

# HANDS ON 04 – EXTREMOS

COMPLEMENTOS DE MATEMÁTICAS

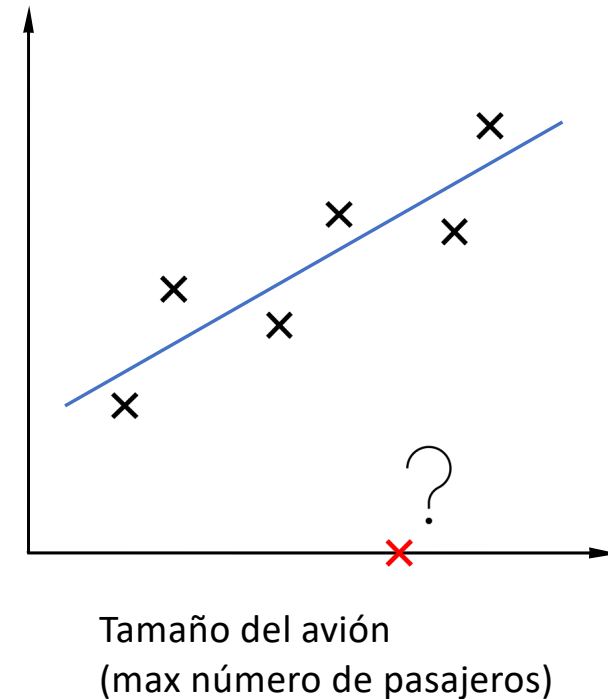
Gonzalo Rubio  
[g.rubio@upm.es](mailto:g.rubio@upm.es)

# Predicción basada en datos



Precio del avión (€)

- ¿Cómo obtengo esa recta?
  - Camino 1: Álgebra lineal
  - Camino 2: Cálculo (de varias variables)
  - Camino 3: Machine learning



# Mínimos cuadrados. Historia.

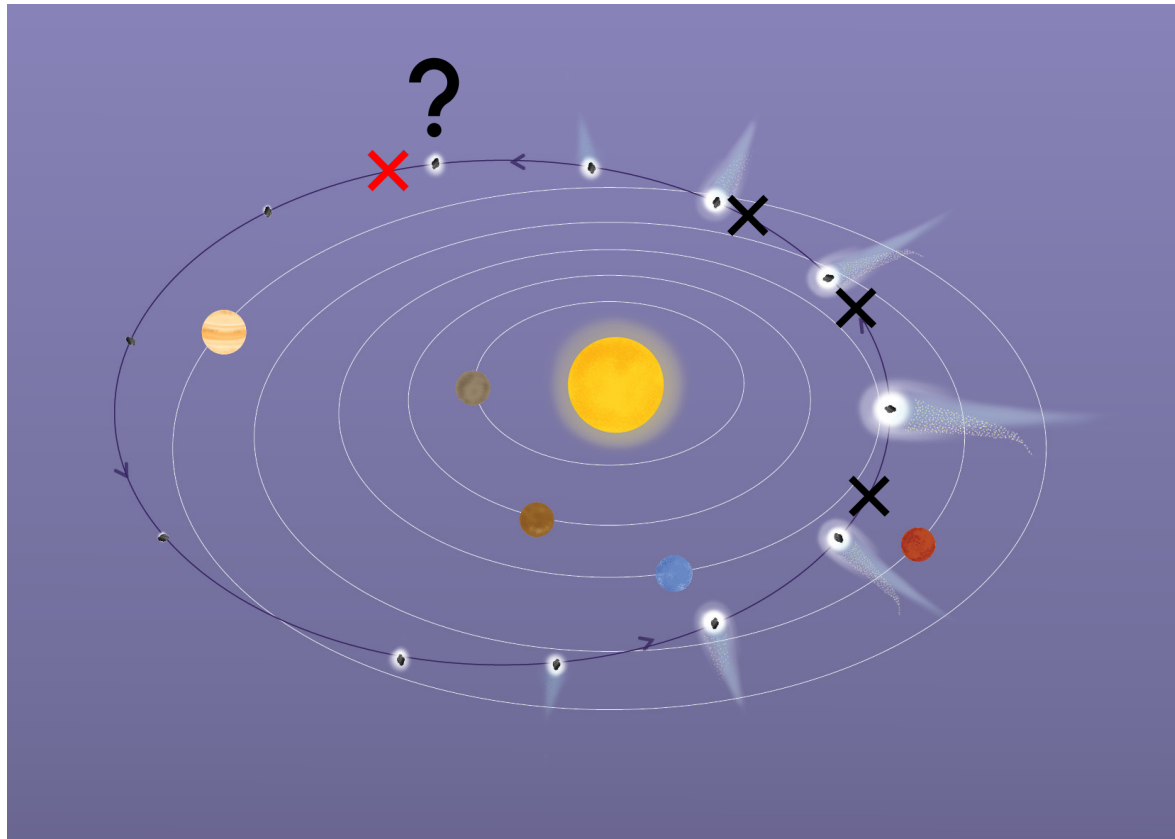
- Principios del siglo XIX
  - Adrien-Marie Legendre introduce el método de mínimos cuadrados
  - Legendre trabajó en la predicción de la trayectoria de proyectiles, midiendo la forma de la tierra y en astronomía.
  - Problema recurrente ajuste de datos contaminados por error.



Retrato usado durante 200 años para representar a Adrien-Marie Legendre

En 2005 se descubrió que que este es en realidad un político llamado Louis Legendre

# Mínimos cuadrados. Historia.



GONZALO RUBIO CALZADO



Solución: tener los errores en cuenta.  
Usar la ecuación orbital que mejor ajustara los datos, teniendo en cuenta que contenían un error.  
Método fue publicado en 1805 y fue todo un éxito. En solo 10 años se había convertido en una herramienta estándar en Francia, Italia y Prusia.

