



# NOCIONES BÁSICAS EN MACHINE LEARNING

Dr. Jorge Hermosillo

Laboratorio de Semántica Computacional



# MODELOS LINEALES DE REGRESIÓN

- **Regresión lineal:** combinación lineal de variables de entrada:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D \text{ donde } \mathbf{x} = (x_1, \dots, x_D)^T$$

- **Regresión lineal con funciones base:** combinaciones lineales de funciones no-lineales de las variables de entrada:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

- Sea  $\phi_0(\mathbf{x}) = 1$ , de tal forma que:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

donde  $\mathbf{w} = (w_0, \dots, w_{M-1})^T$  y  $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$

# FEATURES (CARACTERÍSTICAS) COMUNES

- ▶ Hemos visto la importancia de hacer pre-procesamiento de los datos, también llamado el ***proceso de extracción de características*** (“*feature extraction*”) de las variables de entrada.
- ▶ Si los datos originales están dados por el vector  $\mathbf{x}$ , entonces sus características están dadas por  $\{\phi_j(\mathbf{x})\}$

- ▶ Funciones base Polinomiales (“*polynomial features*”):

$$\phi_j(\mathbf{x}) = \mathbf{x}^j$$

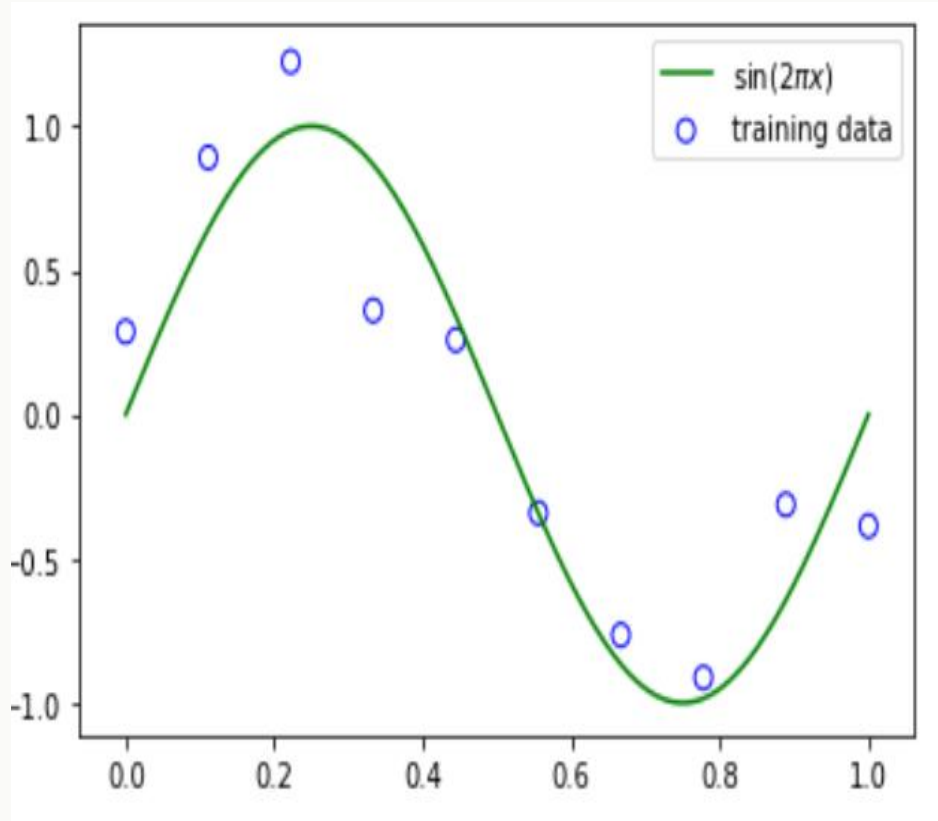
- ▶ Funciones base Gaussianas (“*Gaussian features*”)

$$\phi_j(\mathbf{x}) = \exp \left\{ -\frac{(\mathbf{x} - \mu_j)^2}{2s^2} \right\}$$

- ▶ Funciones base Sigmoidales (“*sigmoidal features*”)

$$\phi_j(\mathbf{x}) = \sigma \left( \frac{\mathbf{x} - \mu_j}{s} \right) \quad \text{donde } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

# REGRESIÓN POLINOMIAL



## DATOS DE ENTRADA:

Un conjunto de  $N$  datos de entrenamiento y sus respectivos valores objetivo:

$$\mathbf{x} \equiv (x_1, x_2, \dots, x_N)^T$$

$$\mathbf{t} \equiv (t_1, t_2, \dots, t_N)^T$$

donde:

$$t_n = \sin(2\pi x_n) + \epsilon$$

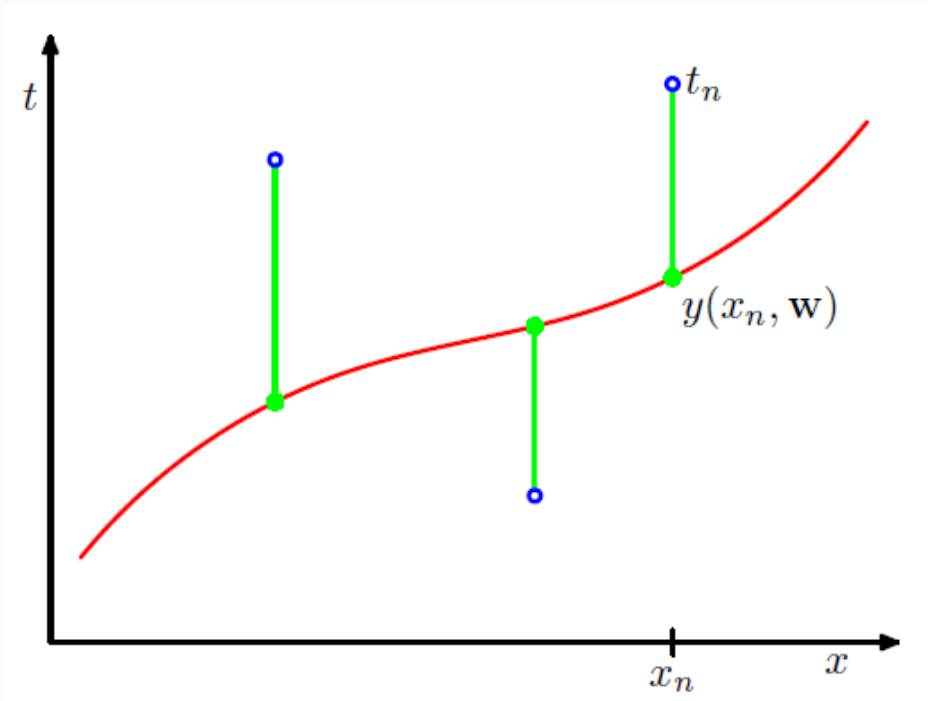
con  $\epsilon \sim \mathcal{N}(0, (0.3)^2)$

## SALIDA:

Un modelo lineal de ajuste polinomial:

$$y(x, \mathbf{w}) = w_0 + w_1 x + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

# FUNCIÓN DE ERROR/PÉRDIDA (LOSS FUNCTION)



$$y(x, \mathbf{w}) = w_0 + w_1x + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

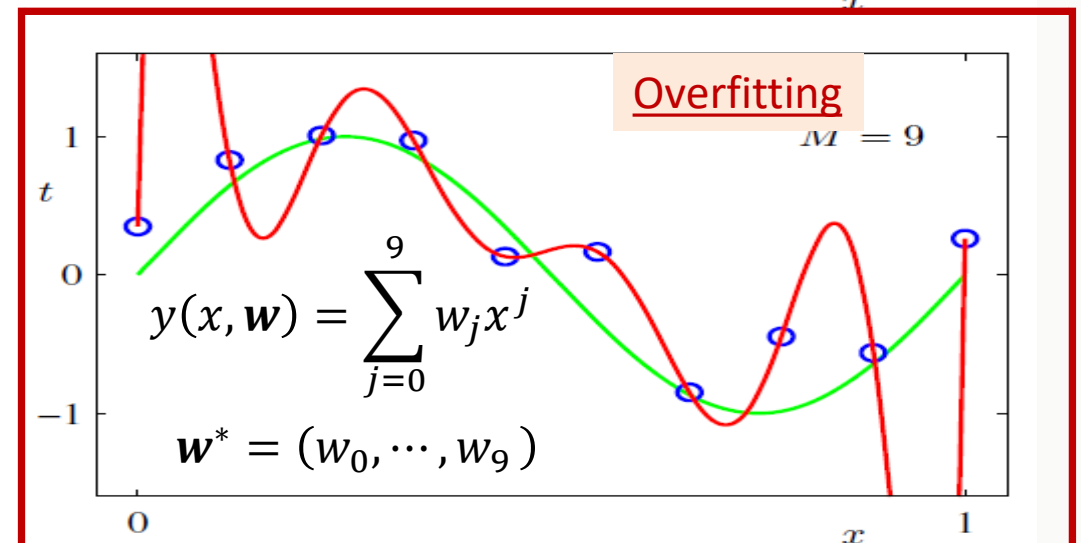
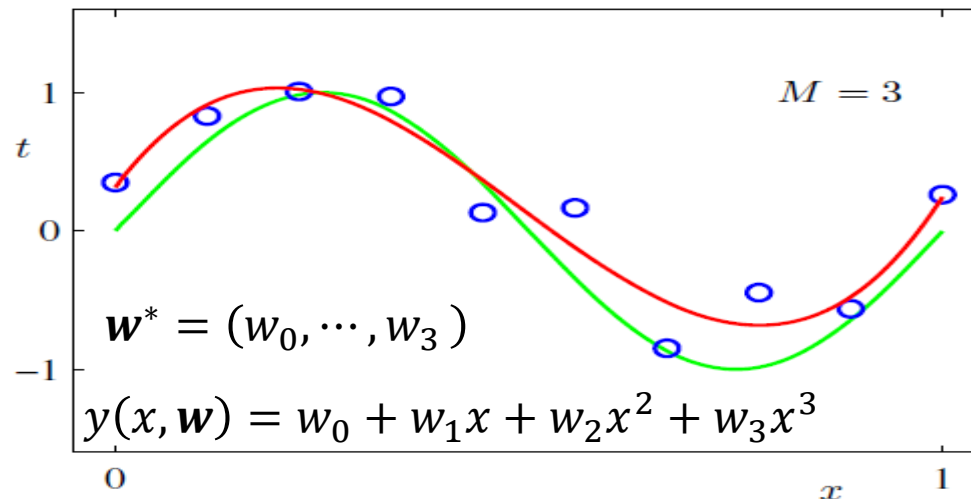
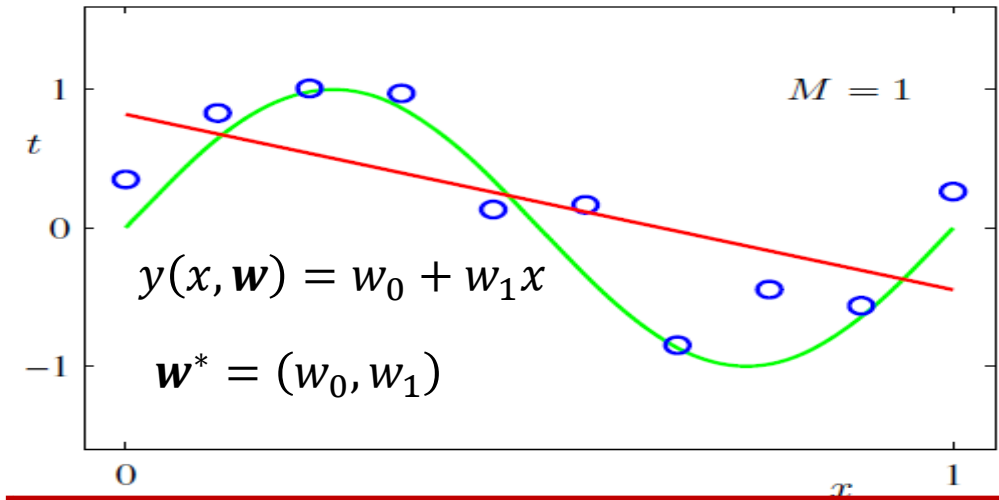
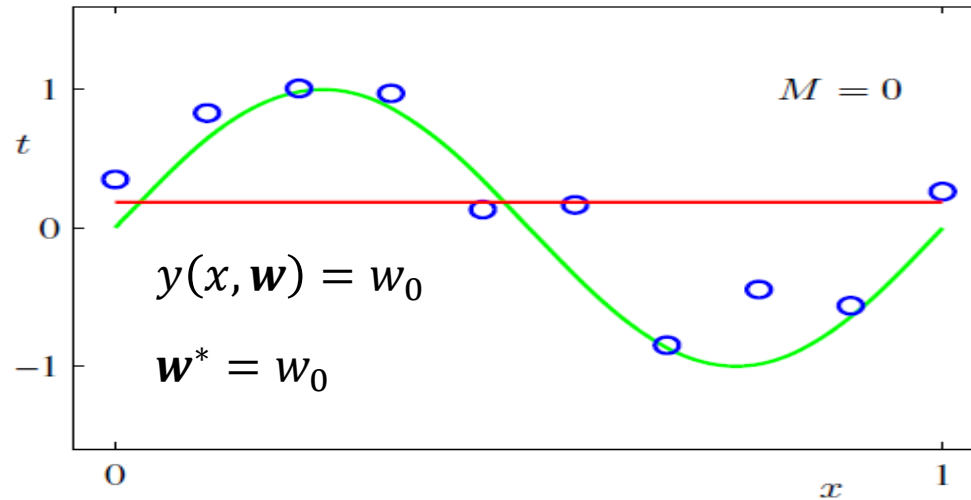
el ajuste se hace mediante la minimización de la función de error:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

Donde  $t_n$  es el objetivo (valor de la muestra – puntos azules en la figura).

Tomada de: Bishop, Christopher M. ( 2006).  
Pattern recognition and machine learning. New  
York. Springer.

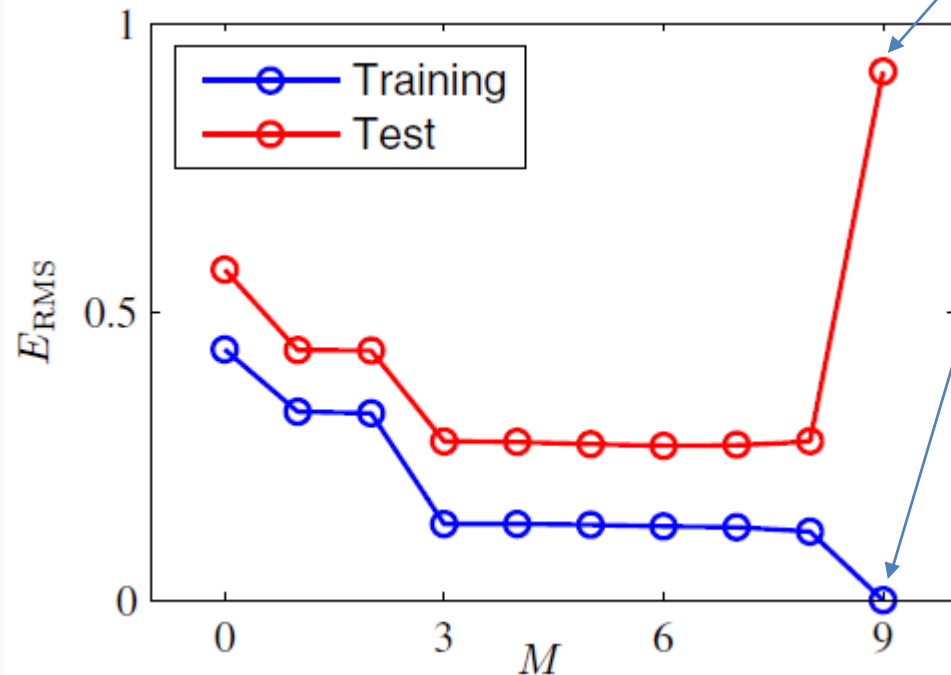
# EFFECTO DE LA COMPLEJIDAD DEL MODELO EN EL AJUSTE



# MÉTRICA DE DESEMPEÑO: $E_{\text{RMS}}$

Error cuadrático medio:

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$



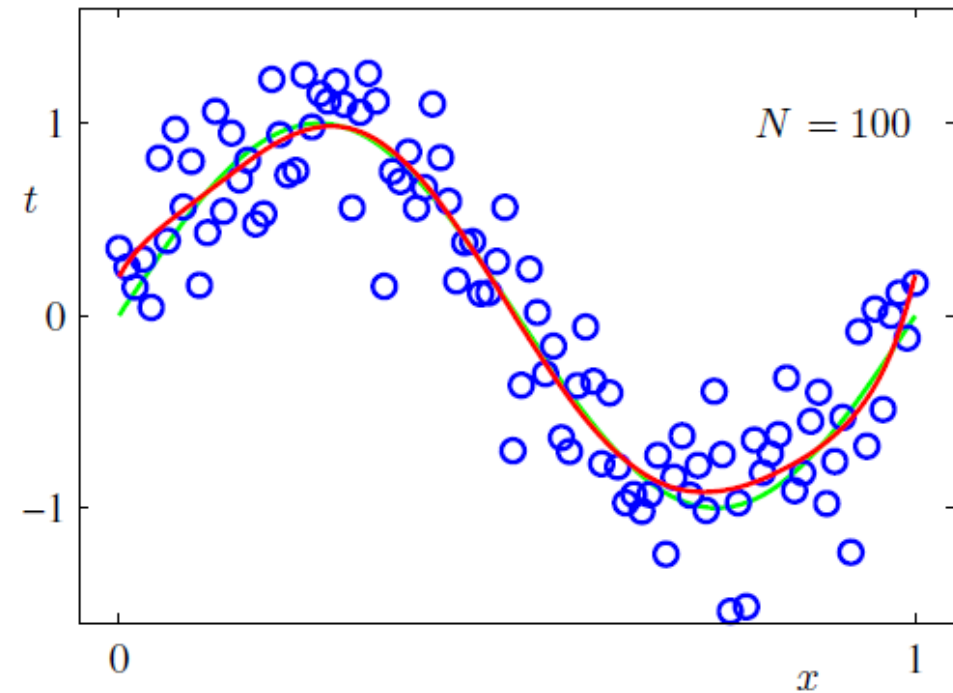
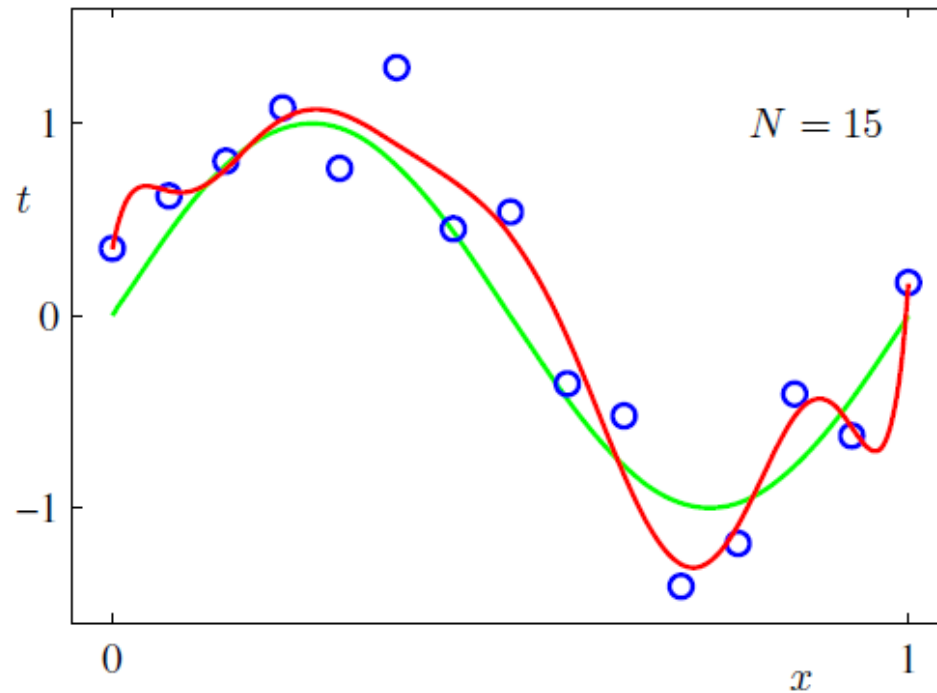
Overfitting

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

Tomadas de: Bishop, Christopher M. ( 2006). Pattern recognition and machine learning. New York. Springer.



