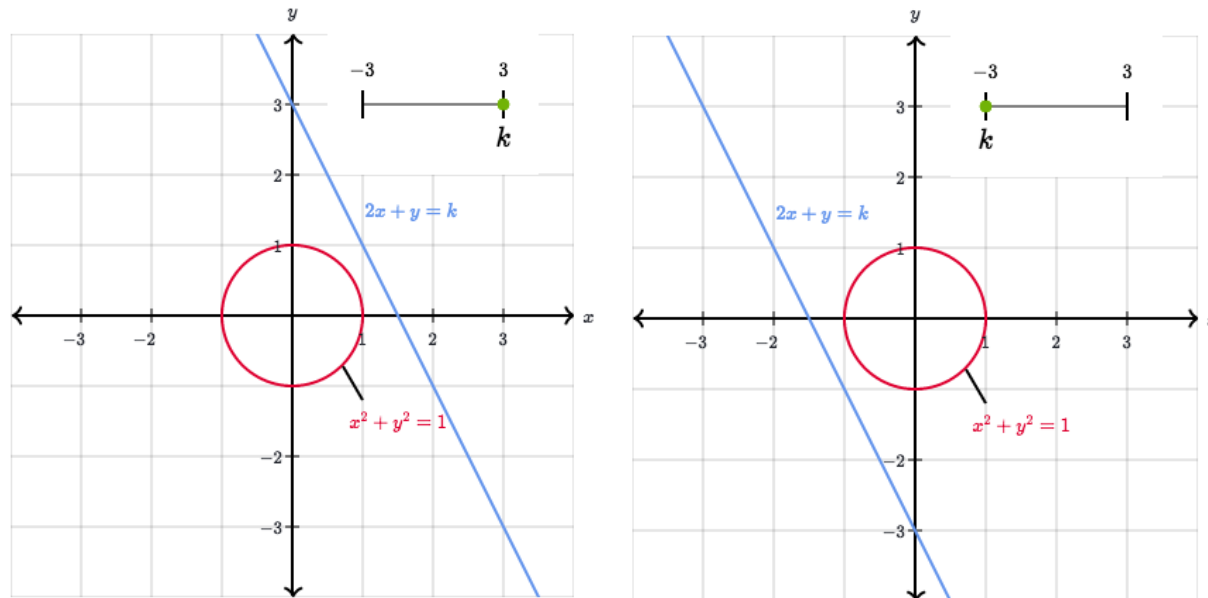
Ejemplo

- Encuentre el valor máximo de las líneas de contorno f :

$$f(x, y) = 2x + y = k$$

con la restricción sobre la entrada (x, y) tal que éstas deben cumplir con:

$$g(x, y) = x^2 + y^2 = 1$$



Ejemplo (Lagrangiano)

- Para este ejemplo el Lagrangiano se escribe:

$$\begin{aligned}\mathcal{L}(x, y, \lambda) &= f(x, y) - \lambda(g(x, y) - 1) \\ &= 2x + y - \lambda(x^2 + y^2 - 1)\end{aligned}$$

- Podemos ver que nuestra restricción se escribe:

$$\partial_{\lambda}\mathcal{L}(x, y, \lambda) = 0 \implies g(x, y) = 1$$

- También vemos que $\nabla\mathcal{L}$ tiene dos componentes:

$$\nabla f = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \text{ y } \nabla g = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

- Por lo que la condición de tangencia $\nabla\mathcal{L} = 0$ se escribe:

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} = \lambda_0 \begin{bmatrix} 2x_0 \\ 2y_0 \end{bmatrix}$$

Ejemplo (planteamiento de la solución)

- Resumiendo buscamos:

- Un punto (x_0, y_0) tal que: $g(x_0, y_0) = 1$ que para el ejemplo eso es:

$$x_0^2 + y_0^2 = 1$$

- Además este punto es tal que:

$$\begin{cases} 2\lambda_0 x_0 = 2 \\ 2\lambda_0 y_0 = 1 \end{cases}$$

- Calcula (x_0, y_0) y λ_0 :