# ¿Cómo obtengo esa recta? Álgebra lineal

$$(1,1), (2,2), (3,2)$$

$$1 = C + D$$

$$2 = C + 2D$$

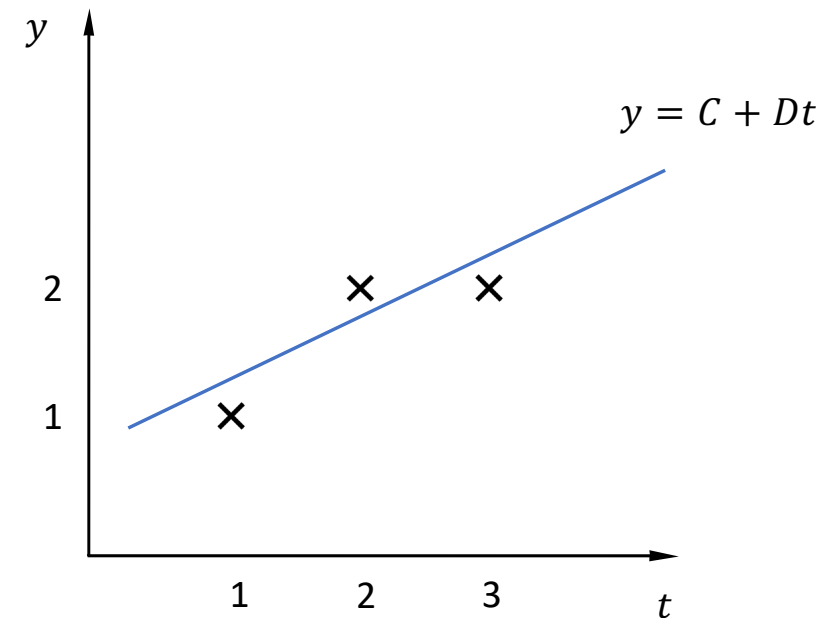
$$2 = C + 3D$$

No parece que este sistema tenga solución

$$Ax = b$$

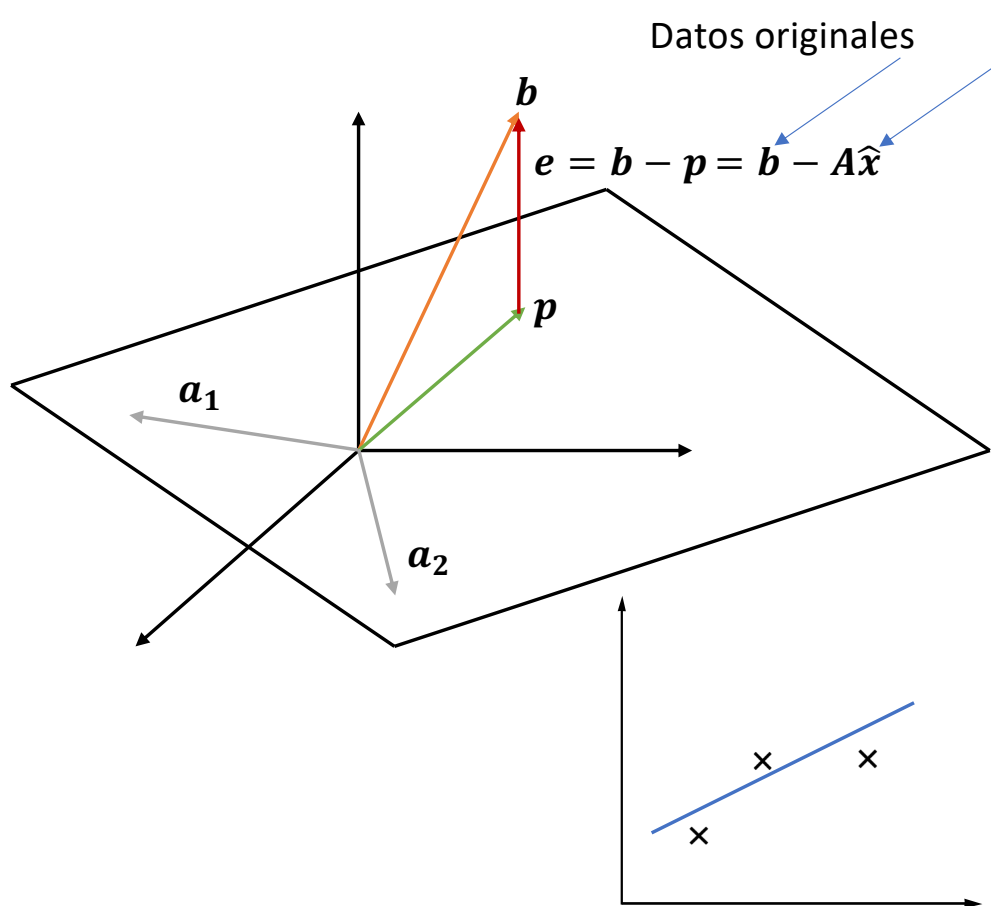
$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

Efectivamente,  $b$  no pertenece a  $C(A)$ ,  
luego el sistema no tiene solución



Hagamos que  $b$  pertenezca a  $C(A)$

# Recordatorio: proyección ortogonal



$$Ax = b$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

$$A = \begin{bmatrix} \vdots & \vdots \\ a_1 & a_2 \\ \vdots & \vdots \end{bmatrix}$$

$C(A)$

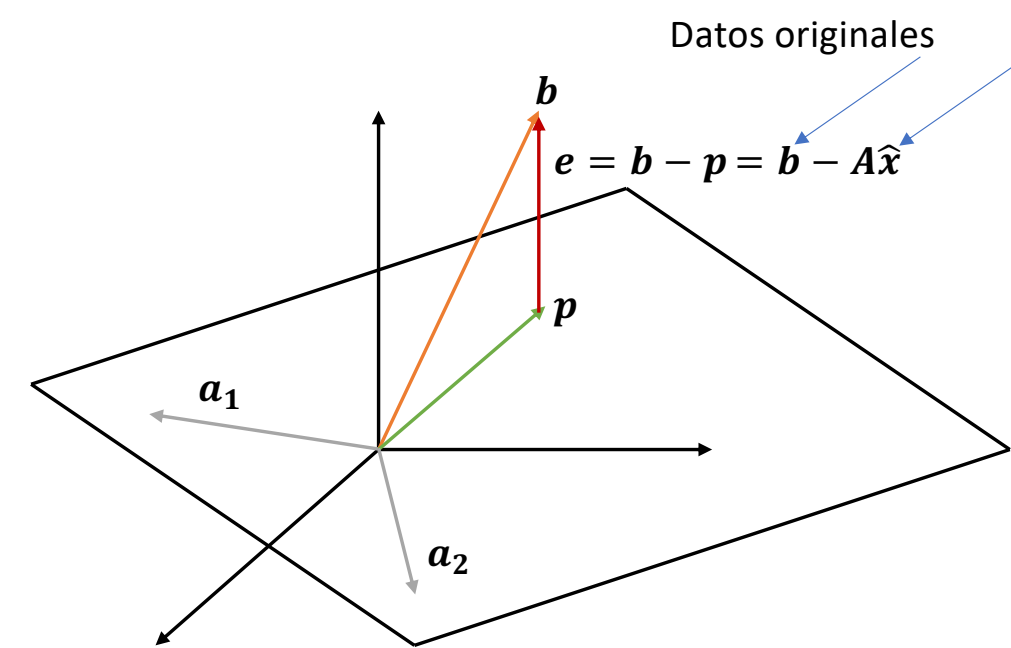
Plano formado por los vectores columna de  $A$

Como  $p$  pertenece a  $C(A)$ , puedo usarlo de RHS del problema original  $Ax = b \rightarrow A\hat{x} = p$  y obtener solución

No puedo pasar por los puntos  $\begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$ .

Intento pasar por unos alineados y que estén lo **más cerca posible** de los originales.

# Recordatorio: proyección ortogonal



$$\mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \quad \|\mathbf{e}\|^2 = e_1^2 + e_2^2 + e_3^2 \quad \text{Mínimo}$$

$$A\mathbf{x} = \mathbf{b}$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

$$A = \begin{bmatrix} \vdots & \vdots \\ \mathbf{a}_1 & \mathbf{a}_2 \\ \vdots & \vdots \end{bmatrix}$$

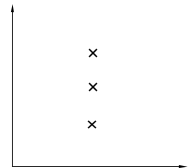
$$\mathbf{p} = \hat{x}_1 \mathbf{a}_1 + \hat{x}_2 \mathbf{a}_2 = A\hat{\mathbf{x}}$$

Si quiero conocer  $\mathbf{p}$ , mi problema es encontrar  $\hat{\mathbf{x}}$

CLAVE:  $\mathbf{e} = \mathbf{b} - \mathbf{p} = \mathbf{b} - A\hat{\mathbf{x}}$  es perpendicular al plano (para que  $\mathbf{b}$  y  $\mathbf{p}$  estén lo más cerca posible)