# EFFECTO DE LA CANTIDAD DE DATOS N EN EL AJUSTE



$$y(x, \mathbf{w}) = \sum_{j=0}^9 w_j x^j$$



# CONCLUSIONES

- ▶ El modelo es lineal en los parámetros, no en los datos. Esto es importante porque la optimización se hace con respecto a los parámetros.
- ▶ A mayor número de pesos, mayor complejidad -> mayor flexibilidad. Pero, modelos muy complejos son propensos al sobreajuste con pocos datos
- ▶ Con pocos datos, se requiere un modelo poco complejo, pero que resulta poco flexible ante el incremento en el número de datos.
- ▶ Deseamos un compromiso entre la complejidad del modelo (número de parámetros) y su flexibilidad (capacidad de ajuste).
- ▶ El sobreajuste se da cuando los parámetros escalan a valores muy altos.
- ▶ Idea: tratar de penalizar valores altos de los parámetros en función de la complejidad => **Regularización**

# REGULARIZACIÓN

- La forma más simple es la penalización sobre la norma del vector de pesos (parámetros).

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Donde:

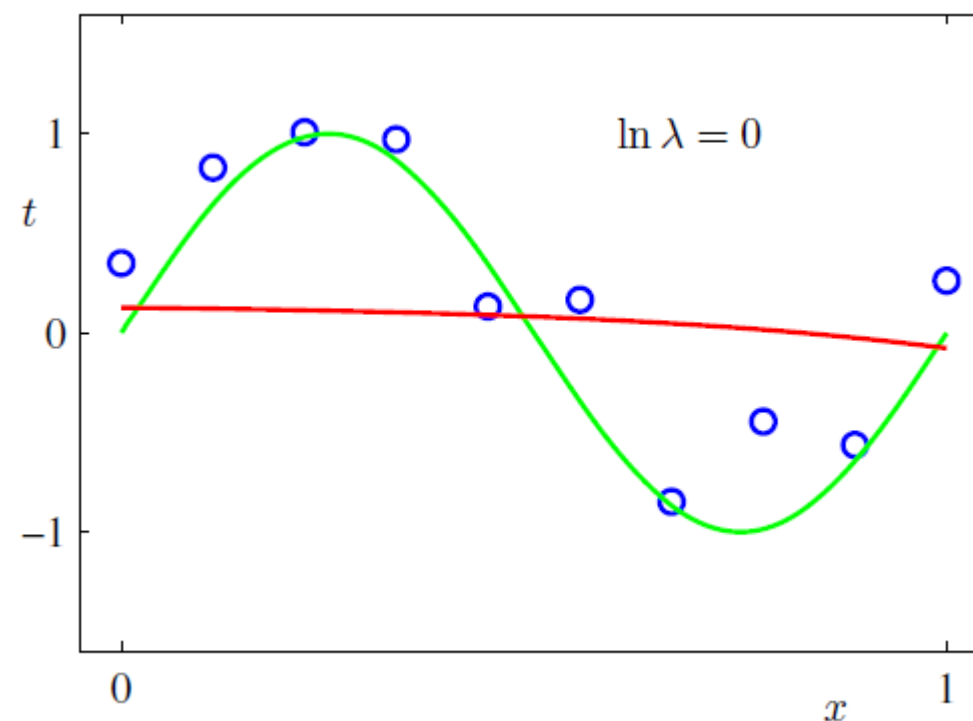
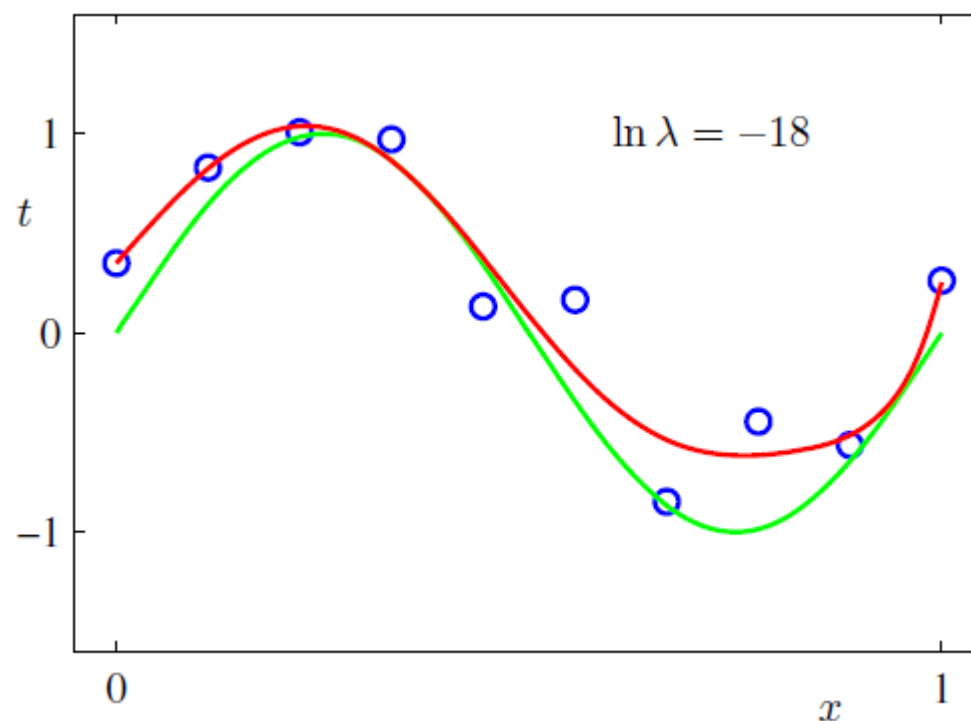
$$\|\mathbf{w}\|^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$$

# EFECTO DE LA REGULARIZACIÓN EN LOS PARÁMETROS

Tabla de pesos  $w$  para  $M = 9$  polinomios con varios valores para el parámetro de regularización  $\lambda$ . Nota que  $\ln \lambda = -\infty$  corresponde a un modelo sin regularización. Vemos que, a medida que el valor de  $\lambda$  aumenta, la magnitud típica de los pesos disminuye.

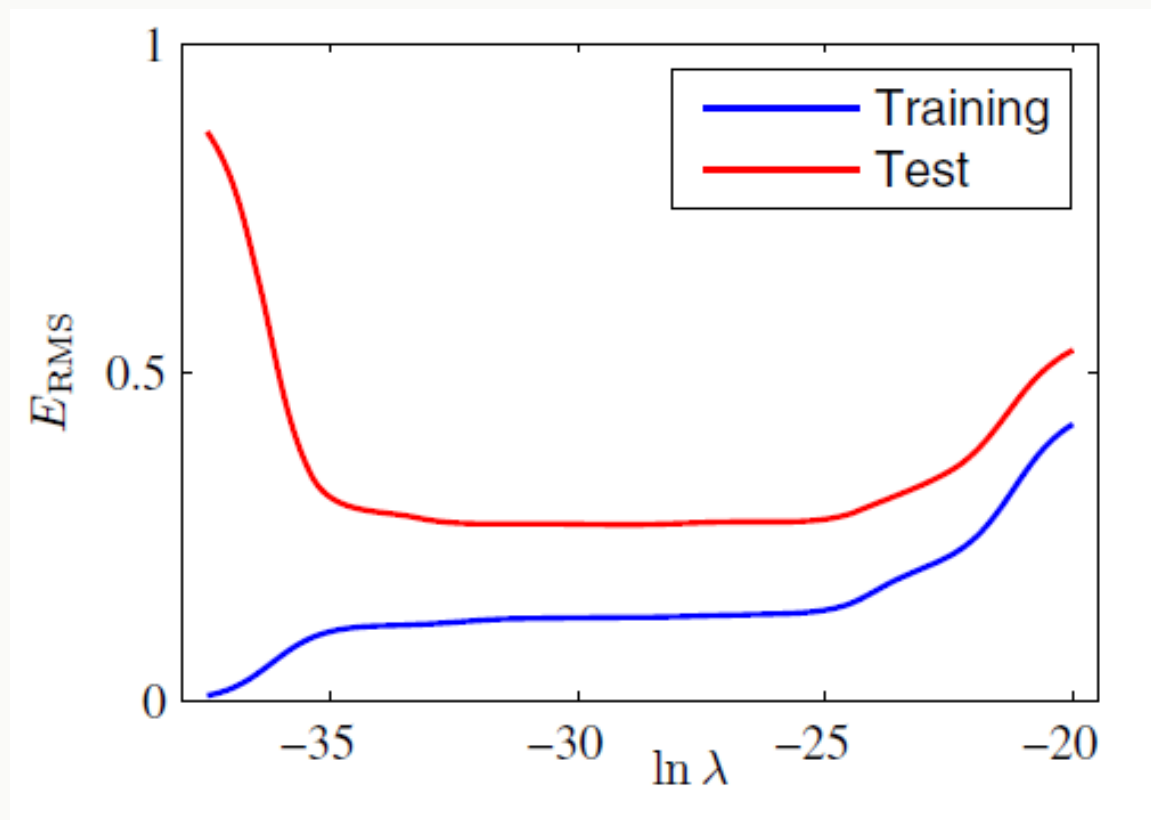
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01

# EFFECTO DE LA REGULARIZACIÓN EN EL AJUSTE



$M = 9$  con  $\ln \lambda = -18$  y  $\ln \lambda = 0$

# EFFECTO DE LA REGULARIZACIÓN EN EL DESEMPEÑO



# **FUNCIONES DE PÉRDIDA (COSTO)**

Nociones básicas

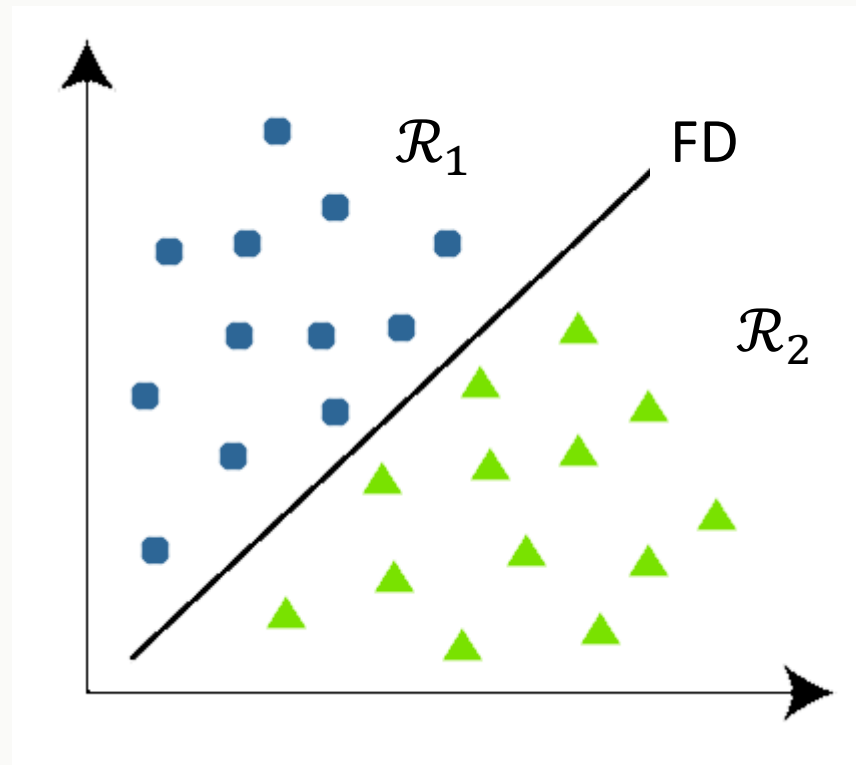
# ¿POR QUÉ FUNCIÓN DE “PÉRDIDA”?

- ▶ En ML, el propósito final se basa en minimizar o maximizar una función llamada "función objetivo".
- ▶ El grupo de funciones que se minimizan se denominan "funciones de pérdida".
- ▶ La función de pérdida se usa como medida de cuán bueno es un modelo de regresión/clasificación en términos de poder predecir el resultado esperado.



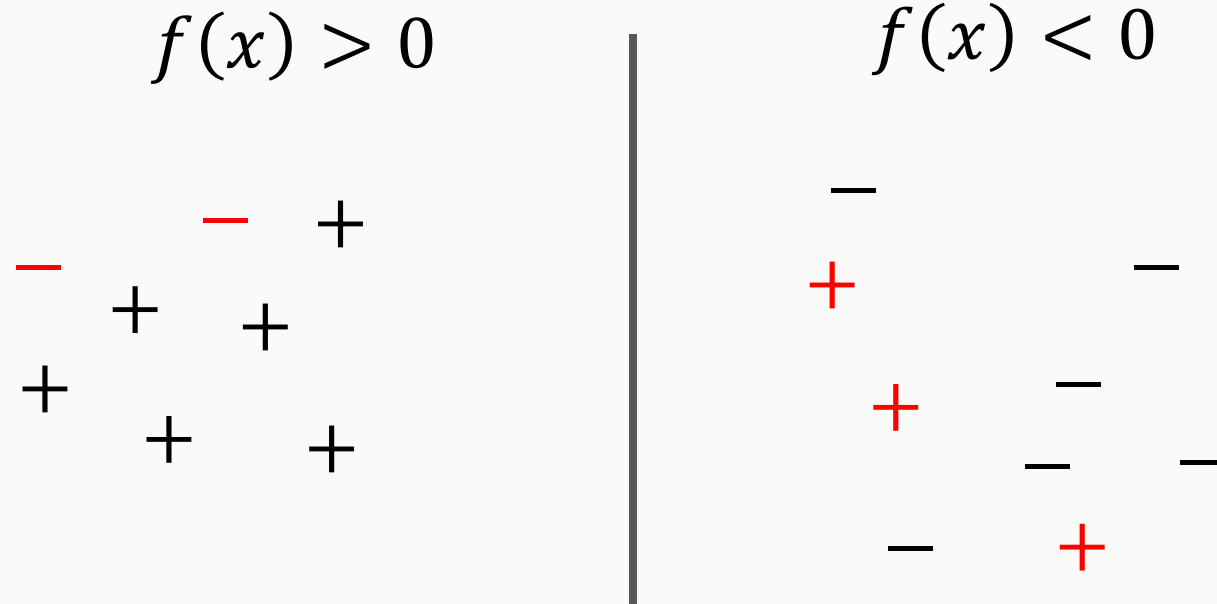
# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

- Un clasificador separa los datos de entrada en **regiones de decisión** cuyos límites llamamos **Fronteras (o superficies) de Decisión (FD)**.



# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

- Supongamos por ahora, que un punto está **bien clasificado** si  $\text{sign}(f(x_i))$  es  $y_i$ , donde  $f(\cdot)$  es la función de clasificación y  $y_i := \{-1, +1\}$ .



Fracción de veces que  $\text{sign}(f(x_i))$  no es  $y_i$ :

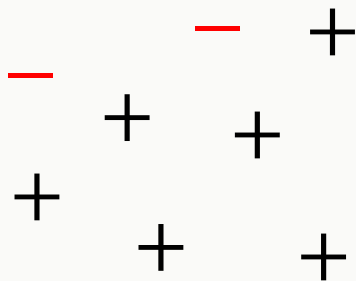
$$\frac{1}{n} \sum_{i=1}^n y_i \neq \text{sign}(f(x_i))$$

Minimizar esta función es computacionalmente muy difícil.