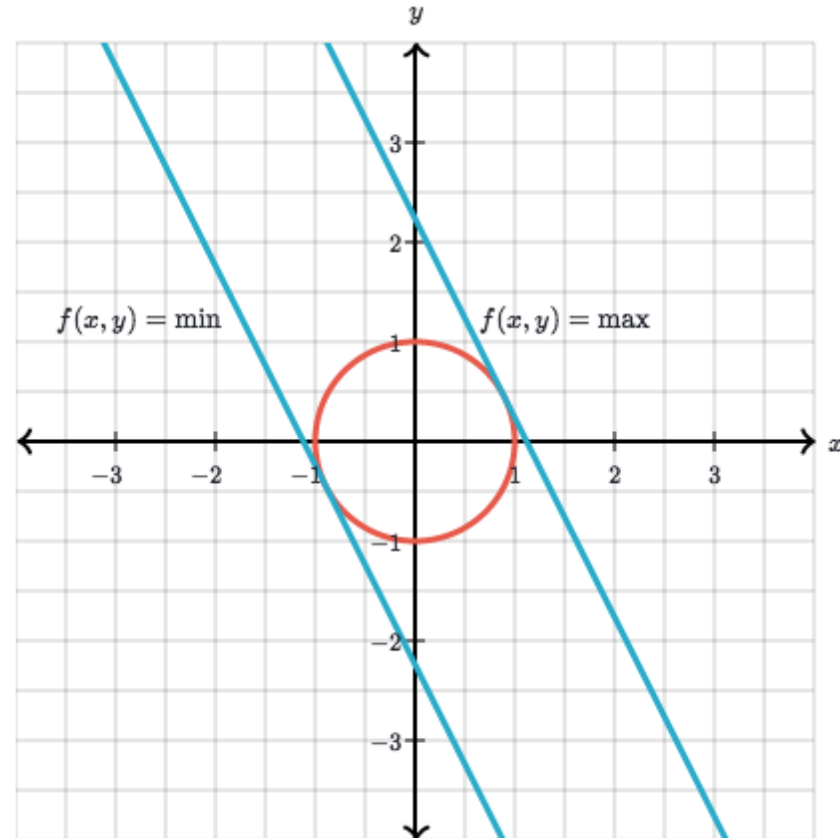
Ejemplo (solución)

$$\lambda_0 = \frac{\pm\sqrt{5}}{2} \quad (x_0, y_0) = \left(\frac{1}{\lambda_0}, \frac{1}{2\lambda_0} \right)$$

$$(x_0, y_0) = \left(\frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}} \right) \Rightarrow f(x_0, y_0) = \max$$

o bien

$$(x_0, y_0) = \left(\frac{-2}{\sqrt{5}}, \frac{-1}{\sqrt{5}} \right) \Rightarrow f(x_0, y_0) = \min$$



Método de optimización por multiplicadores de Lagrange

- Paso 1: Define el lagrangiano

$$\mathcal{L}(x, y, \dots, \lambda) = f(x, y, \dots) - \lambda(g(x, y, \dots) - c)$$

- Paso 2: Plantea las ecuaciones de soluciones óptimas

$$\nabla \mathcal{L}(x, y, \dots, \lambda) = 0$$

- Paso 3:
 - Considera cada solución, que se verá como $(^1x_0, ^2x_0, \dots, ^Nx_0, \lambda_0)$.
 - Introduce cada punto sin λ en f como entrada.
 - Cualquiera que dé el mayor (o más pequeño) valor es el punto máximo (o mínimo) que estás buscando.

El problema de optimización en SVM's

- Nuestro objetivo es:

$$\mathbf{w}^*, t^* = \operatorname{argmin}_{\mathbf{w}, t} \frac{1}{2} \|\mathbf{w}\|^2$$

- Sujeto a las siguientes N restricciones:

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \geq 1, \quad 1 \leq i \leq N$$

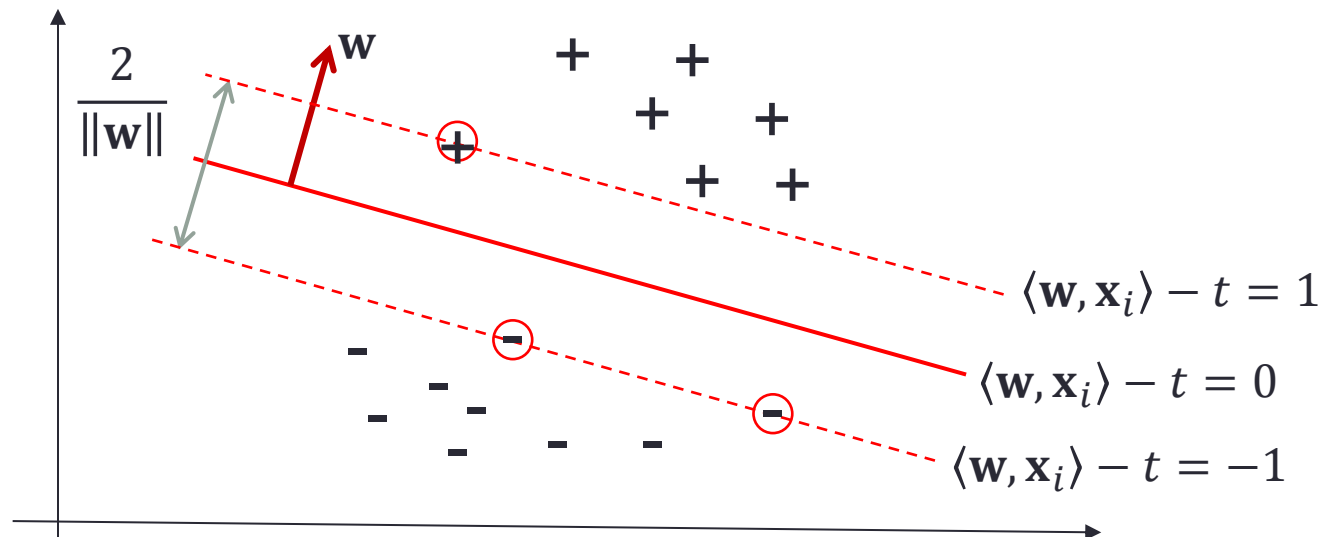
El problema de optimización en SVM's

- Nuestro objetivo es:

$$\mathbf{w}^*, t^* = \operatorname{argmin}_{\mathbf{w}, t} \frac{1}{2} \|\mathbf{w}\|^2$$

- Sujeto a las siguientes N restricciones:

$$y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle - t) \geq 1, \quad 1 \leq i \leq N$$



El problema de optimizacion en SVM's

- Definimos el lagrangiano

$$\begin{aligned}\mathcal{L}_P(\mathbf{w}, t, \alpha_1, \dots, \alpha_N) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - t) - 1) \\ &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + \sum_{i=1}^N \alpha_i y_i t + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle - \langle \mathbf{w}, \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \rangle + t \left(\sum_{i=1}^N \alpha_i y_i \right) + \sum_{i=1}^N \alpha_i\end{aligned}$$

- Para un t óptimo $\partial_t \mathcal{L}_P = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$
- Para pesos óptimos $\partial_{\mathbf{w}} \mathcal{L}_P = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$

Modelo dual de optimización

- Reinsertando estas expresiones en \mathcal{L}_P obtenemos \mathcal{L}_D el lagrangiano del problema dual:

$$\begin{aligned}\mathcal{L}_D(\alpha_1, \dots, \alpha_N) &= -\frac{1}{2} \left\langle \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right\rangle + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i\end{aligned}$$

Modelo dual de optimización

- El problema de optimización dual es el siguiente:

$$\alpha_1^*, \dots, \alpha_N^* = \operatorname{argmax}_{\alpha_1, \dots, \alpha_N} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i$$

- Sujeto a las restricciones:

$$\alpha_i > 0, \quad 1 \leq i \leq N \quad \text{y} \quad \sum_{i=1}^N \alpha_i y_i = 0$$

Aspectos importantes

- La forma dual del problema de optimización de las SVM ilustra dos aspectos importantes :

$$\alpha_1^*, \dots, \alpha_N^* = \underset{\alpha_1, \dots, \alpha_N}{\operatorname{argmax}} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i$$

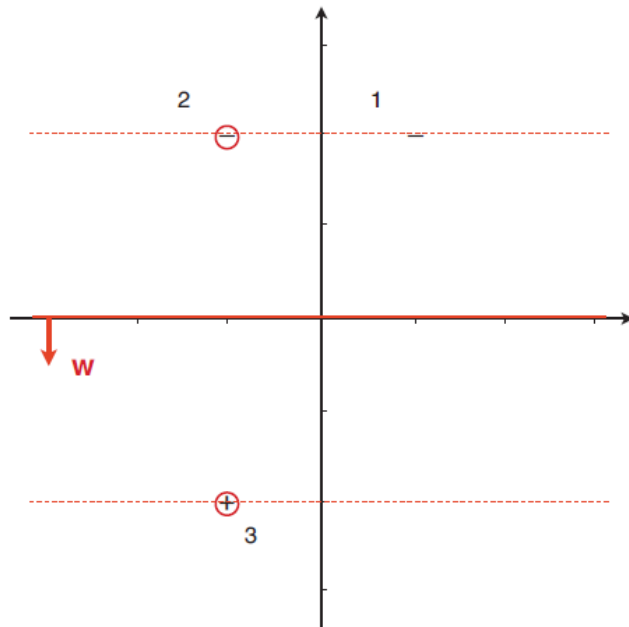
1. Maximizar el margen es equivalente a encontrar los vectores de soporte; es decir, los puntos para los cuales los multiplicadores de Lagrange son no nulos.
2. El problema de optimización está completamente definido por el producto punto de pares de instancias de entrenamiento: las entradas de la matriz Gram.

Ejercicio

Let the data points and labels be as follows (see Figure):

$$\mathbf{X} = \begin{pmatrix} 1 & 2 \\ -1 & 2 \\ -1 & -2 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -1 \\ -1 \\ +1 \end{pmatrix} \quad \mathbf{X}' = \begin{pmatrix} -1 & -2 \\ 1 & -2 \\ -1 & -2 \end{pmatrix}$$

The matrix \mathbf{X}' on the right incorporates the class labels; i.e., the rows are $y_i \mathbf{x}_i$.



1. Encuentre la matriz Gram
2. Expresé el problema de optimización Dual
3. Encuentre los vectores de soporte