Luego debe ser perpendicular a  $\mathbf{a}_1$  y también a  $\mathbf{a}_2$

# Recordatorio: proyección ortogonal



$$\begin{aligned} a_1^T(b - A\hat{x}) &= 0 \\ a_2^T(b - A\hat{x}) &= 0 \end{aligned} \longrightarrow A^T(b - A\hat{x}) = 0$$

Reordenando términos:

$$A^T A \hat{x} = A^T b$$

Luego  $\hat{x}$

$$\hat{x} = (A^T A)^{-1} A^T b$$

Y  $p$

$$p = A(A^T A)^{-1} A^T b$$

Y la ecuación a resolver es:

$$A\hat{x} = A(A^T A)^{-1} A^T b$$

$(A^T A)$  es invertible si  
 $A$  tiene columnas  
independientes



$$Ax = b$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

$$A = \begin{bmatrix} \vdots & \vdots \\ a_1 & a_2 \\ \vdots & \vdots \end{bmatrix}$$

$$p = \hat{x}_1 a_1 + \hat{x}_2 a_2 = A\hat{x}$$

Si quiero conocer  $p$ , mi problema es  
encontrar  $\hat{x}$

CLAVE:  $e = b - p = b - A\hat{x}$  es perpendicular al  
plano (**para que  $b$  y  $p$  estén lo más cerca posible**)

Luego debe ser perpendicular a  $a_1$  y  
también a  $a_2$



# PEQUEÑO INCISO

$(A^T A)$  es invertible si  $A$  tiene columnas independientes

Demostración:

$A^T A$  es invertible implica que  $A^T A x = 0$  si y solo si  $x = 0$

IDEA FELIZ

$$x^T A^T A x = 0$$

$$(Ax)^T Ax = 0$$

$$||Ax||^2 = 0$$

Las columnas de  $A$  son independientes luego esto solo se cumple si  $x = 0$

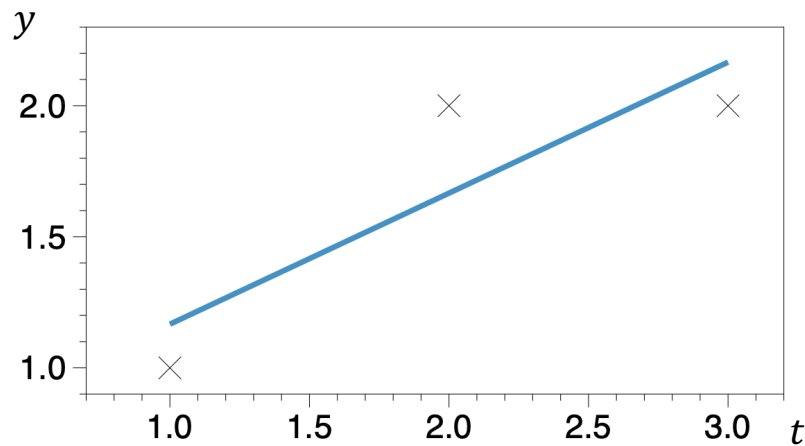
# Recta de regresión lineal

$$Ax = b$$

$$A^T A \hat{x} = A^T b$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} \hat{C} \\ \hat{D} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$



$$\begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix} \begin{bmatrix} \hat{C} \\ \hat{D} \end{bmatrix} = \begin{bmatrix} 5 \\ 11 \end{bmatrix}$$

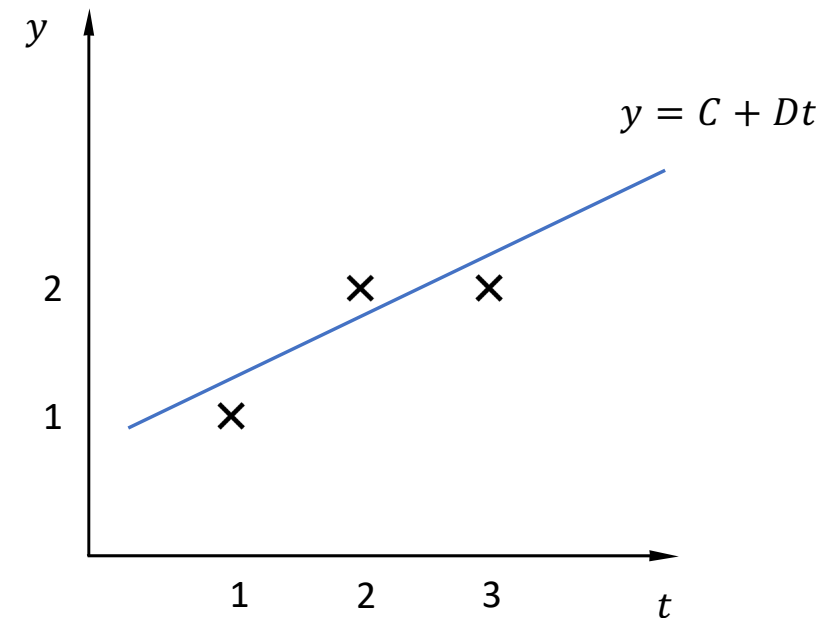
$$y = 2/3 + 1/2t$$

Ecuaciones normales

# Recta de regresión lineal. Un poco de Código.

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \quad \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ 1 & t_3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

```
8  import numpy as np
9
10
11  A = np.matrix([[1, 1], [1, 2], [1, 3]])
12  b = np.matrix([1, 2, 2])
13  b = b.transpose()
14
15  AT = A.transpose()
16
17  Ahat = np.matmul(AT, A)
18  bhat = np.matmul(AT, b)
19
20  x = np.linalg.solve(Ahat, bhat)     $A^T A \hat{x} = A^T b$ 
```

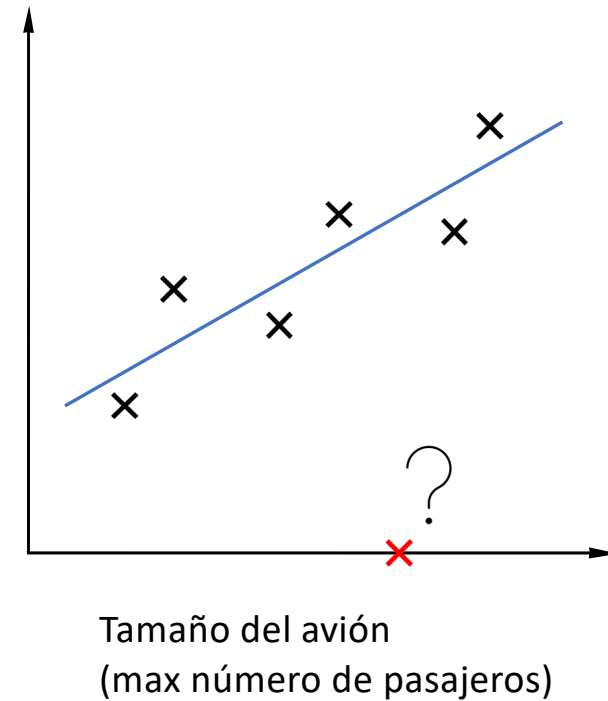


# Predicción basada en datos



Precio del avión (€)

- ¿Cómo obtengo esa recta?
  - Camino 1: Álgebra lineal
  - Camino 2: Cálculo (de varias variables)
  - Camino 3: Machine learning