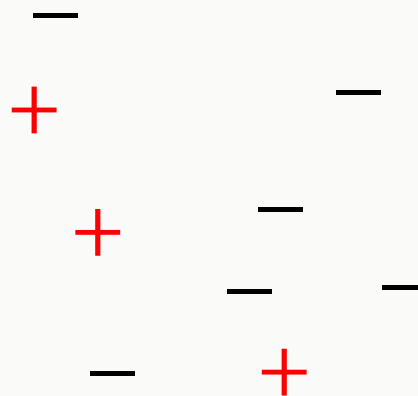
# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

- Supongamos por ahora, que un punto está **bien clasificado** si  $\text{sign}(f(x_i))$  es  $y_i$ , donde  $f(\cdot)$  es la función de clasificación y  $y_i := \{-1, +1\}$ .

$$yf(x) \approx 0$$



$$yf(x) \approx 0$$

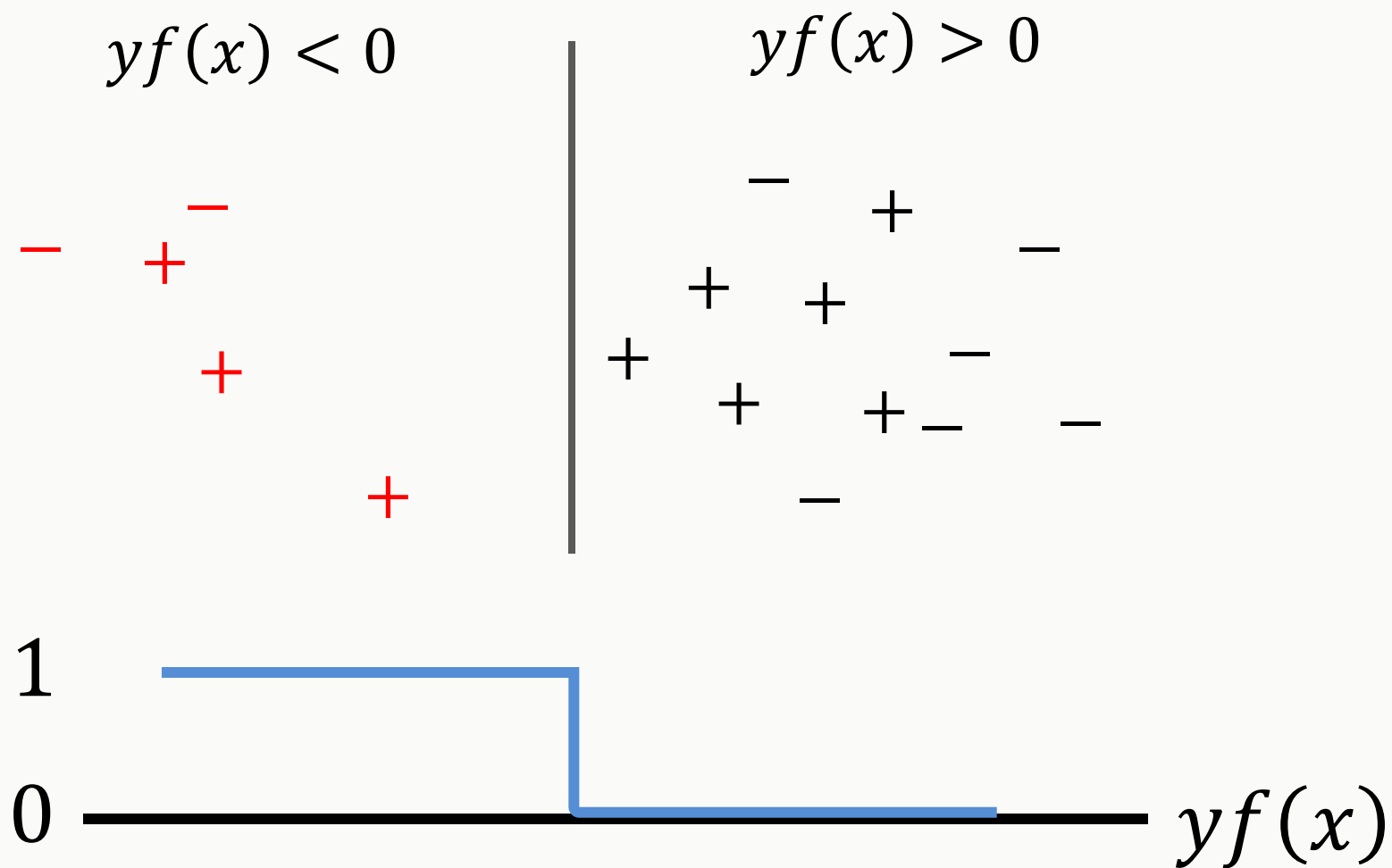


Fracción de veces que  $\text{sign}(f(x_i))$  no es  $y_i$ :

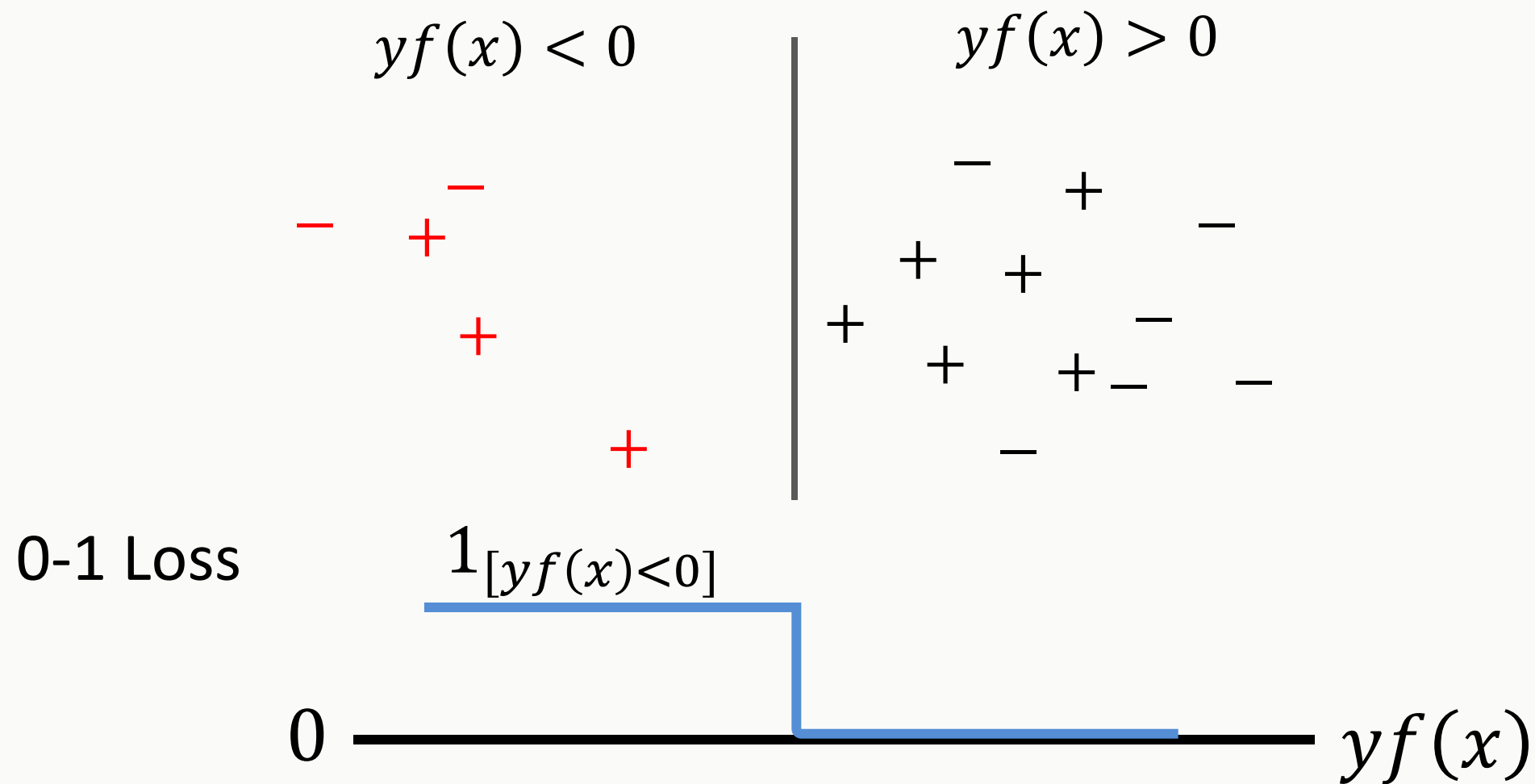
$$\frac{1}{n} \sum_{i=1}^n y_i \neq \text{sign}(f(x_i))$$

Minimizar esta función es computacionalmente muy difícil.

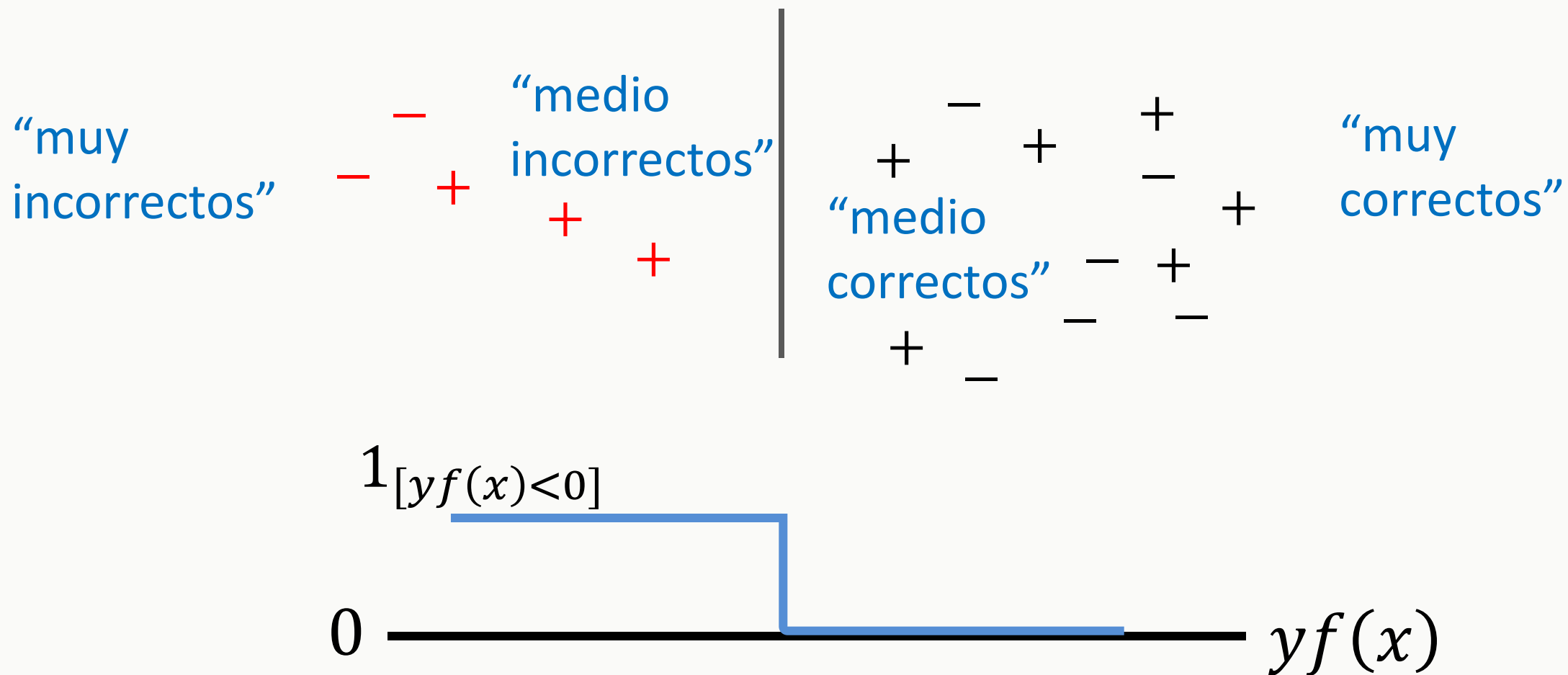
# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN



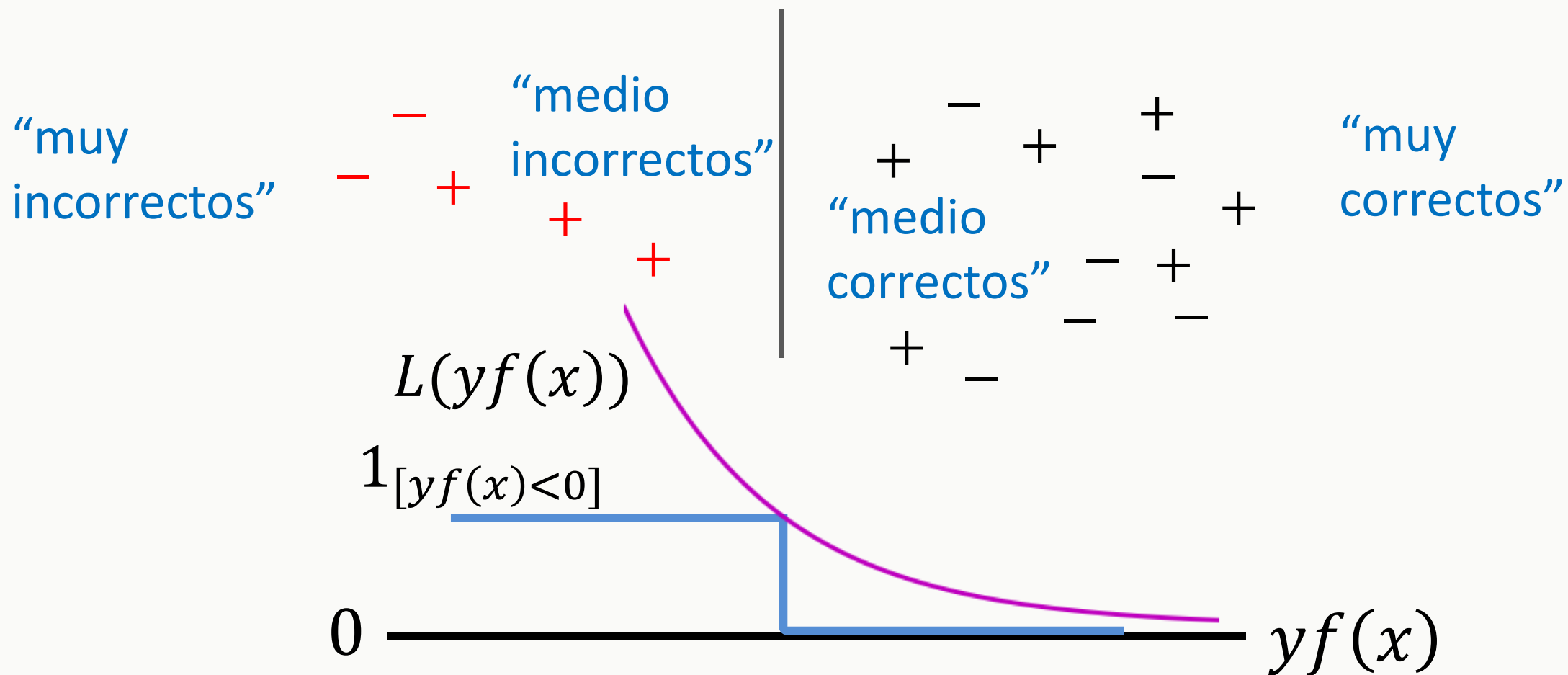
# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN



# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

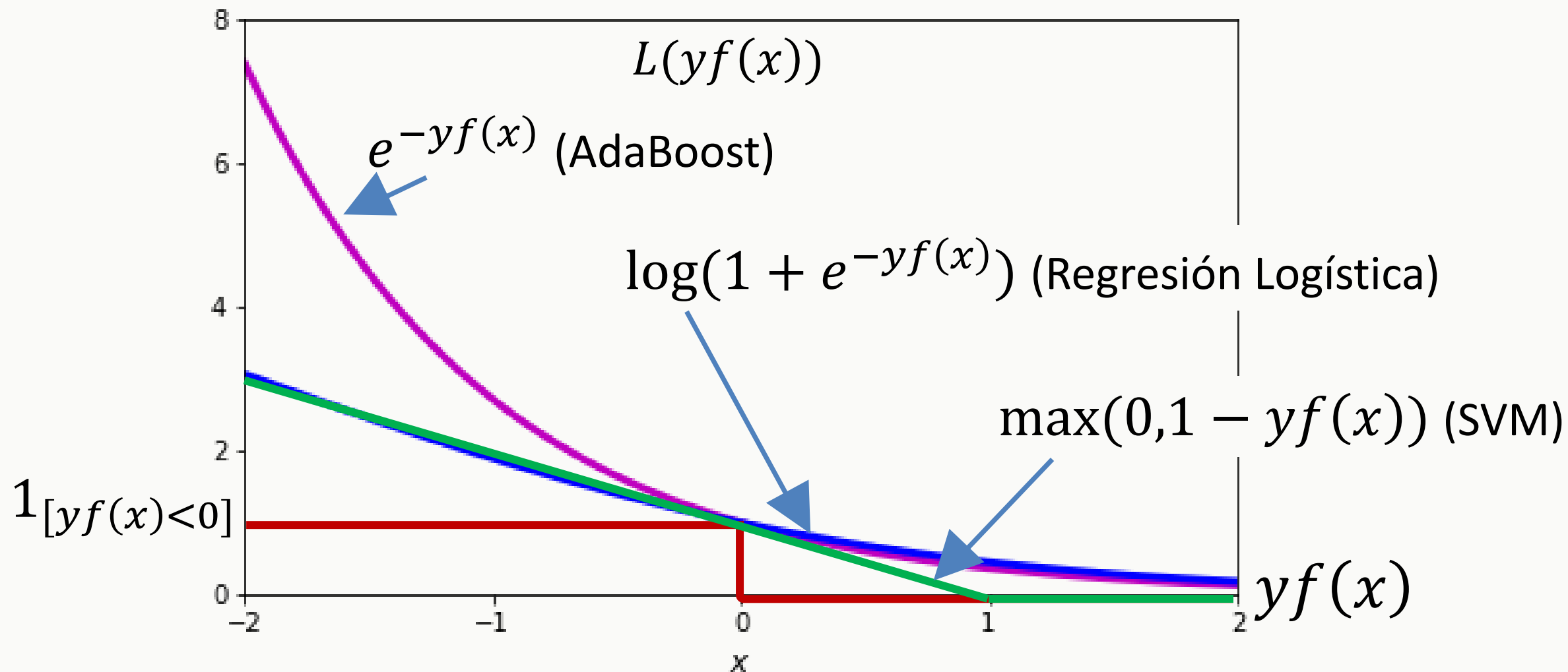


# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN





# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN



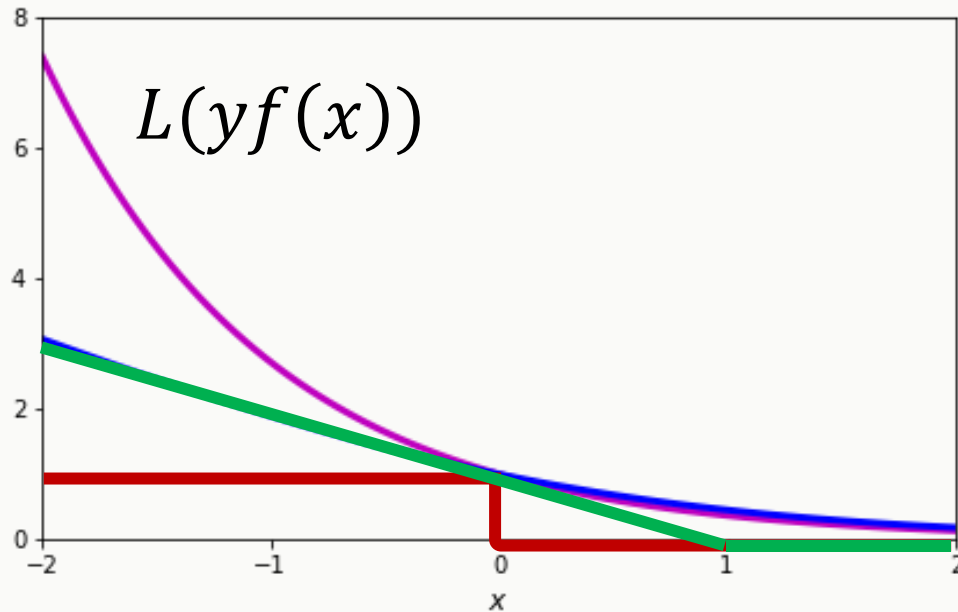
# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

Fracción de veces que  $\text{sign}(f(x_i))$  no es  $y_i$

$$= \frac{1}{n} \sum_{i=1}^n 1_{[y_i \neq \text{sign}(f(x_i))]}$$

$$= \frac{1}{n} \sum_{i=1}^n 1_{[y_i f(x_i) < 0]}$$

$$\leq \frac{1}{n} \sum_{i=1}^n L(y_i f(x_i))$$



# FUNCIÓN DE PÉRDIDA PARA CLASIFICACIÓN

Fracción de veces que  $\text{sign}(f(x_i))$  no es  $y_i$   $= \frac{1}{n} \sum_{i=1}^n 1_{[y_i \neq \text{sign}(f(x_i))]}$

Objetivo

$$\min_{\text{modelos } f} \frac{1}{n} \sum_{i=1}^n L(y_i f(x_i))$$

$$= \frac{1}{n} \sum_{i=1}^n 1_{[y_i f(x_i) < 0]}$$

$$\leq \frac{1}{n} \sum_{i=1}^n L(y_i f(x_i))$$

# DESCENSO DE GRADIENTE

Nociones básicas

# Ejemplo de problema

**Conjunto de  
entrenamiento** de precios  
de casas (Portland, Oregon,  
USA)

Tamaño en pies <sup>2</sup> (x)	Precio (\$) en miles (y)
2104	460
1416	232
1534	315
852	178
...	...

# Modelación del problema (I)

**Conjunto de  
entrenamiento** de precios  
de casas (Portland, Oregon,  
USA)

	Tamaño en pies <sup>2</sup> (x)	Precio (\$) en miles (y)
N=47	2104	460
	1416	232
	1534	315
	852	178
	...	...

Notación:

N : número de instancias (ejemplos) de entrenamiento

x : variable de “entrada” / **atributos** (*features*)

y : variable de “salida” / variable “objetivo”

# Modelación del problema (II)

**Conjunto de  
entrenamiento** de precios  
de casas (Portland, Oregon,  
USA)

Tamaño en pies <sup>2</sup> (x)	Precio (\$) en miles (y)
2104 $x^{(1)}$	460 $y^{(1)}$
1416 $x^{(2)}$	232
1534	315
852	178
...	...

N=47

Notación:

N : número de instancias (ejemplos) de entrenamiento

x : variable de “entrada” / **atributos** (*features*)

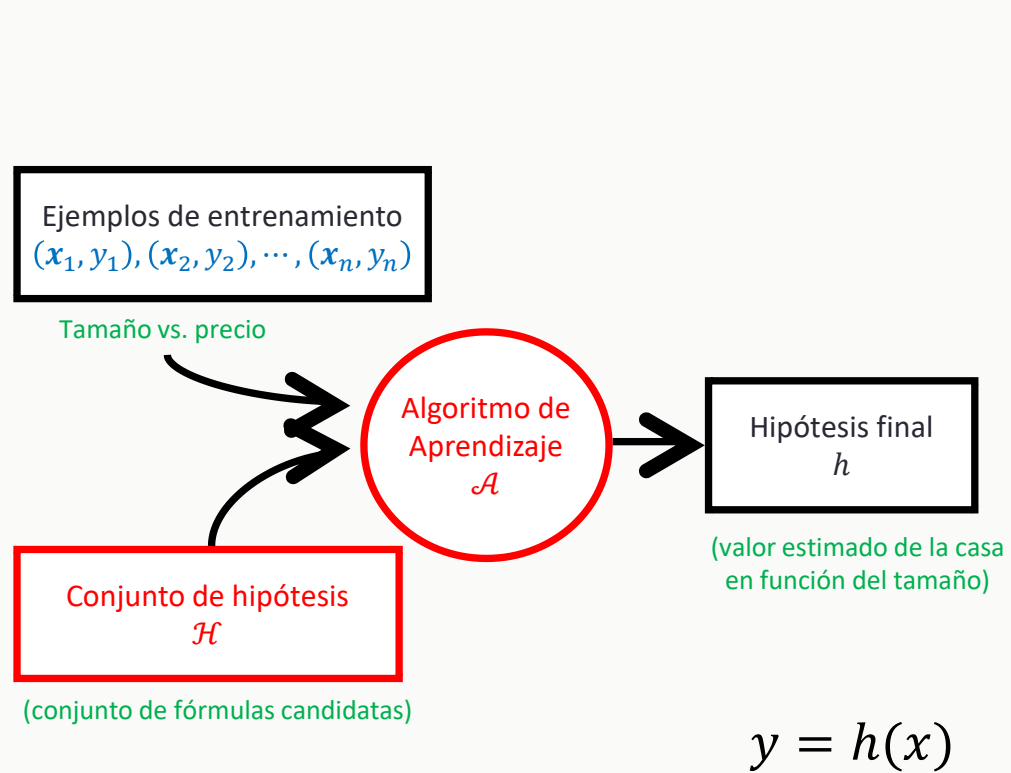
y : variable de “salida” / variable “objetivo”

$(x, y)$  – un ejemplo de  
entrenamiento

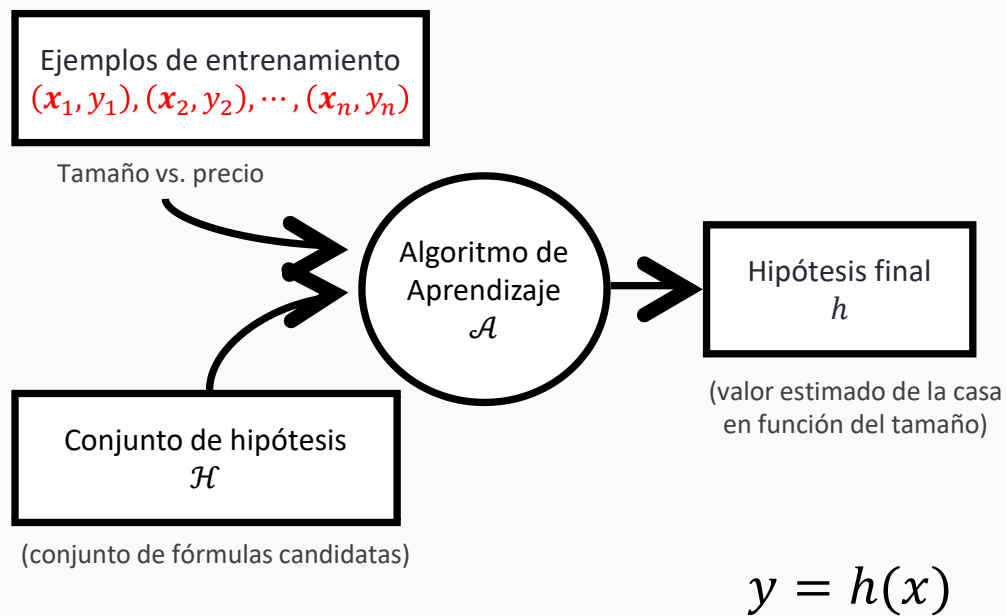
$(x^{(i)}, y^{(i)})$  – i-ésimo ejemplo de  
entrenamiento



# Problema de Aprendizaje

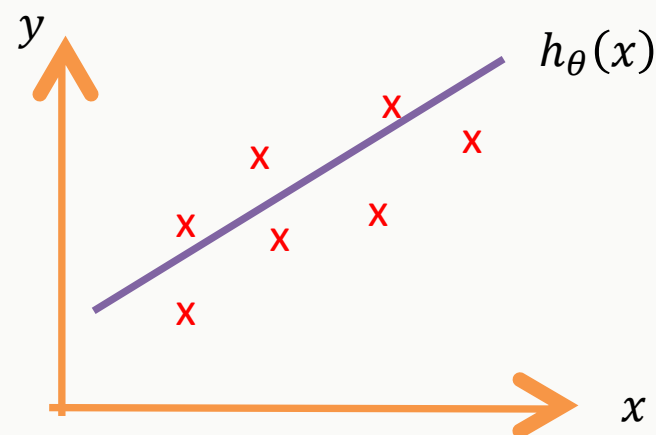


# Problema de Aprendizaje



¿Cómo representamos  $h$ ?

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



Regresión lineal de una variable (univariada).

# FUNCIÓN DE COSTO

---

# Problema de aprendizaje

Conjunto de entrenamiento

	Tamaño en pies <sup>2</sup> (x)	Precio (\$) en miles (y)
N=47	2104	460
	1416	232
	1534	315
	852	178
	...	...

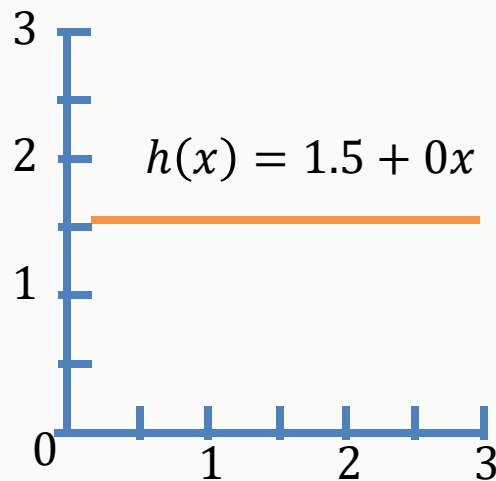
Hipótesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

$\theta_i$ 's : Parámetros

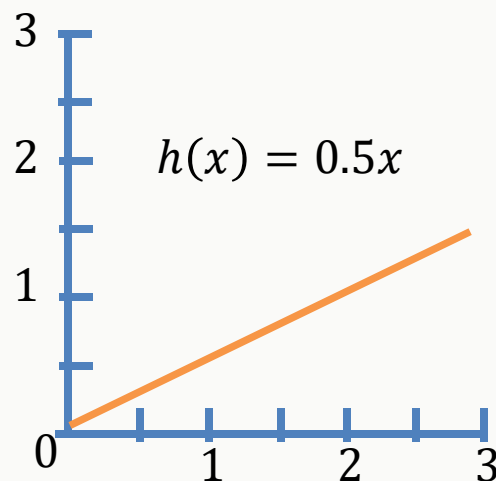
¿Cómo elegir los  $\theta_i$ 's?

# Hipótesis en función de los parámetros

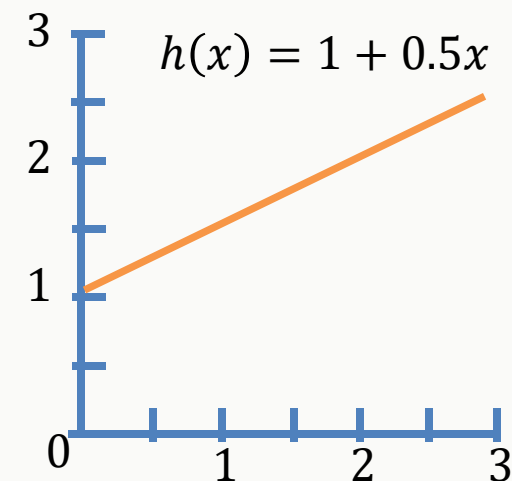
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



$$\begin{aligned}\theta_0 &= 1.5 \\ \theta_1 &= 0\end{aligned}$$

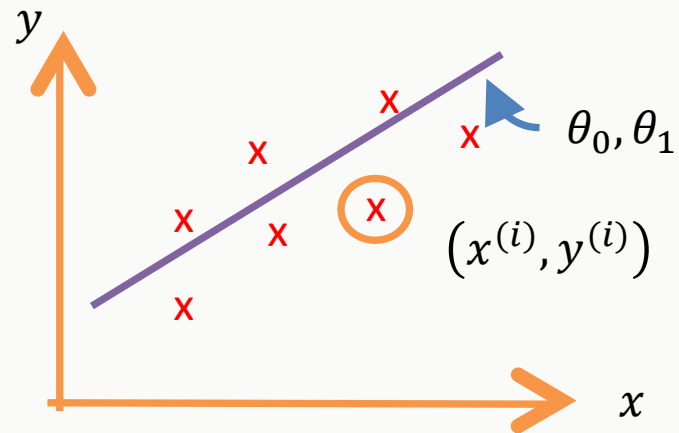


$$\begin{aligned}\theta_0 &= 0 \\ \theta_1 &= 0.5\end{aligned}$$



$$\begin{aligned}\theta_0 &= 1 \\ \theta_1 &= 0.5\end{aligned}$$

# Función de costo

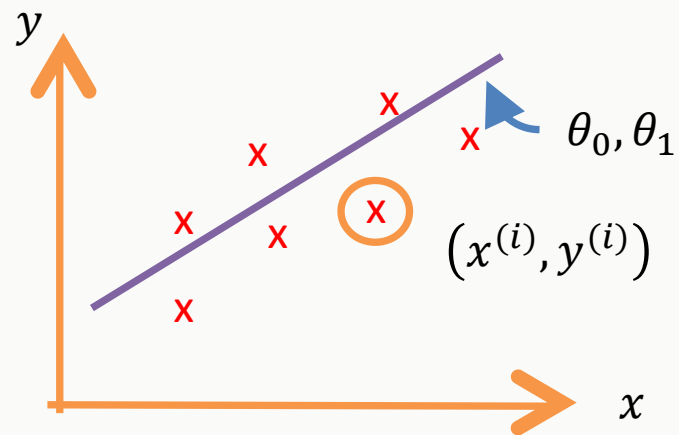


Idea: Elige  $\theta_0, \theta_1$  de tal forma que  $h_\theta(x)$  esté cerca de  $y$  para nuestros ejemplos de entrenamiento  $(x, y)$

Función de error cuadrático:

$$\Rightarrow J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^N \underbrace{(h_\theta(x^{(i)}) - y^{(i)})^2}_{h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}}$$

# Minimización de costo



Idea: Elige  $\theta_0, \theta_1$  de tal forma que  $h_\theta(x)$  esté cerca de  $y$  para nuestros ejemplos de entrenamiento  $(x, y)$

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)})^2$$

Función de costo

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

$$\min_{\theta_0, \theta_1} \frac{1}{2N} \sum_{i=1}^N (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$



# FUNCIÓN DE COSTO

---

Intuición I

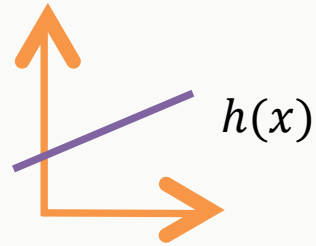
# Problema de minimización

Hipótesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Parámetros:

$$\theta_0, \theta_1$$



Función de costo:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Objetivo:  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

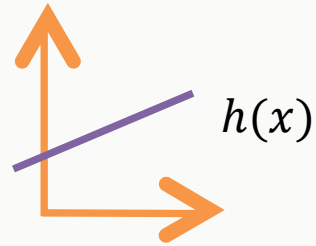
# Minimización simplificada

Hipótesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Parámetros:

$$\theta_0, \theta_1$$



Función de costo:

$$J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

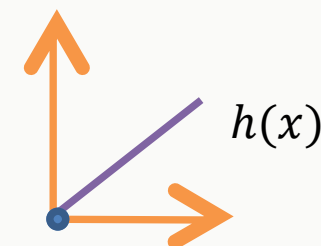
Objetivo:  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Simplificación:

$$h_{\theta}(x) = \theta_1 x$$

Parámetros:

$$\theta_1$$

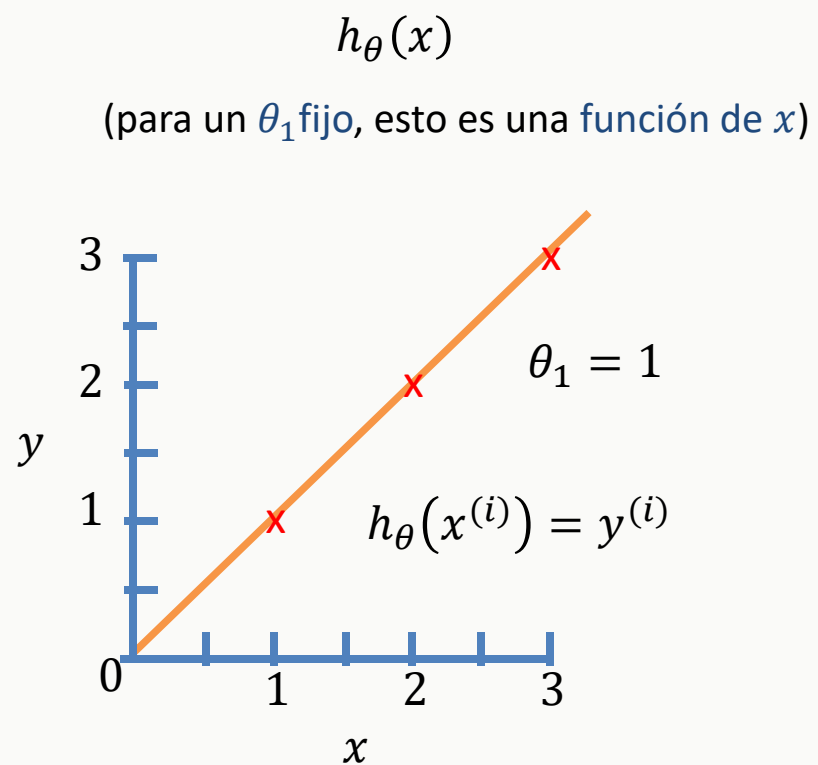


Función de costo:

$$J(\theta_1) = \frac{1}{2N} \sum_{i=1}^N \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\theta_1 x^{(i)}}$$