# ¿Cómo obtengo esa recta? Cálculo

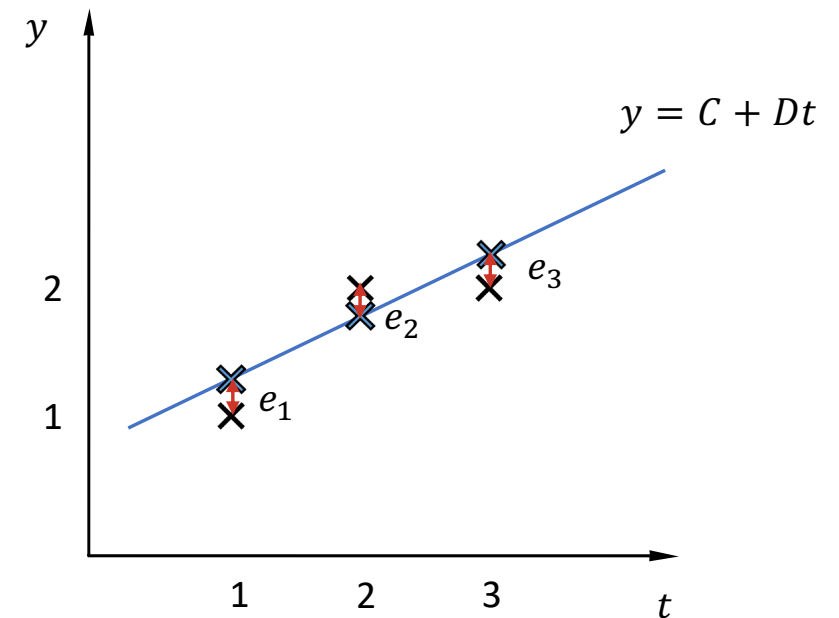
$$\begin{aligned} e_1 &= C + D - 1 \\ e_2 &= C + 2D - 2 \\ e_3 &= C + 3D - 2 \end{aligned} \quad \mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix}$$

Para que la recta sea “buena”  $\|\mathbf{e}\|$  debe ser pequeño

Es útil trabajar con  $\|\mathbf{e}\|^2$ . Buscamos su mínimo

$$\|\mathbf{e}\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

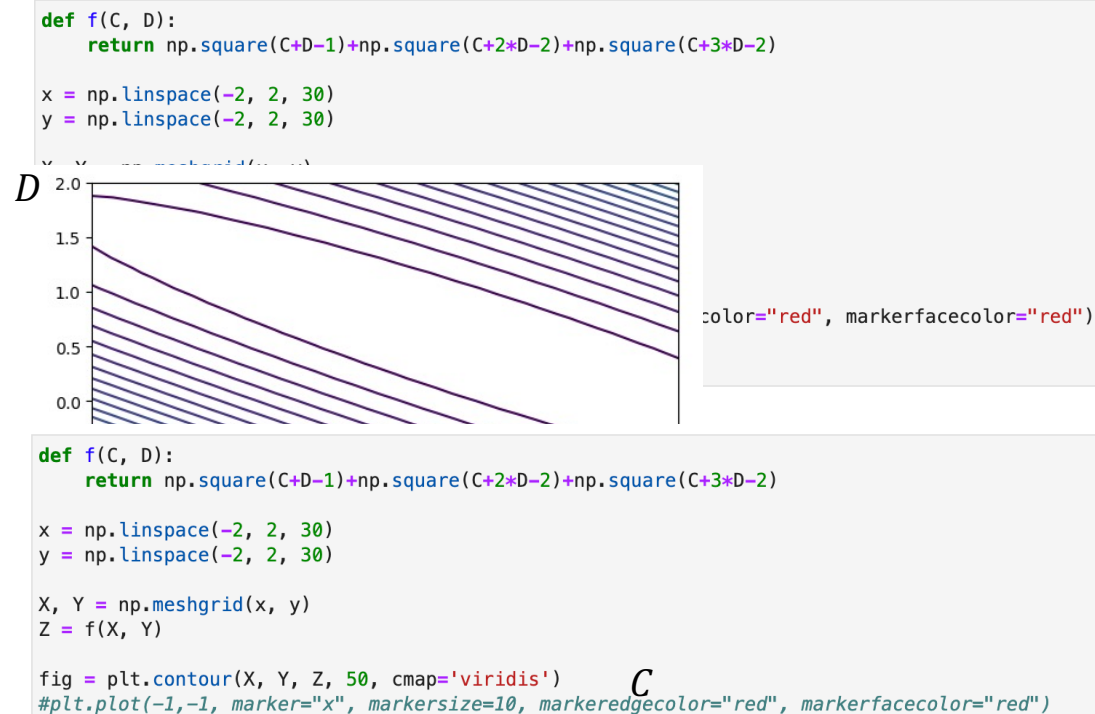
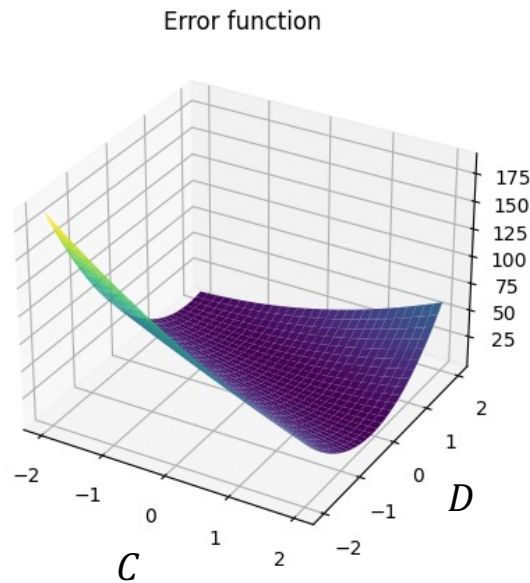
Esto es una función real de varias variables  $\|\mathbf{e}\|^2 = f(C, D)$



# El error es una función de varias variables

$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

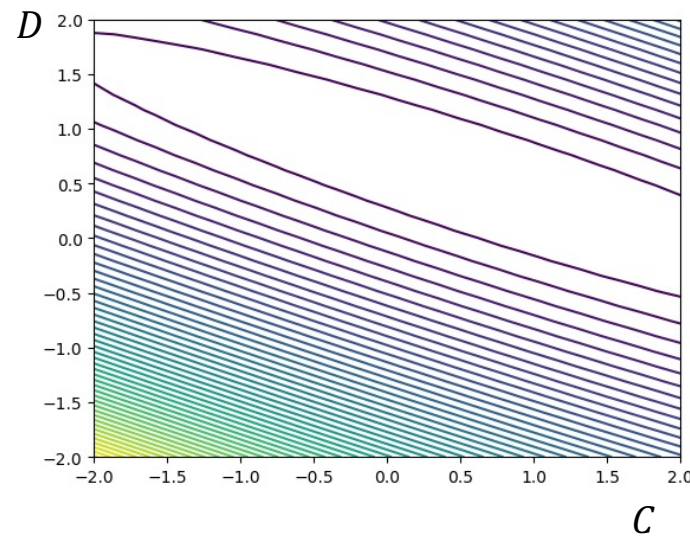
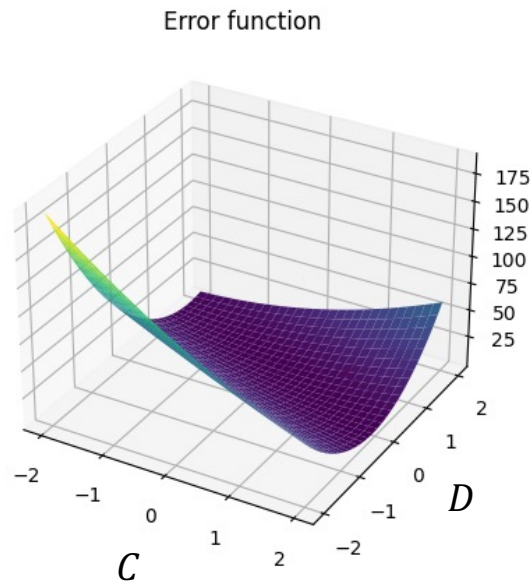
Esto es una función real de varias variables  $\|e\|^2 = f(C, D)$