Objetivo:  $\min_{\theta_1} J(\theta_1)$

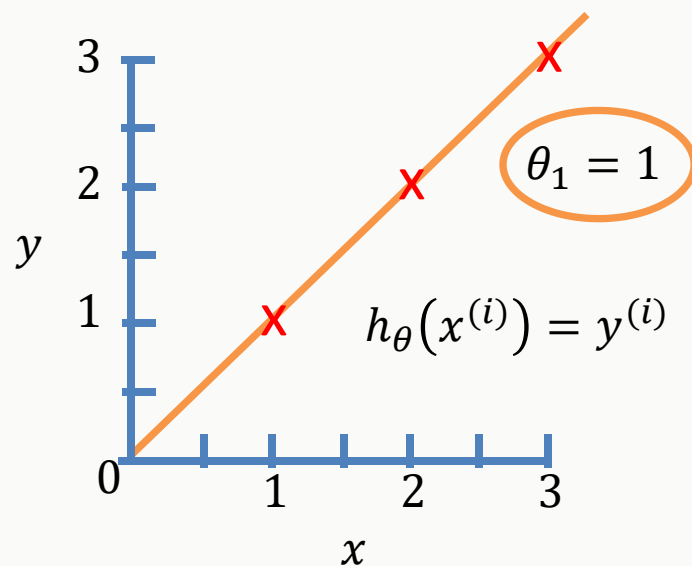
# Hipótesis vs. Costo



# Hipótesis vs. Costo

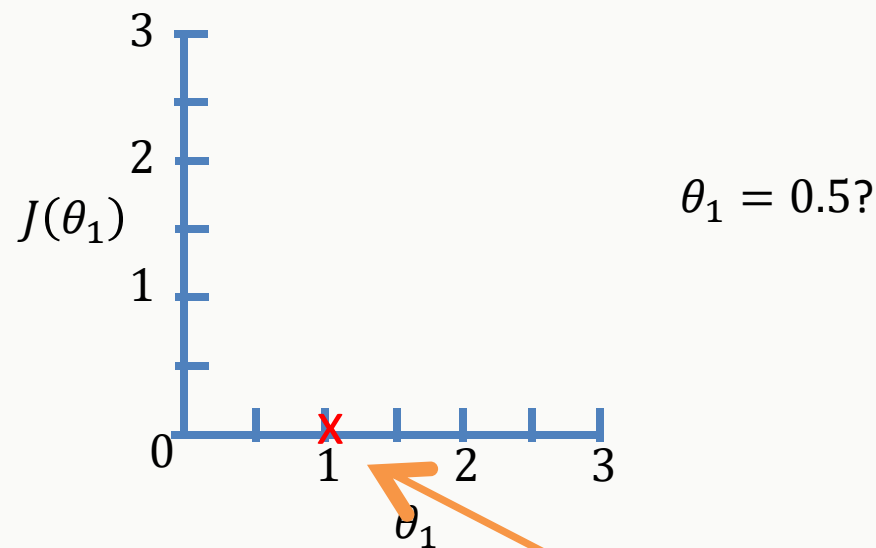
$h_{\theta}(x)$

(para un  $\theta_1$  fijo, esto es una función de  $x$ )



$J(\theta_1)$

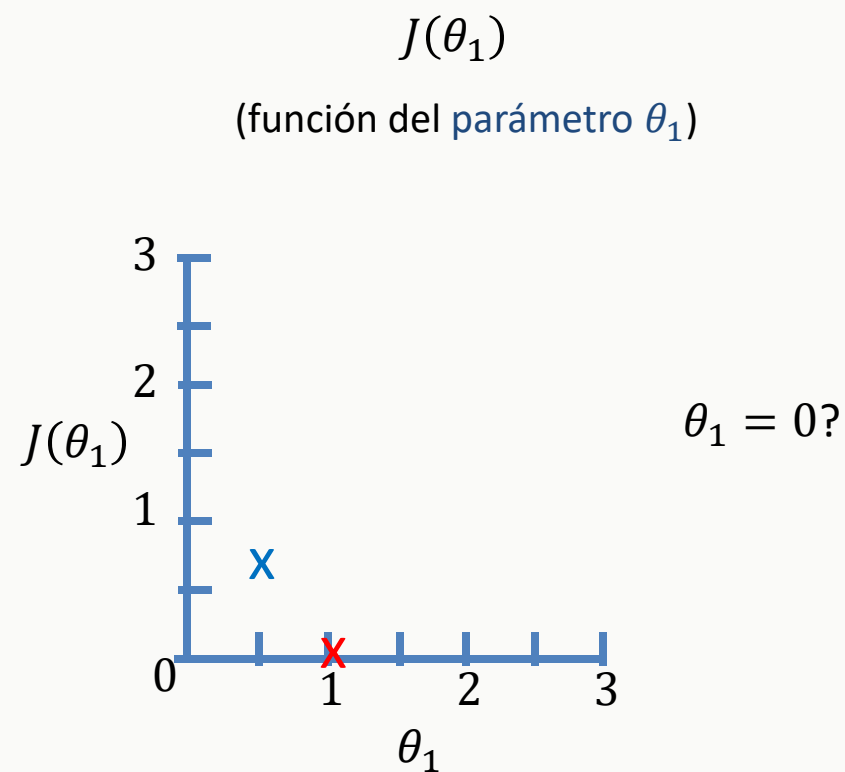
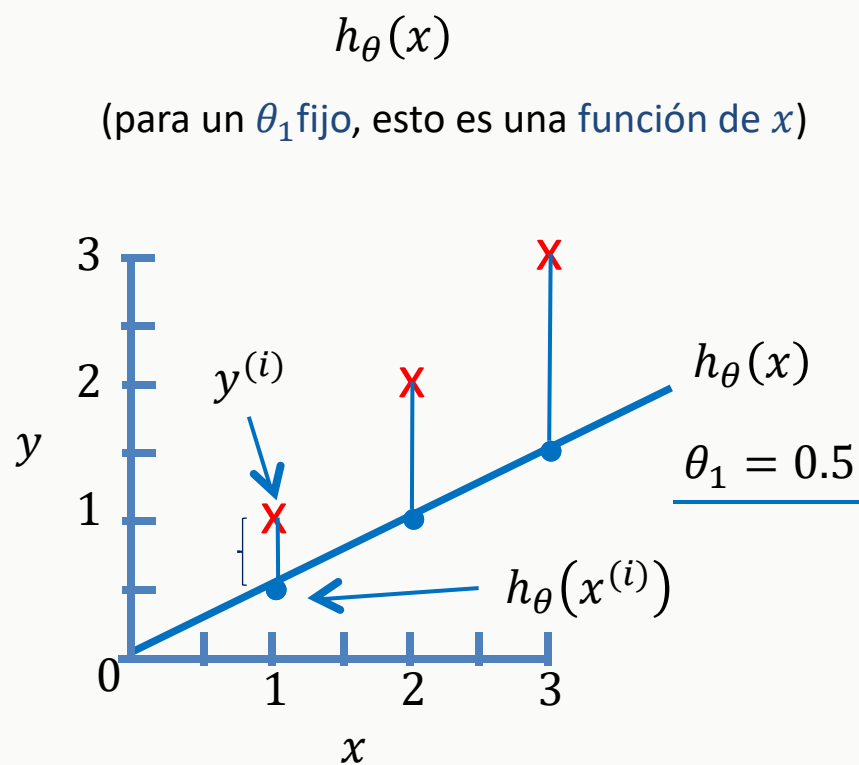
(función del parámetro  $\theta_1$ )



$$J(\theta_1) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2N} \sum_{i=1}^N (\theta_1 x^{(i)} - y^{(i)})^2$$

$$J(1) = \frac{1}{2 \cdot 3} (0^2 + 0^2 + 0^2) = 0$$

# Hipótesis vs. Costo

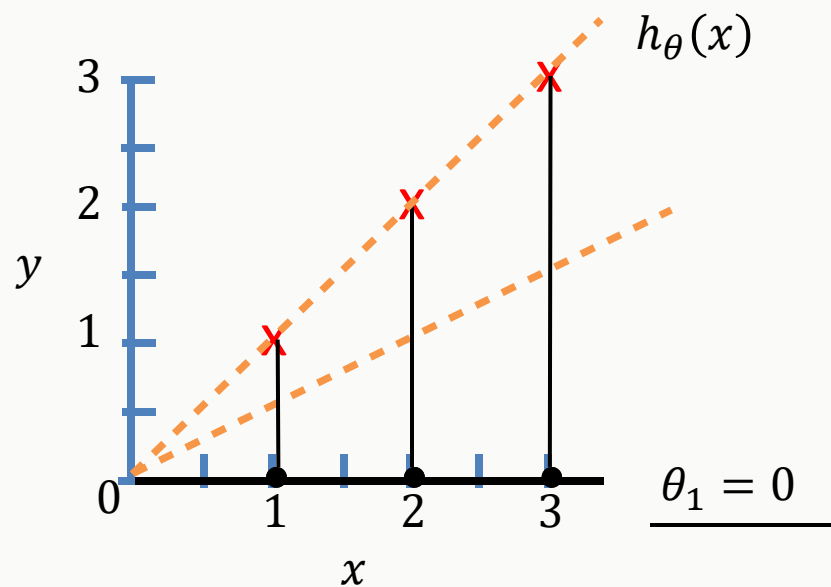


$$J(0.5) = \frac{1}{2 \cdot 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \cong 0.58$$

# Hipótesis vs. Costo

$h_{\theta}(x)$

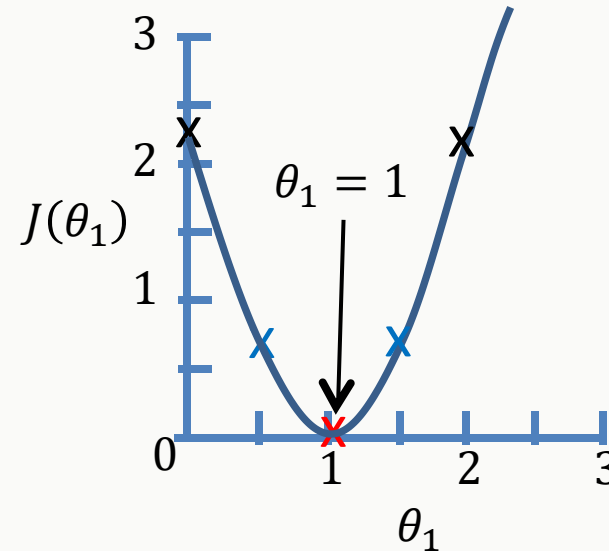
(para un  $\theta_1$  fijo, esto es una función de  $x$ )



$$J(0) = \frac{1}{2 \cdot 3} [(1)^2 + (2)^2 + (3)^2] \cong 2.3$$

$J(\theta_1)$

(función del parámetro  $\theta_1$ )



# FUNCIÓN DE COSTO

---

## Intuición II



# Problema de minimización de costo

Hipótesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parámetros:  $\theta_0, \theta_1$

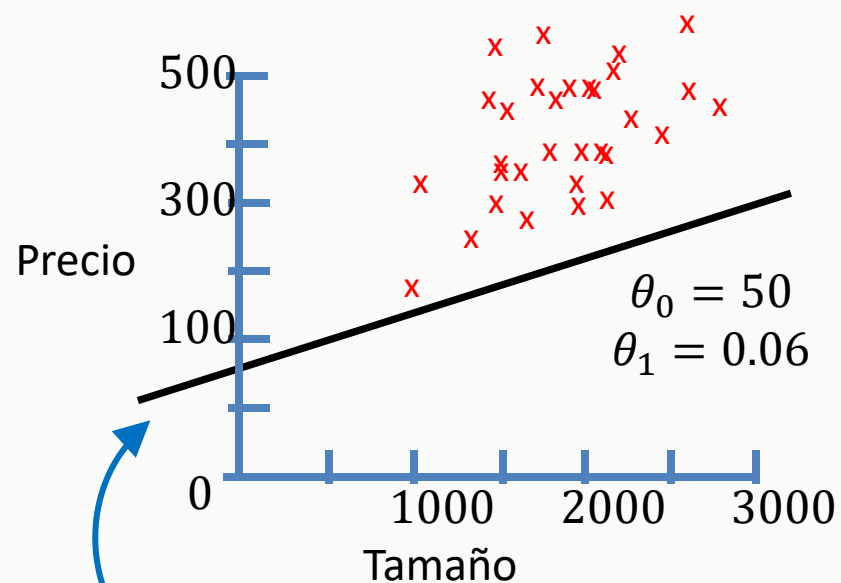
Función de costo:  $J(\theta_0, \theta_1) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Objetivo:  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

# Hipótesis vs. Costo

$$h_{\theta}(x)$$

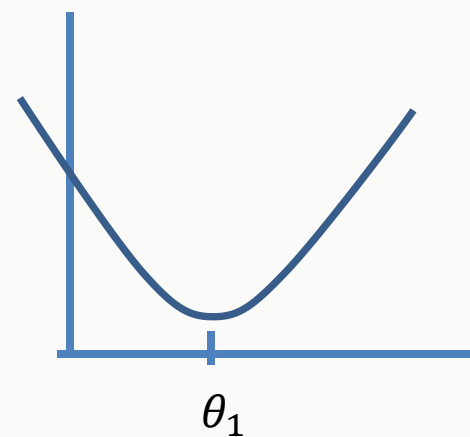
(para un  $\theta_1, \theta_2$  fijos, esto es una función de  $x$ )



$$h_{\theta}(x) = 50 + 0.06x$$

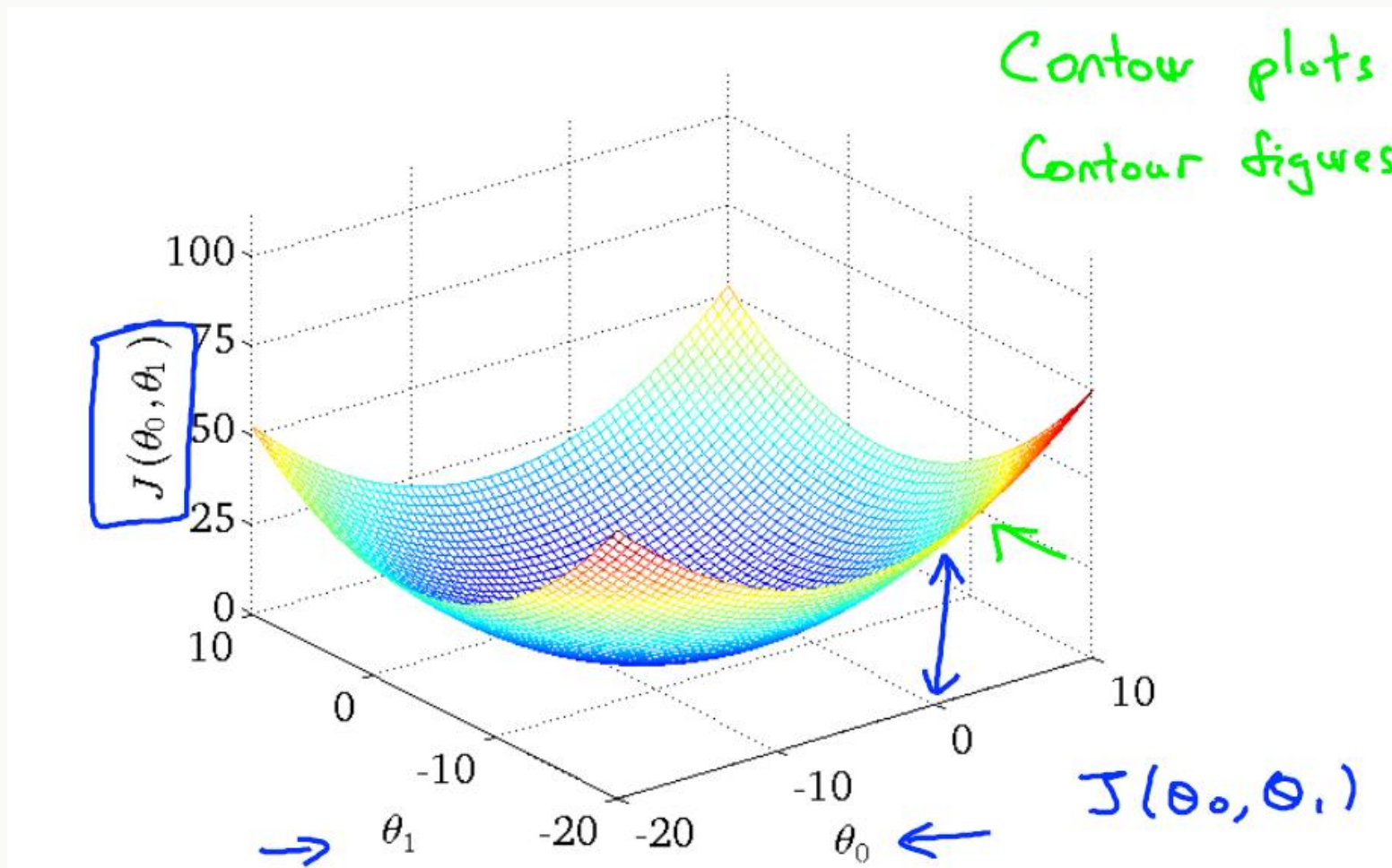
$$J(\theta_0, \theta_1)$$

(función de los parámetros  $\theta_0, \theta_1$ )



$$\theta_0, \theta_1$$

# Gráficos de contorno



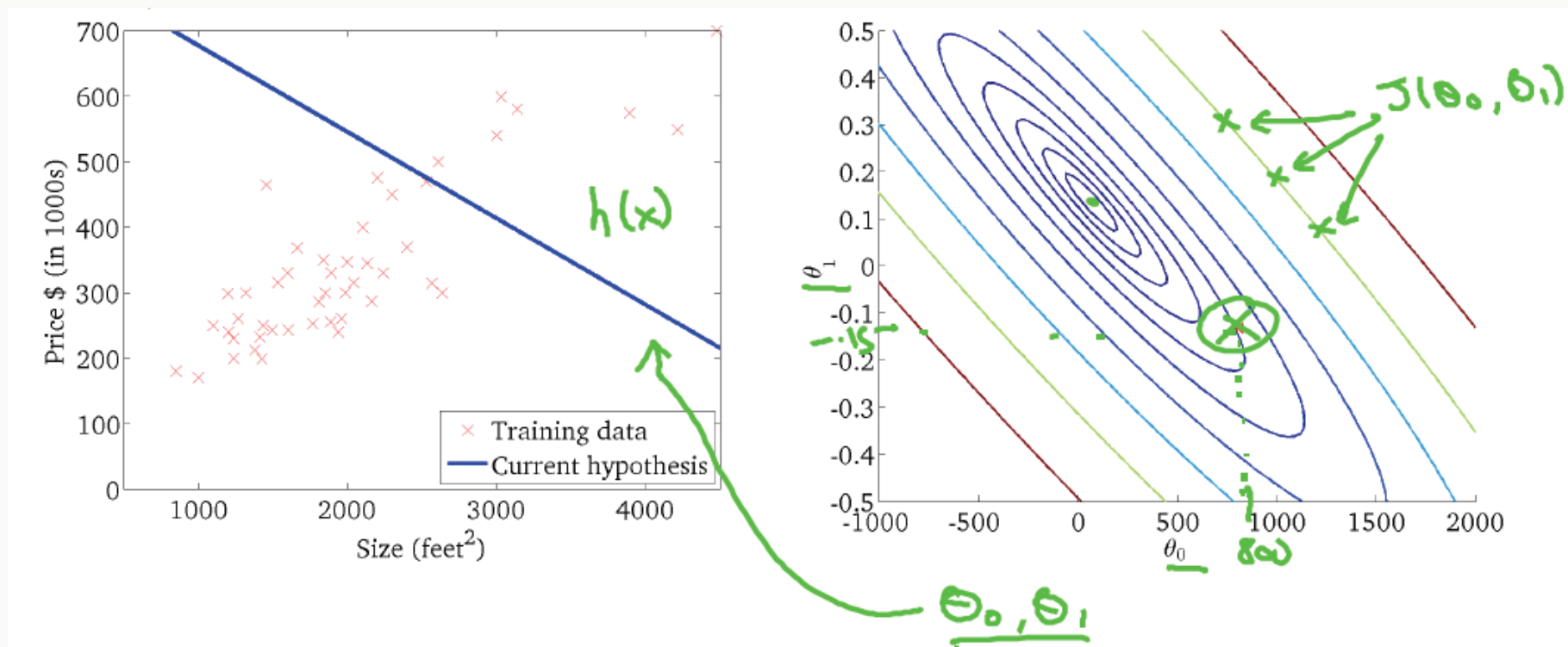
# Hipótesis vs. Costo

$$h_{\theta}(x)$$

(para un  $\theta_1, \theta_2$  fijos, esto es una función de  $x$ )

$$J(\theta_0, \theta_1)$$

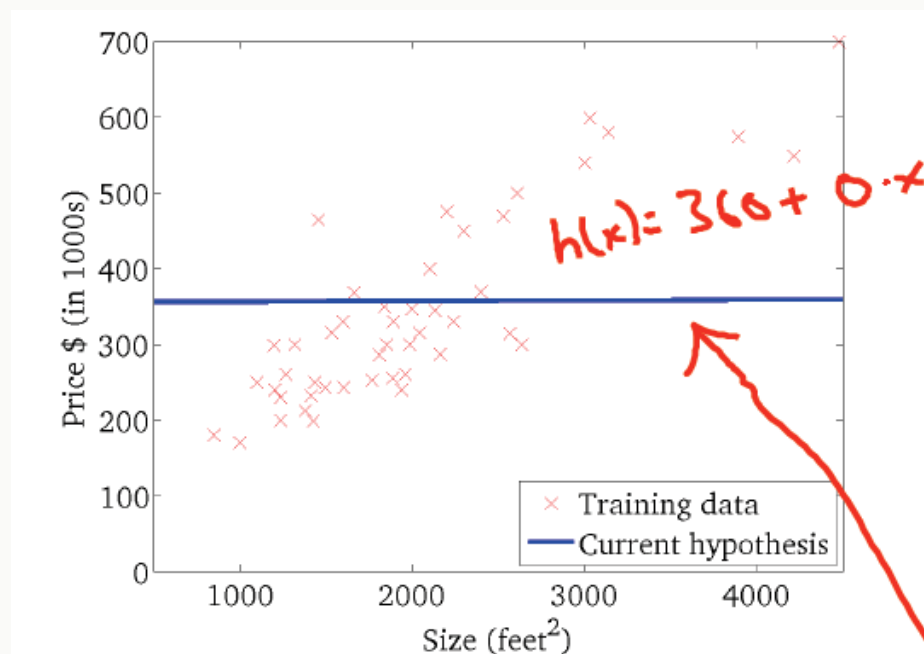
(función de los parámetros  $\theta_0, \theta_1$ )



# Hipótesis vs. Costo

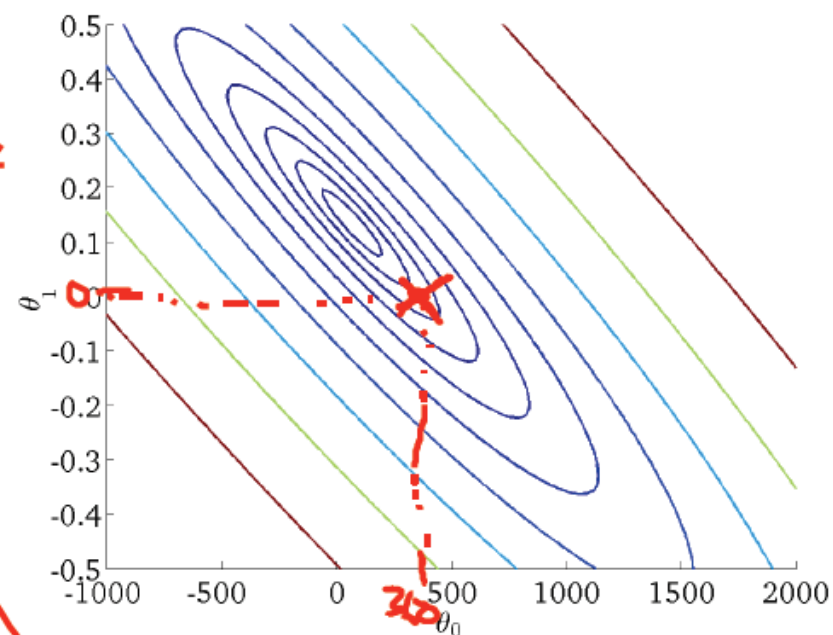
$$h_{\theta}(x)$$

(para un  $\theta_1, \theta_2$  fijos, esto es una función de  $x$ )



$$J(\theta_0, \theta_1)$$

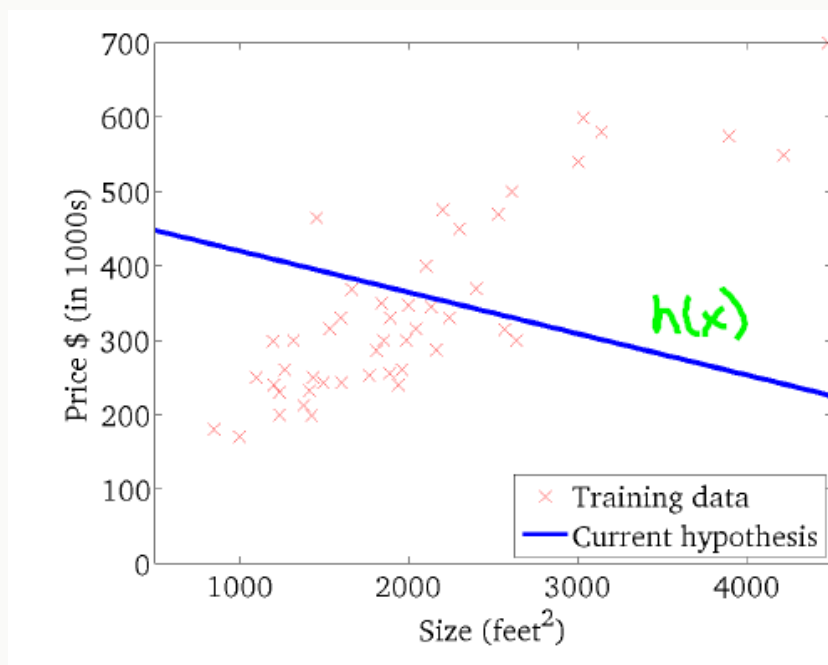
(función de los parámetros  $\theta_0, \theta_1$ )



# Hipótesis vs. Costo

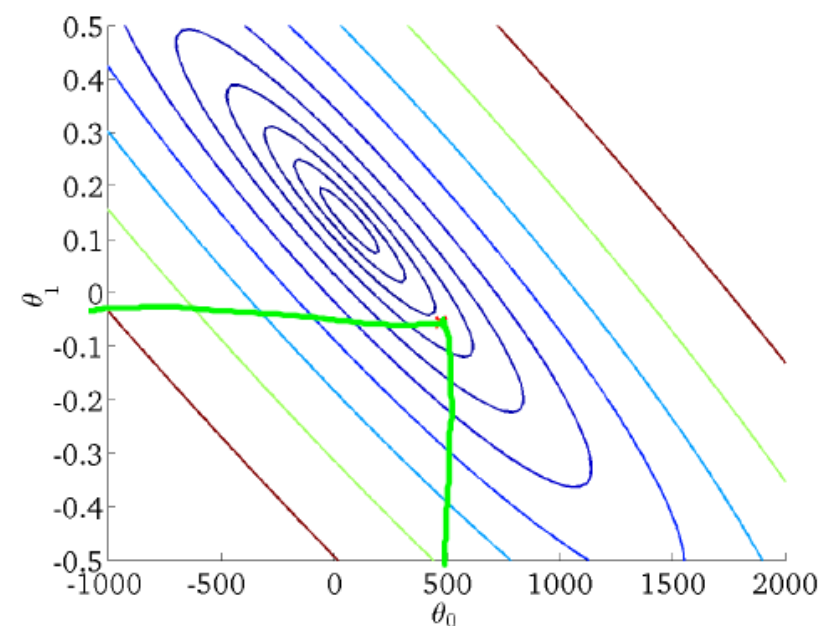
$$h_{\theta}(x)$$

(para un  $\theta_1, \theta_2$  fijos, esto es una función de  $x$ )



$$J(\theta_0, \theta_1)$$

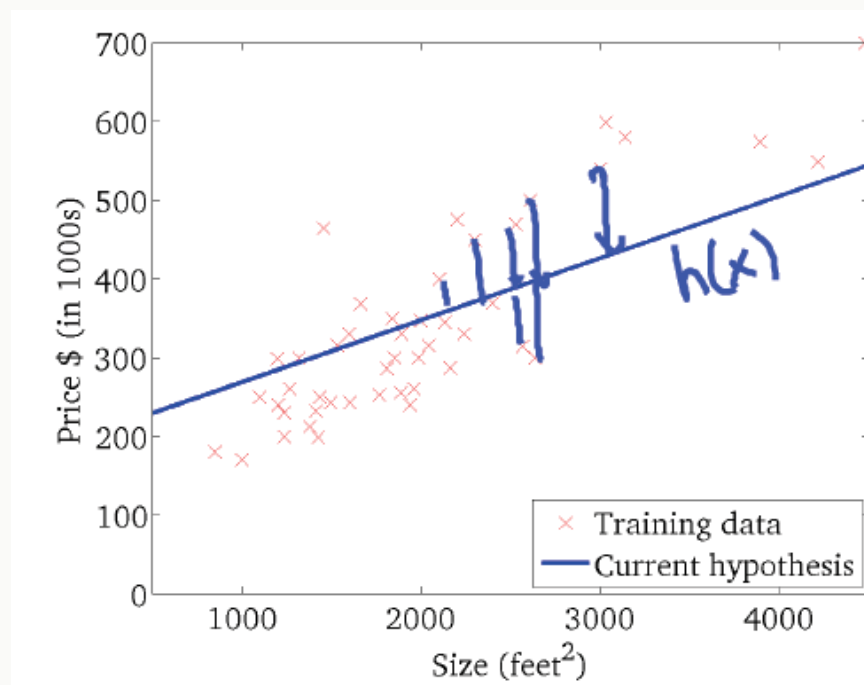
(función de los parámetros  $\theta_0, \theta_1$ )



# Hipótesis vs. Costo

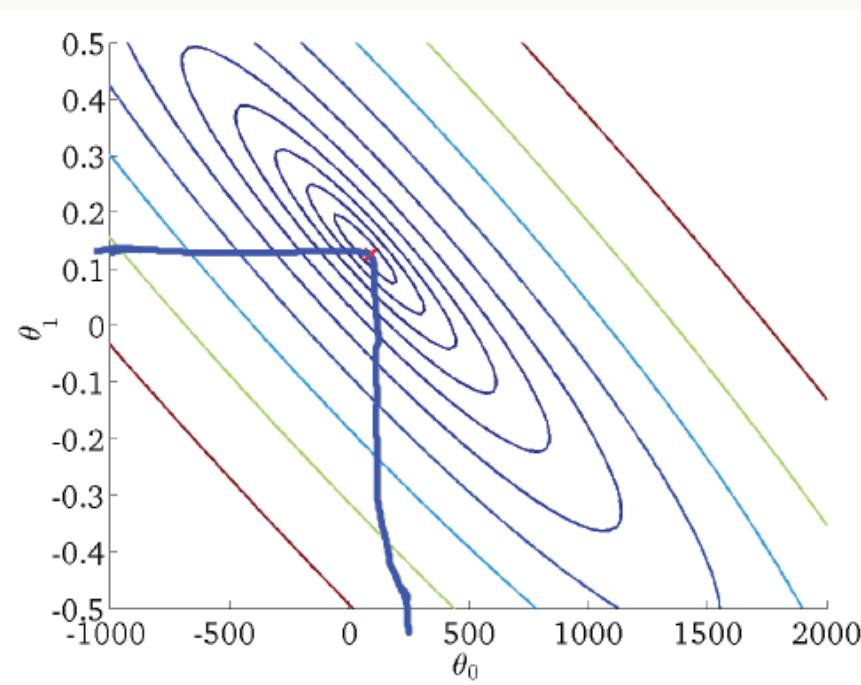
$$h_{\theta}(x)$$

(para un  $\theta_1, \theta_2$  fijos, esto es una función de  $x$ )



$$J(\theta_0, \theta_1)$$

(función de los parámetros  $\theta_0, \theta_1$ )



# DESCENSO DE GRADIENTE

---



# Problema de minimización

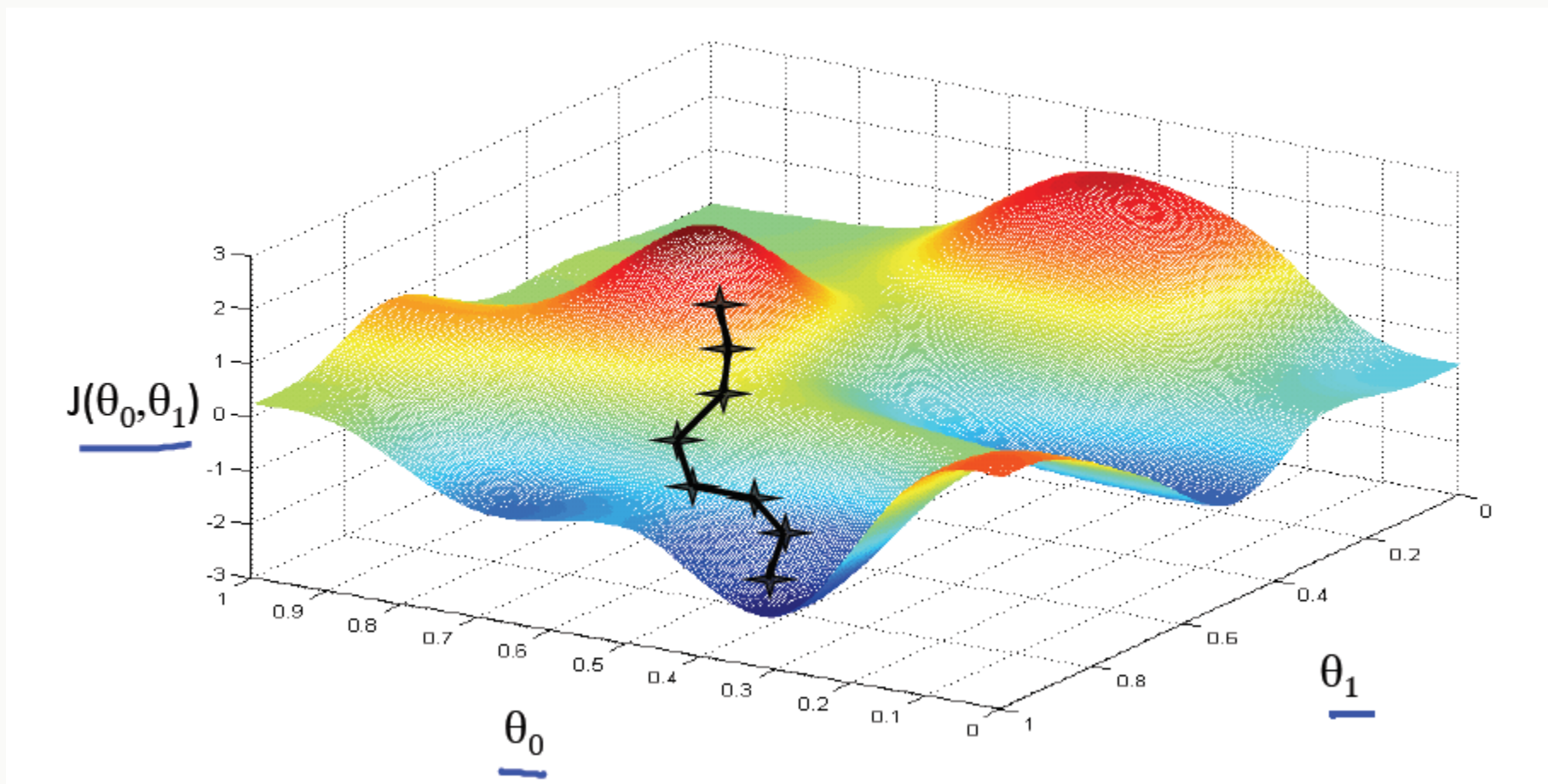
Función de costo:  $J(\theta_0, \theta_1)$