# Relación entre función de error y recta

$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

Esto es una función real de varias variables  $\|e\|^2 = f(C, D)$



# Relación entre función de error y recta

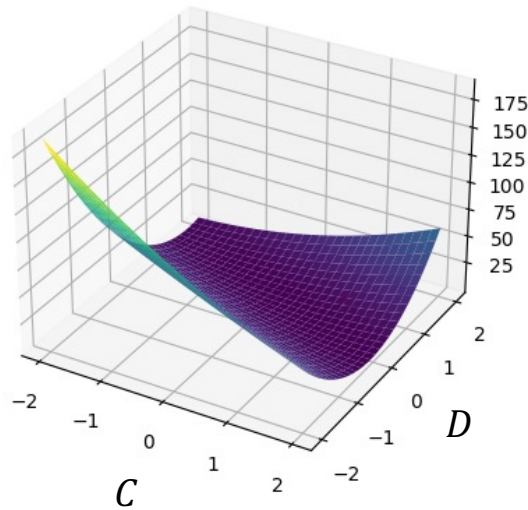
$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 -$$

Esto es una función real de varias variable

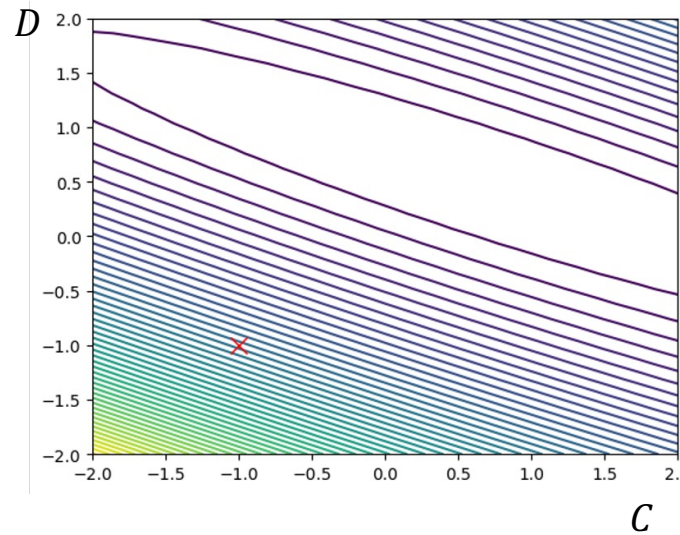
```
x = np.linspace(0, 3, 100)
y = -1-1*x
plt.plot(x, y)
plt.plot(1,1,marker="x", markersize=10, markeredgecolor="red", markerfacecolor="red")
plt.plot(2,2,marker="x", markersize=10, markeredgecolor="red", markerfacecolor="red")
plt.plot(3,2,marker="x", markersize=10, markeredgecolor="red", markerfacecolor="red")
plt.show()
```

$$y = C + Dt$$

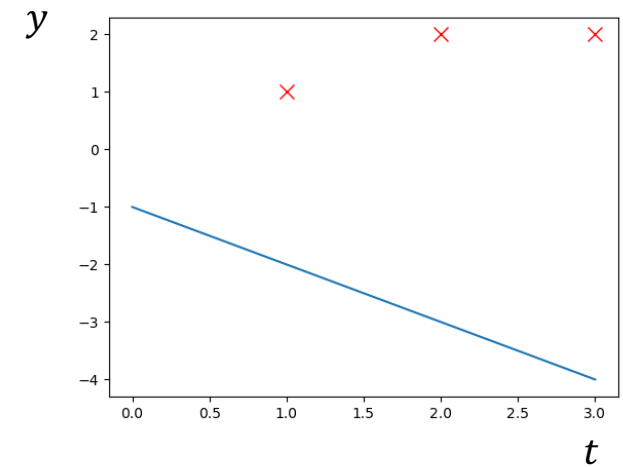
Error function



$$f(-1, -1) = 70$$



$$y = -1 - t$$





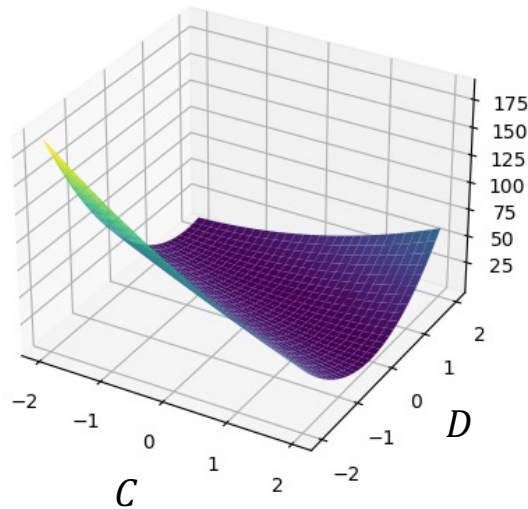
# Relación entre función de error y recta

$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

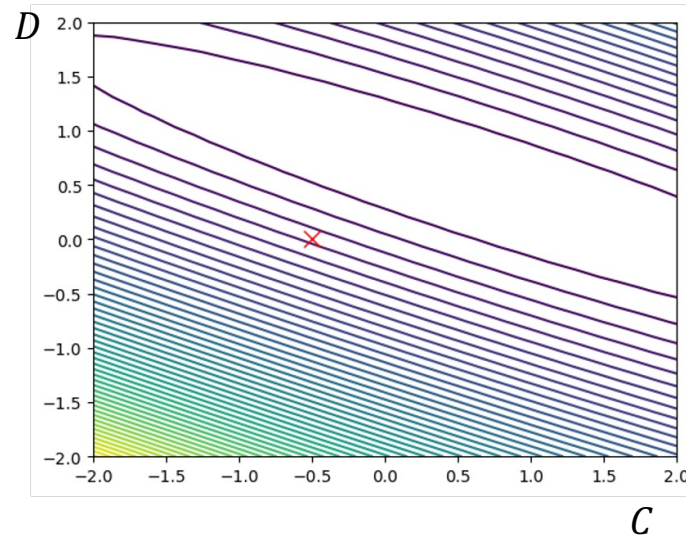
Esto es una función real de varias variables  $\|e\|^2 = f(C, D)$

$$y = C + Dt$$

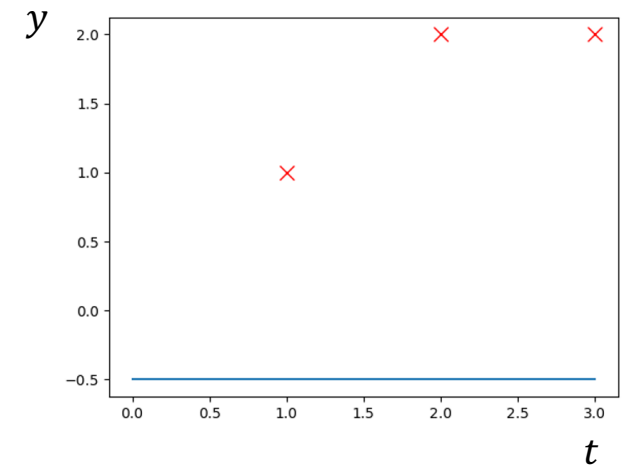
Error function



$$f(-0.5, 0) = 14.75$$



$$y = -0.5$$



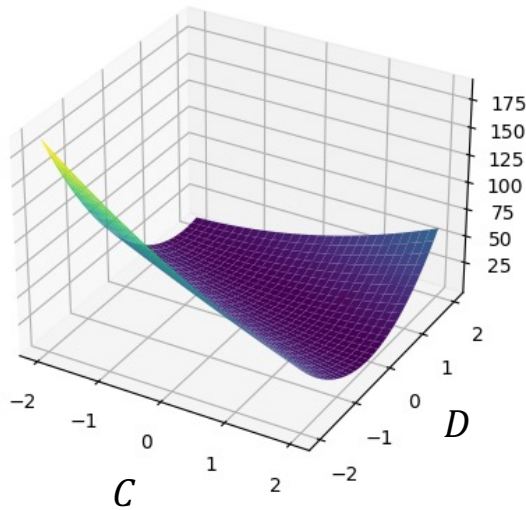
# Relación entre función de error y recta

$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

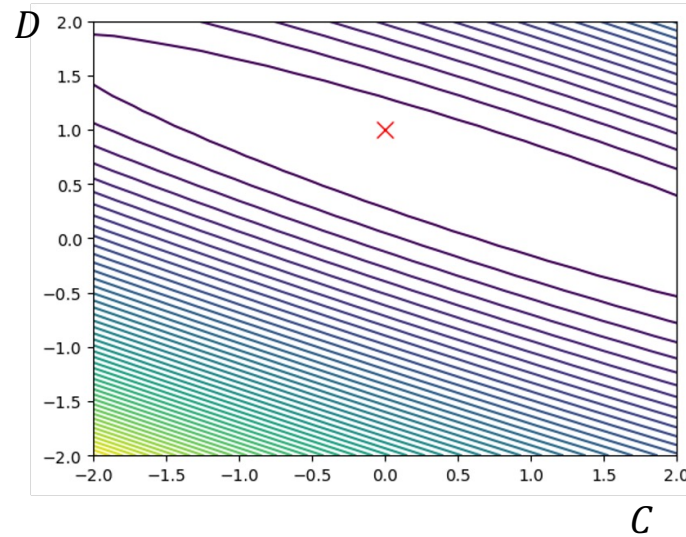
Esto es una función real de varias variables  $\|e\|^2 = f(C, D)$

$$y = C + Dt$$

Error function

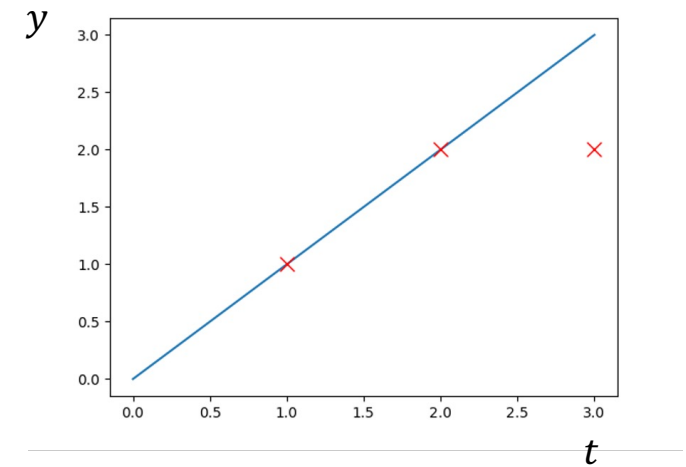


$$f(0,1) = 1$$



GONZALO RUBIO CALZADO

$$y = t$$



18

# Recordatorio: cálculo de extremos

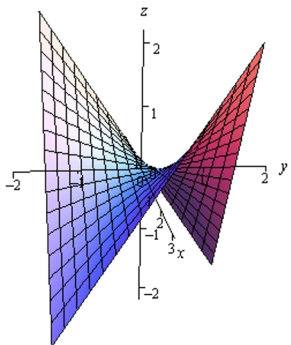
Buscamos candidatos a extremos en funciones de varias variables.

**1. Puntos críticos (o estacionarios):** Sea  $f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$  y  $x_0 \in \overset{\circ}{D}$  (interior de  $D$ )

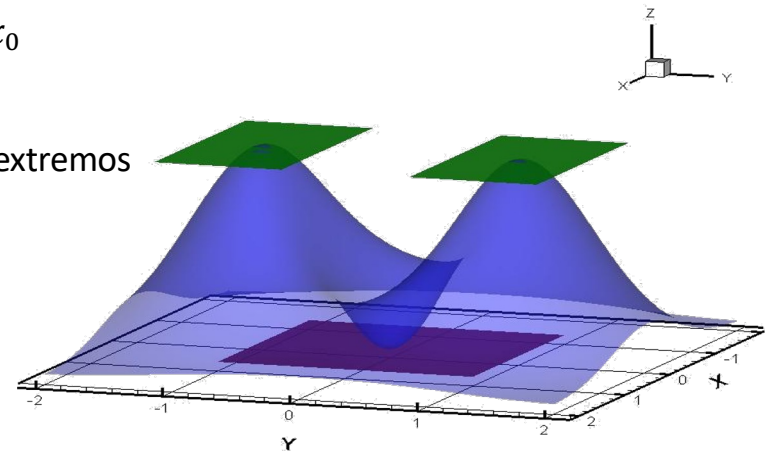
$x_0$  es punto crítico de  $f$  si  $f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$  es diferenciable en  $x_0$  y  $\nabla f(x_0) = \mathbf{0} \Leftrightarrow df(x_0) = 0$

En  $\mathbb{R}^2$  esta condición es equivalente a que el plano tangente sea horizontal en  $x_0$

Sin embargo la existencia de puntos críticos no garantiza la existencia de extremos



Ejemplo: la función  $f(x, y) = xy$  tienen un punto crítico en  $(x, y) = (0, 0)$  pero no es un extremo local



# Recordatorio: cálculo de extremos

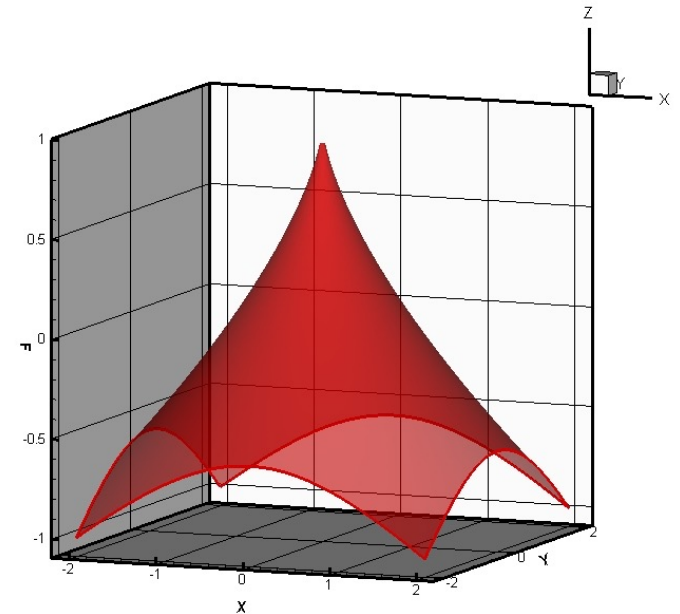
Buscamos candidatos a extremos en funciones de varias variables.