Objetivo:  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

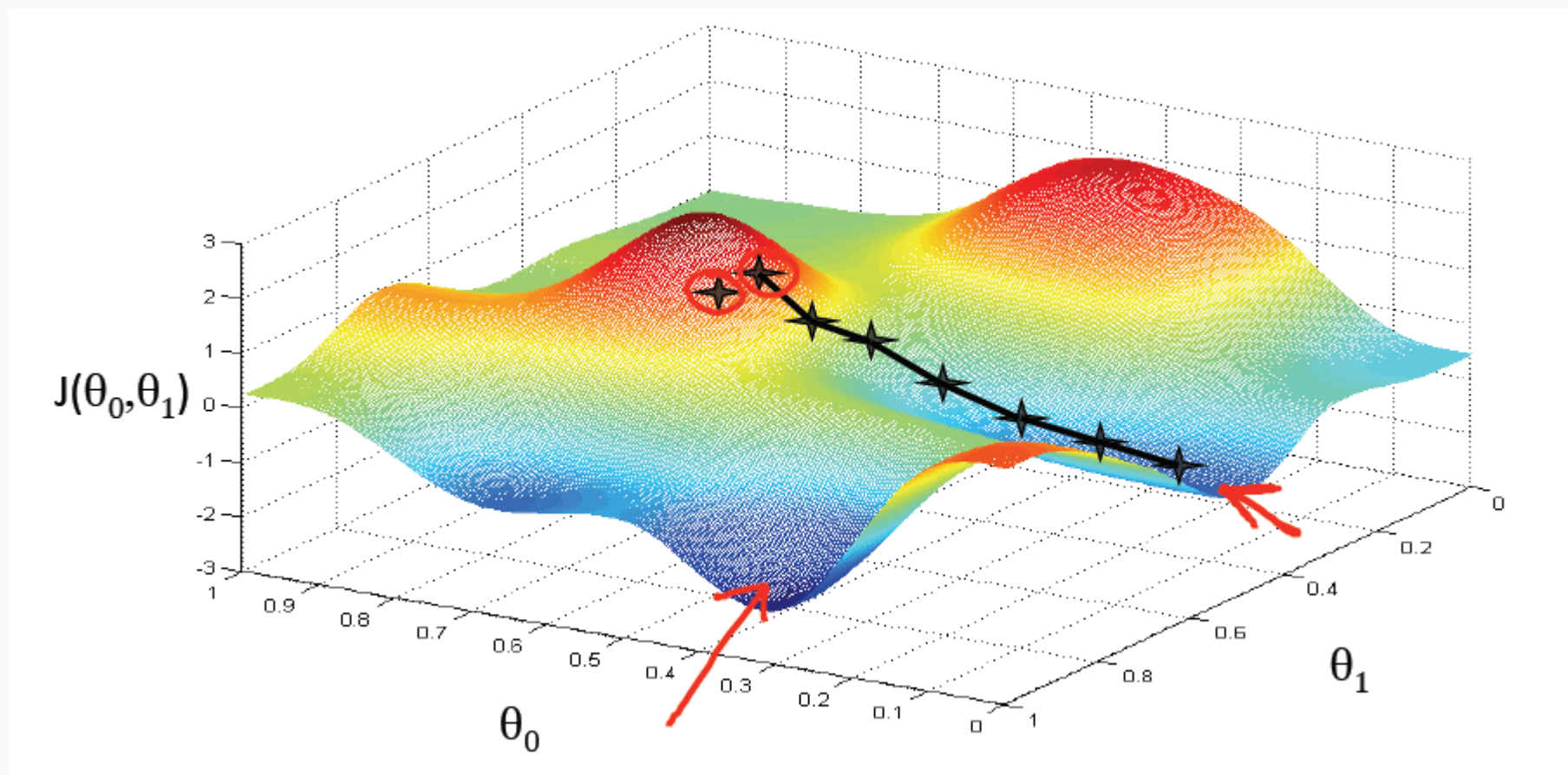
Descripción resumida:

- Comienza con algún  $\theta_0, \theta_1$
- Cambia  $\theta_0, \theta_1$  para reducir  $J(\theta_0, \theta_1)$  hasta que (ojalá) lleguemos a un mínimo.

# Topología de la función de costo



# Topología de la función de costo



# Algoritmo de descenso de gradiente

Repite hasta la convergencia {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{para } j = 0 \text{ y } j = 1)$$

}

Asignación simultánea!!!

Correcto: actualización simultánea

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrecto:

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp1}$$

# INTUICIÓN DEL DESCENSO DE GRADIENTE

---

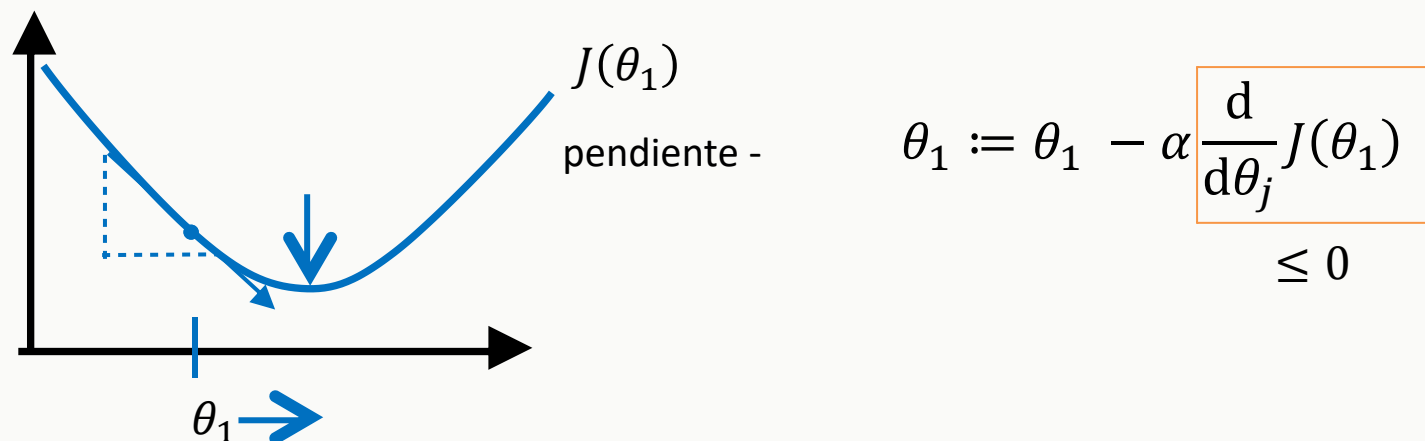
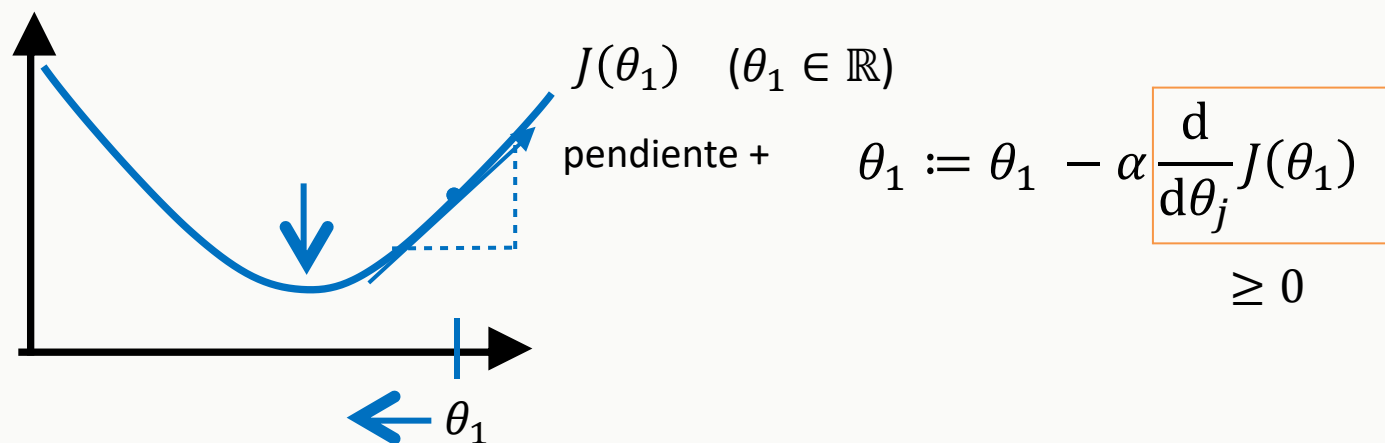
# Algoritmo de descenso de gradiente

Repite hasta la convergencia {

$$\underbrace{\theta_j}_{-} := \underbrace{\theta_j}_{-} - \underbrace{\alpha}_{\text{tasa de aprendizaje}} \underbrace{\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)}_{\text{derivada}} \quad \text{(actualiza simultáneamente } j = 0 \text{ y } j = 1)$$

$$\min_{\theta_1} J(\underbrace{\theta_1}_{-}) \quad \theta_1 \in \mathbb{R}$$

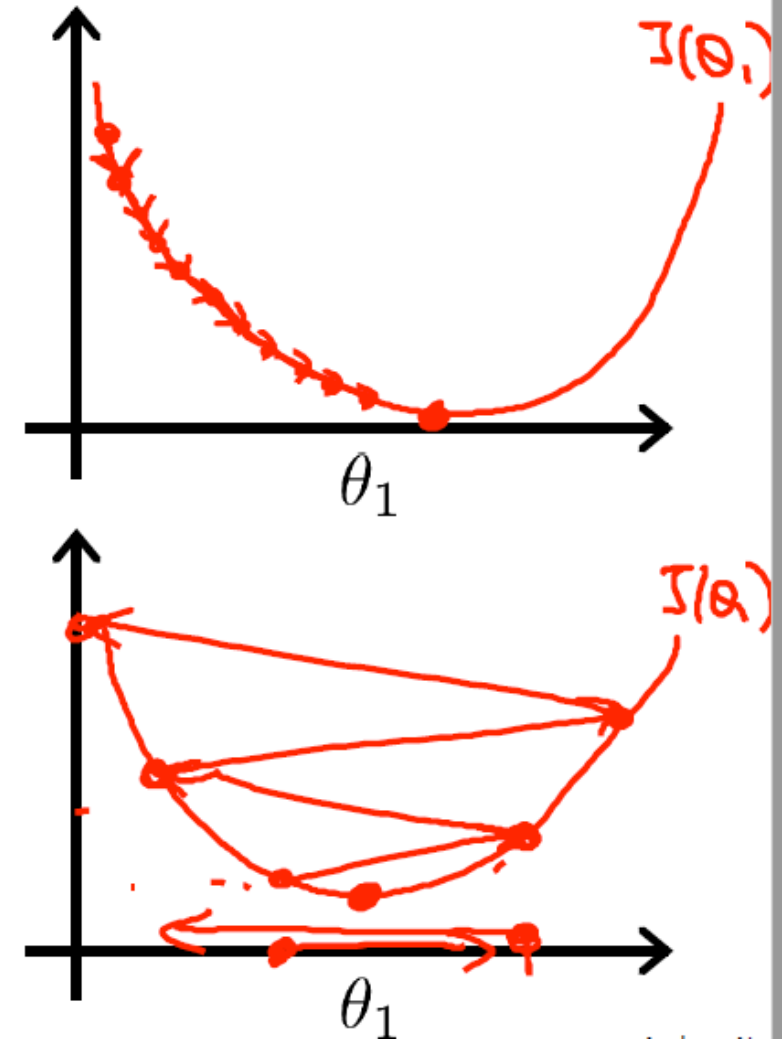
# Actualización en función de la pendiente



$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

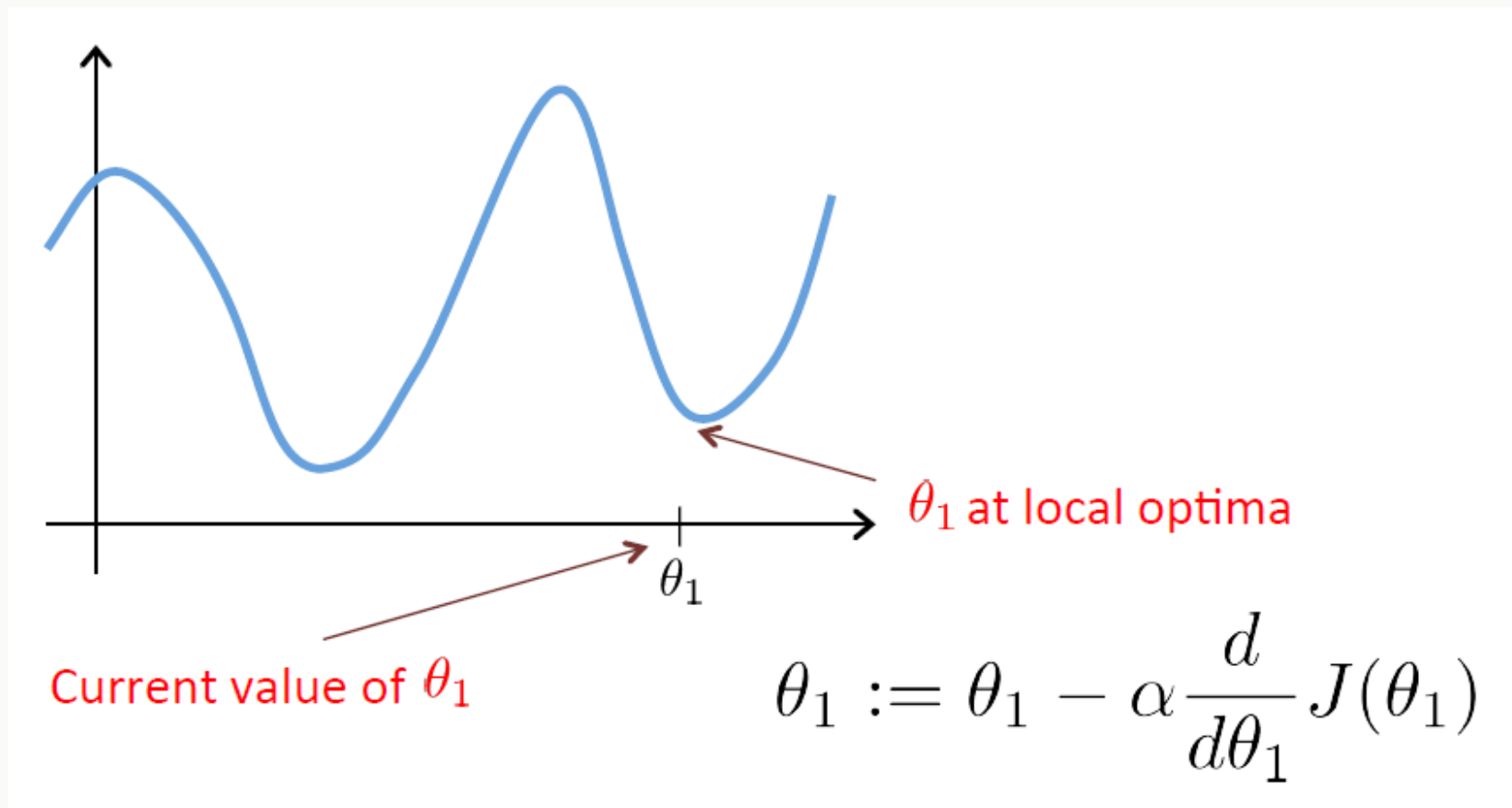
If  $\alpha$  is too small, gradient descent can be slow.

If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.





# Mínimos locales



# Convergencia

- ▶ El descenso de gradiente puede converger a un mínimo local, aún con una tasa de aprendizaje  $\alpha$  fija.

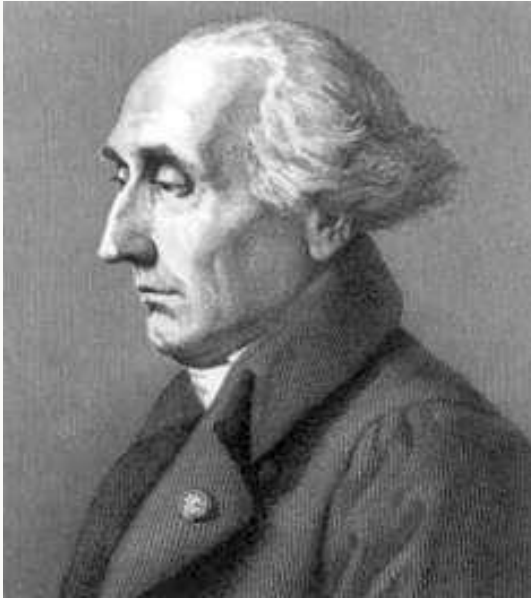
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_j} J(\theta_1)$$

- ▶ Conforme nos acercamos a un mínimo local, el descenso de gradiente toma automáticamente pasos más pequeños.
- ▶ De esta forma, no es necesario decrementar  $\alpha$  en el tiempo.

# OPTIMIZACIÓN POR MULTIPLICADORES DE LAGRANGE

Nociones básicas

# Joseph Louis de Lagrange



Bautizado como **Giuseppe Lodovico Lagrangia**, fue un Físico, Matemático y Astrónomo italiano naturalizado francés.