**2. Puntos singulares:** Sea  $f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$  y  $x_0 \in \overset{\circ}{D}$  (interior de  $D$ )

$x_0$  es un punto singular de  $f$  si no es diferenciable en  $x_0$

**3. Puntos frontera del dominio.**

Si  $D$  incluye al menos a parte de su frontera, puede ocurrir que existan extremos de la función en la frontera del dominio.



# Condición de mínimo

$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

Esto es una función real de varias variables  $f(C, D)$

Mínimos relativos cumplirán  $\nabla f(C, D) = \mathbf{0}$  (Función de error  $f(C, D)$  es siempre creciente)

$$\frac{\partial f}{\partial C} = 0 \longrightarrow \frac{\partial f}{\partial C} = 2(C + D - 1) + 2(C + 2D - 2) + 2(C + 3D - 2) = 0 \longrightarrow 6C + 12D - 10 = 0$$

$$\frac{\partial f}{\partial D} = 0 \longrightarrow \frac{\partial f}{\partial D} = 2(C + D - 1) + 4(C + 2D - 2) + 6(C + 3D - 2) = 0 \longrightarrow 12C + 28D - 22 = 0$$

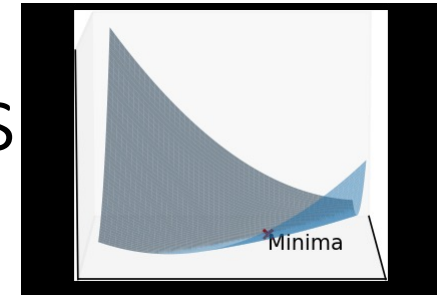
$$\begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 5 \\ 11 \end{bmatrix}$$

¡Igual que por el camino anterior!

# Recta de regresión. Mínimos cuadrados

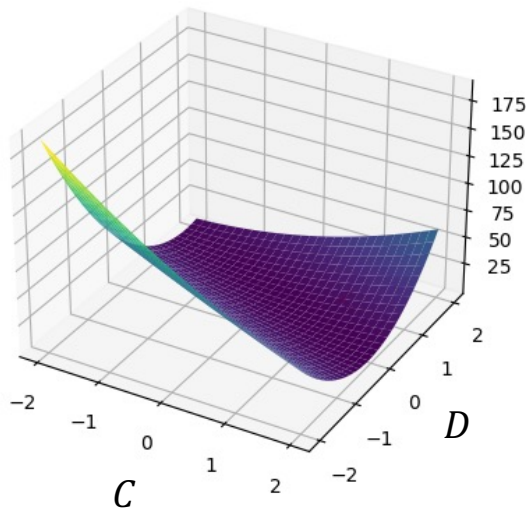
$$\|e\|^2 = (C + D - 1)^2 + (C + 2D - 2)^2 + (C + 3D - 2)^2$$

Esto es una función real de varias variables  $f(C, D)$

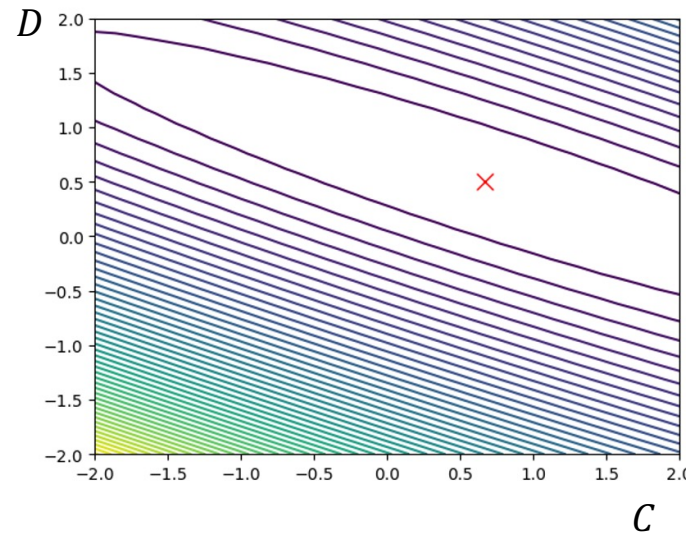


$$y = C + Dt$$

Error function

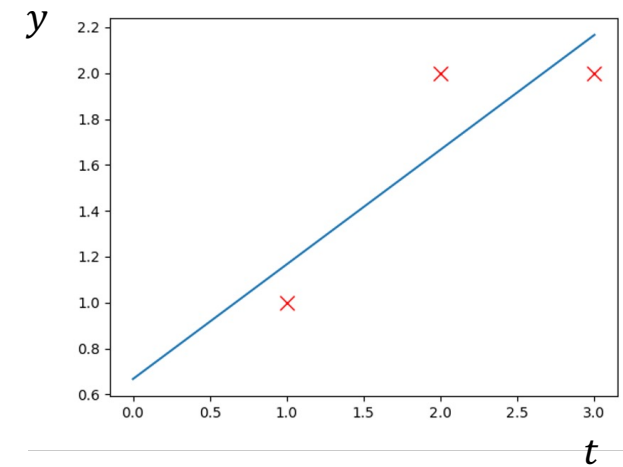


$$f(0,1) = 0.167$$



GONZALO RUBIO CALZADO

$$y = 2/3 + 1/2t$$

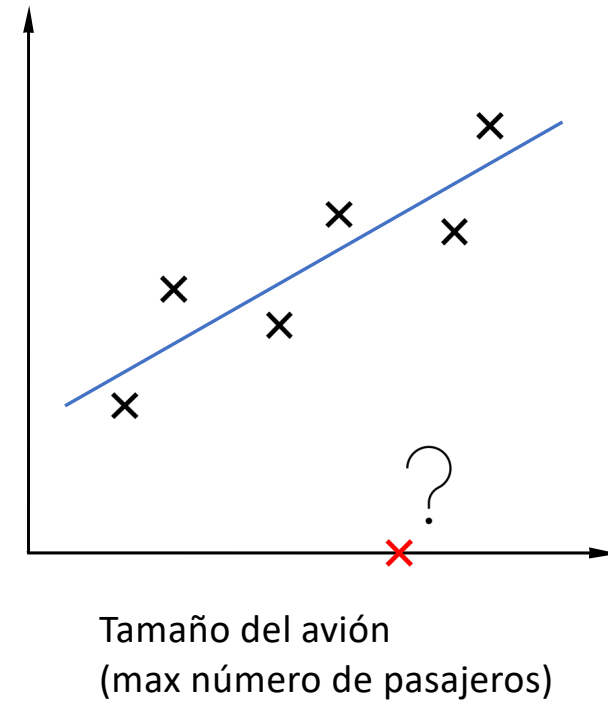


# Predicción basada en datos



Precio del avión (€)