JOSEPH LOUIS LAGRANGE. Senador. Conde del Imperio. Gran Oficial de la Legión de Honor. Gran Cruz de la Imperial Orden de la Reunión. Miembro del Instituto y la Oficina de Longitudes. Nacido en Turín el 25 de enero de 1736. Muerto en París el 10 de abril de 1813.

Profesor de Joseph Fourier quien, aparentemente, no lo tenía en muy buen concepto.

Fuente: Wikipedia

# 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:

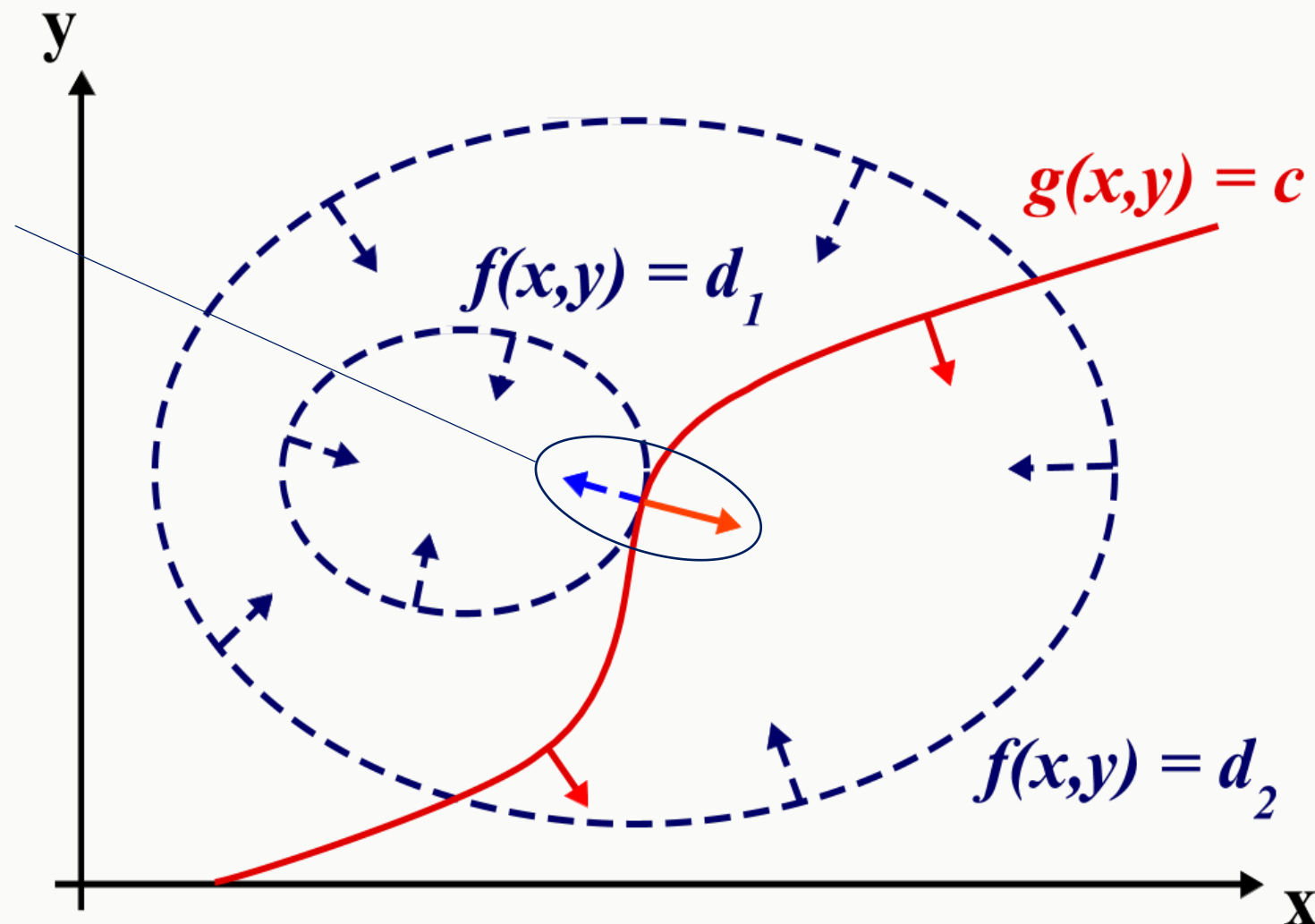
$$g(x_1, x_2, \dots, x_N) = c \quad c \in \mathbb{R}$$

- 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.

# CONDICIÓN DE TANGENCIA

Puntos donde los vectores de gradiente de ambas curvas son paralelos.

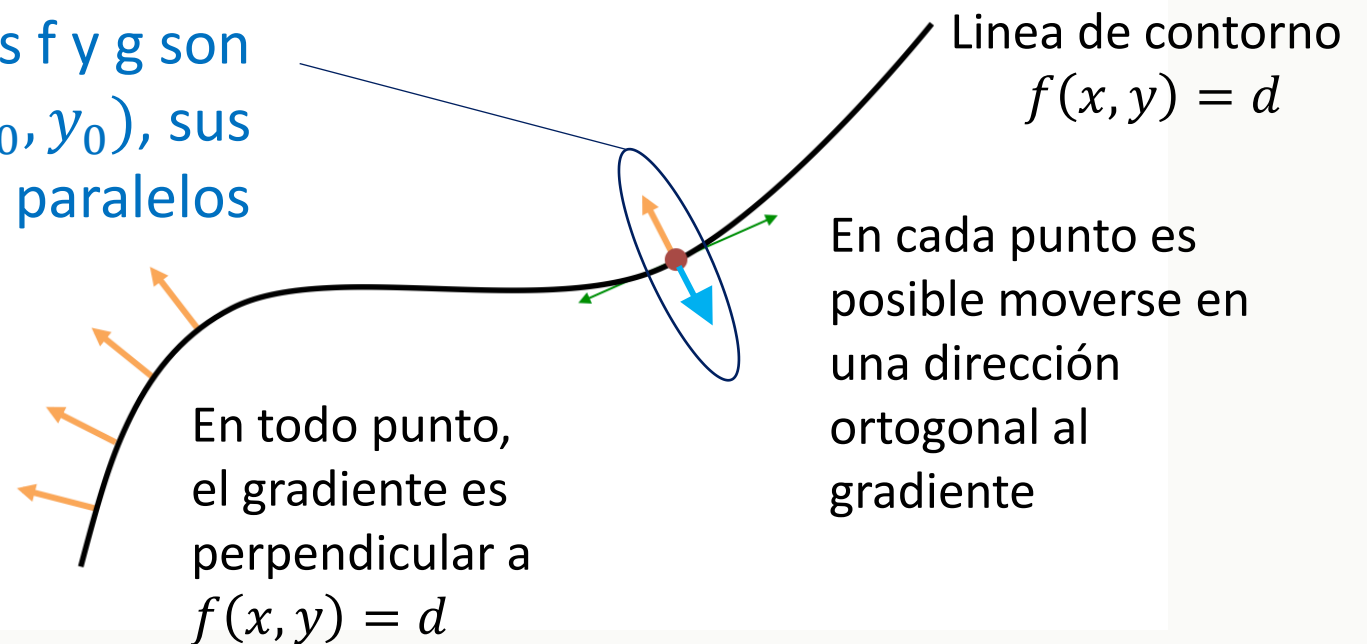
**¡Veamos el gradiente!**



# CONDICIÓN DE TANGENCIA

- Como hemos mencionado el gradiente  $\nabla f$  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  $f(x_1, x_2, \dots, x_N)$ :

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



# MULTIPLICADOR DE LAGRANGE

► 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  $\lambda$  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

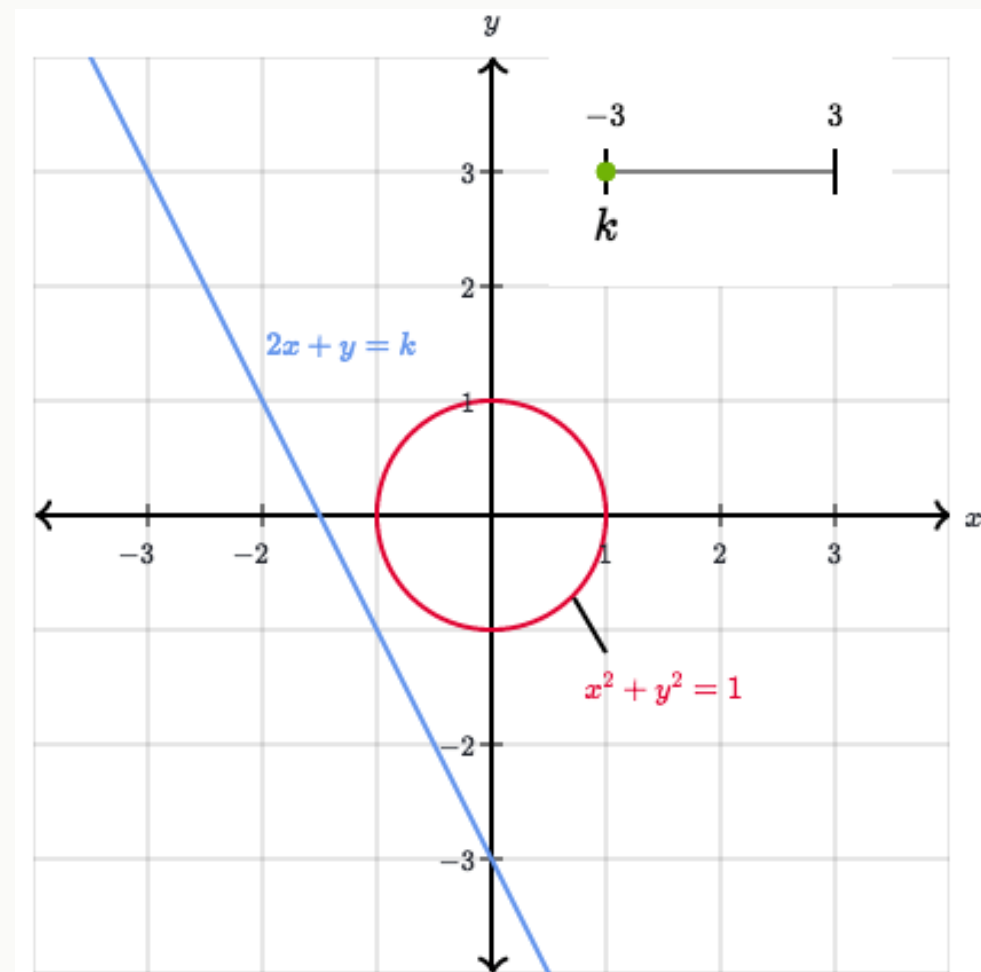
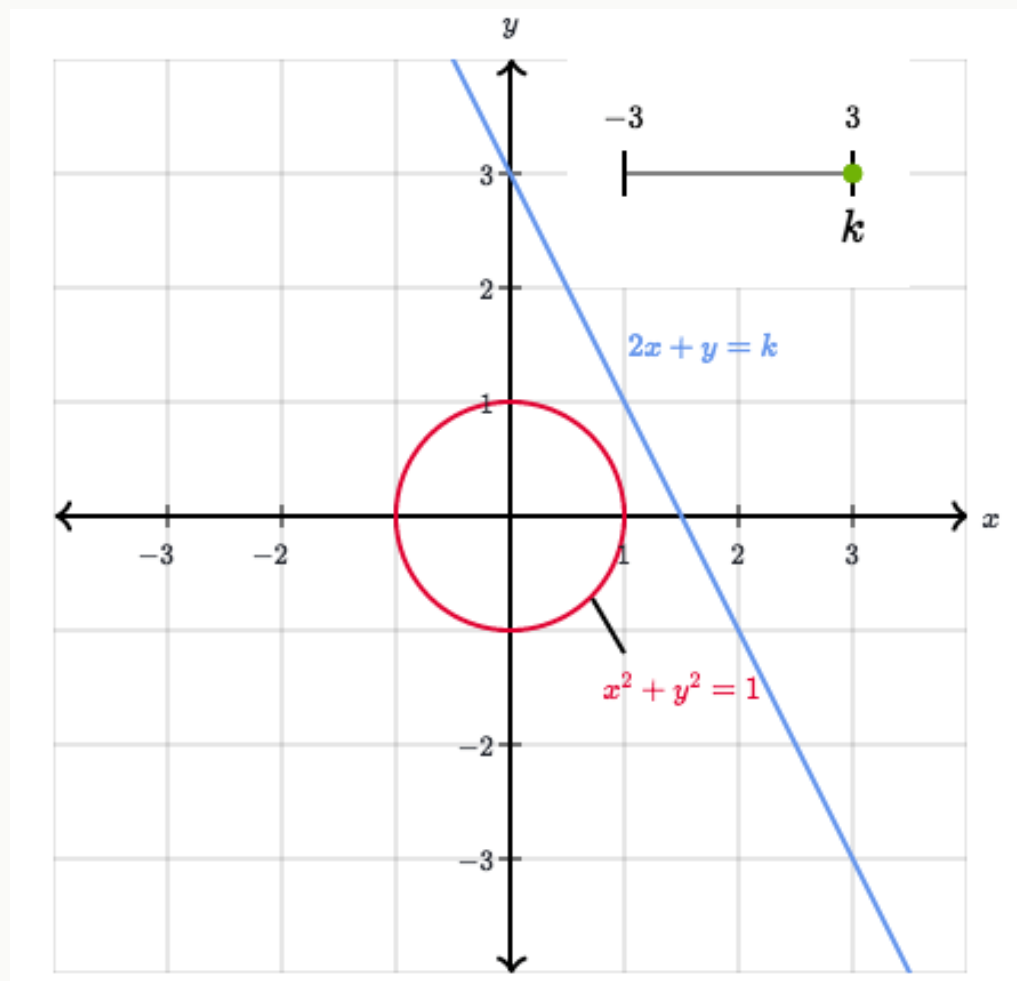
- Encuentre el valor máximo de las líneas de contorno  $f$ :

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

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

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

# EJEMPLO



## 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:*

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

## EJEMPLO (LAGRANGIANO)

- *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 (EL PROBLEMA)

► 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:*

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

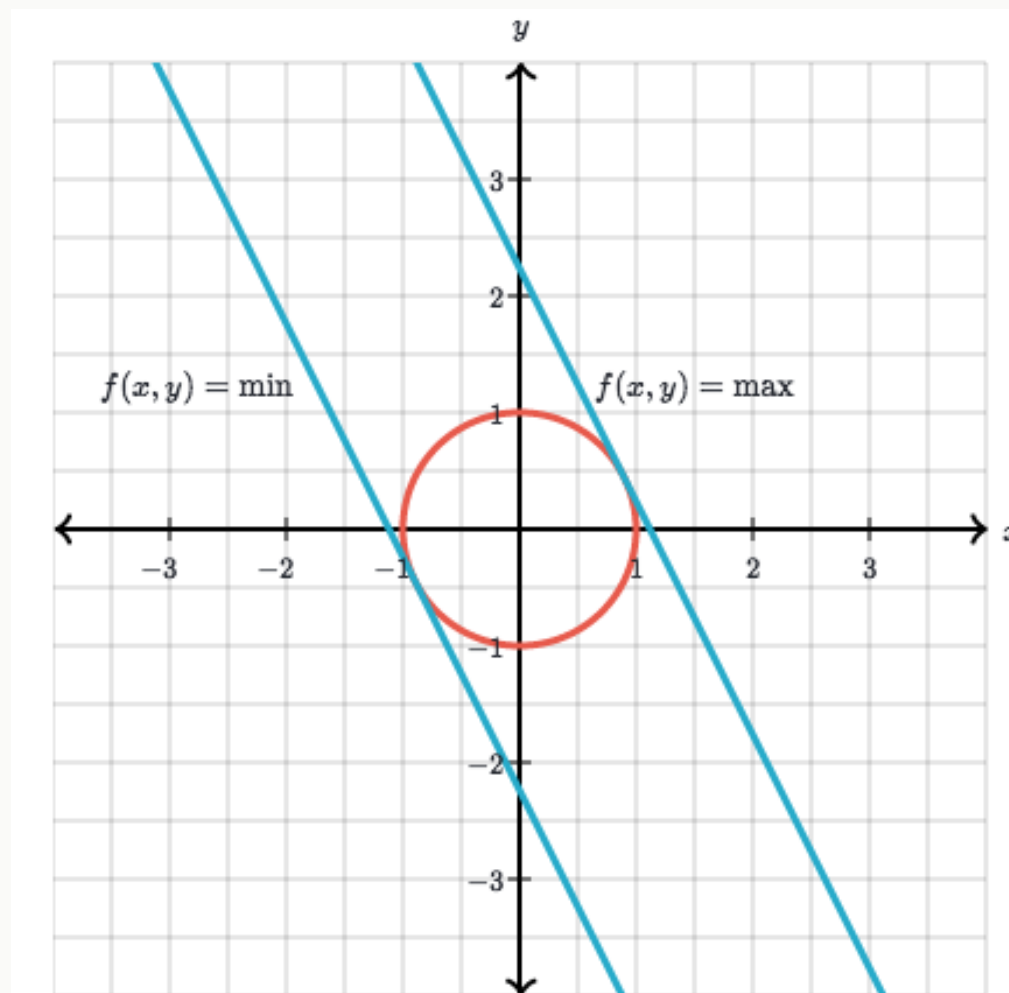
## EJEMPLO (LA SOLUCIÓN)

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

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

o bien

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





# OPTIMIZACIÓN CON RESTRICCIONES DE IGUALDAD

- 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: Resuelve las ecuaciones.

# OPTIMIZACIÓN CON RESTRICCIONES DE DESIGUALDAD

- ▶ Se consideran ahora restricciones del tipo  $g(x) \leq 0, g(x) \geq 0$ .
- ▶ Básicamente la estrategia es la misma que la descrita anteriormente, minimizando el lagrangiano con condiciones de igualdad y luego aplicando una de las siguientes reglas:

$$g(x) \geq 0 \Rightarrow \lambda \geq 0$$

$$g(x) \leq 0 \Rightarrow \lambda \leq 0$$

$$g(x) = 0 \Rightarrow \lambda \text{ no está restringida}$$