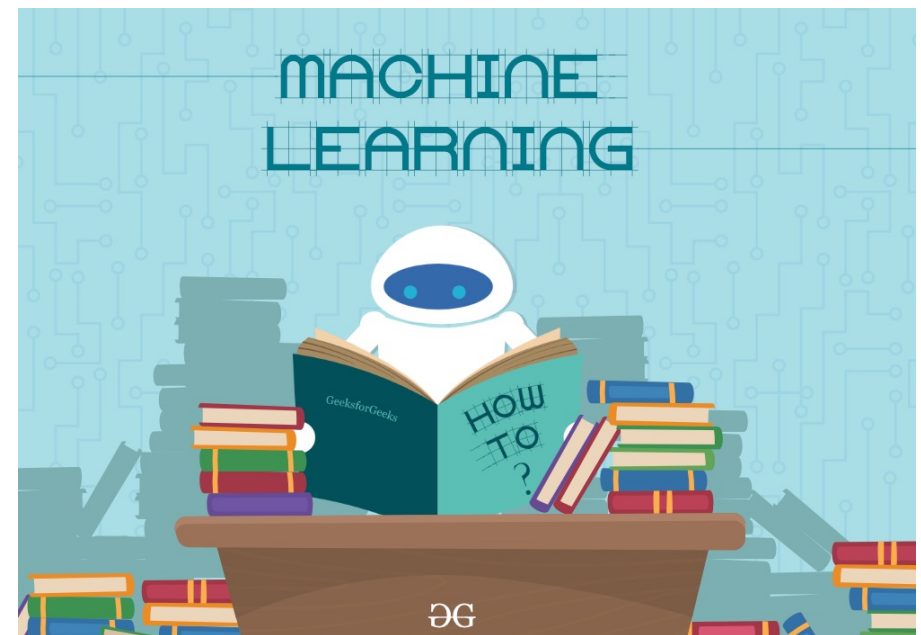
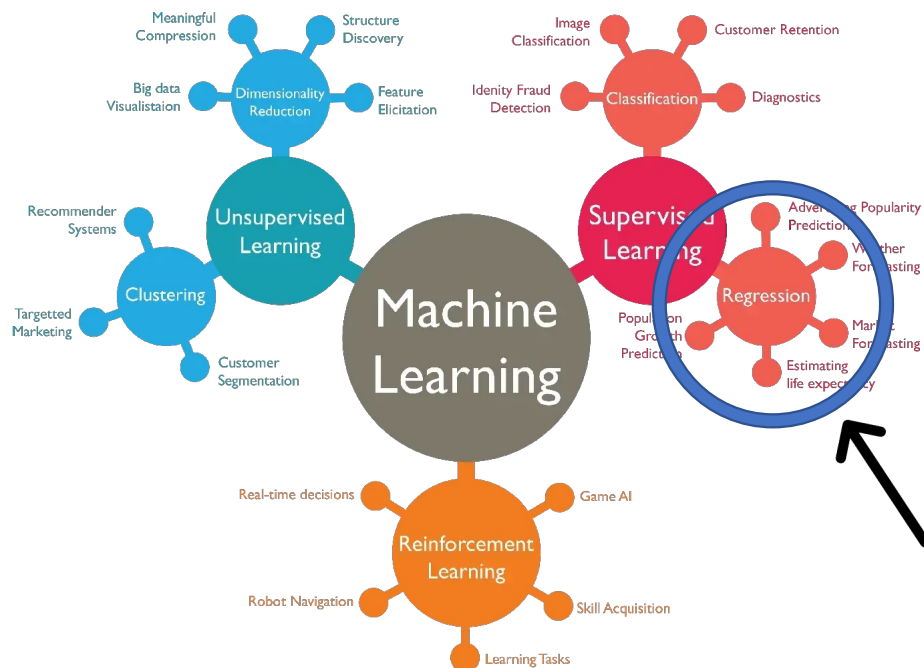
- ¿Cómo obtengo esa recta?
  - Camino 1: Álgebra lineal
  - Camino 2: Cálculo (de varias variables)
  - Camino 3: Machine learning





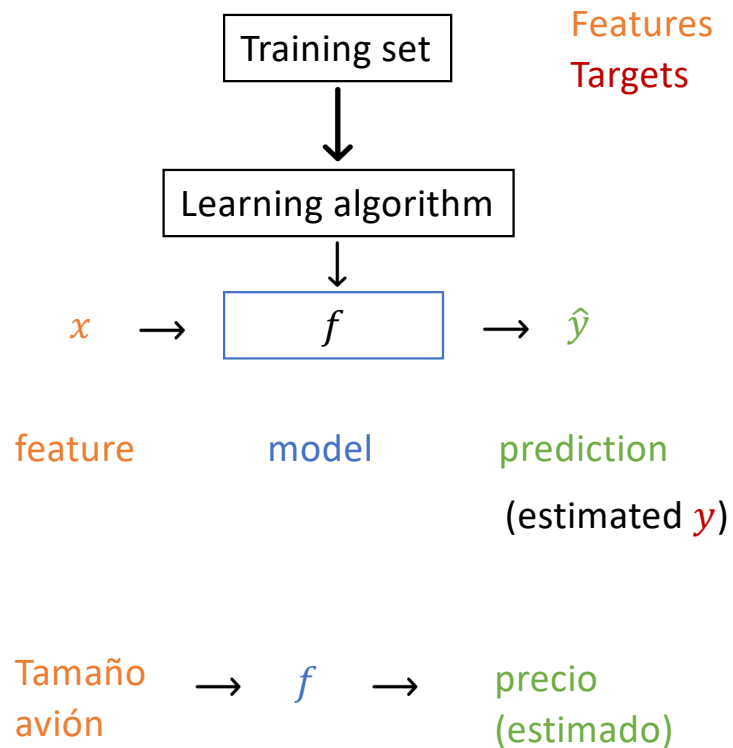
# ¿Cómo obtengo esa recta? Machine learning

Campo de estudio que permite que un ordenador “aprenda” sin necesidad de programarlo de forma explícita.

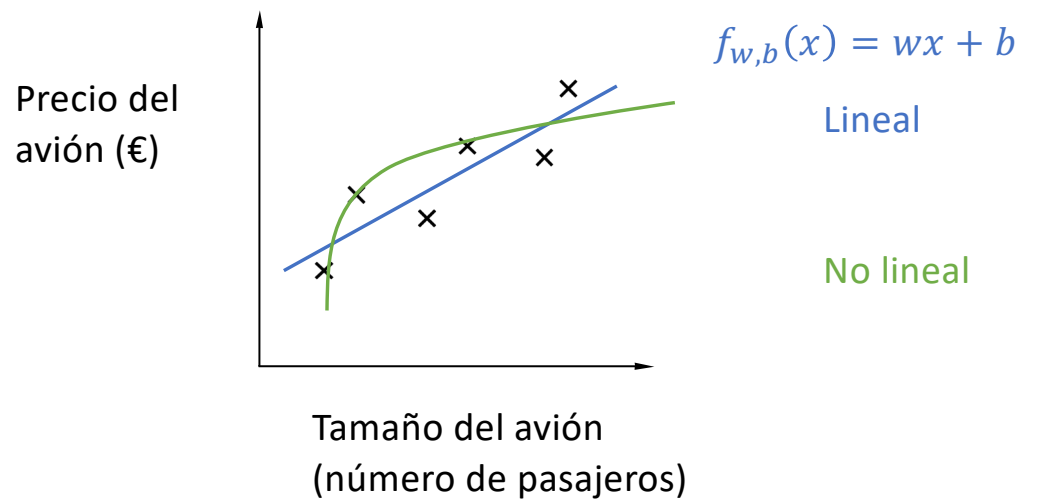




# Mínimos cuadrados - Machine learning

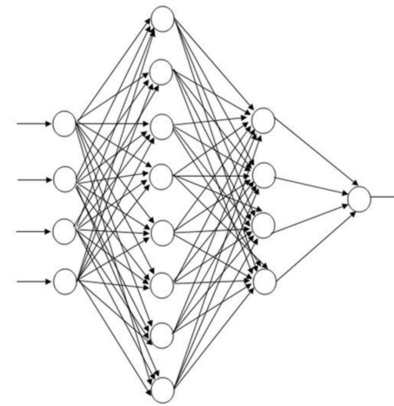
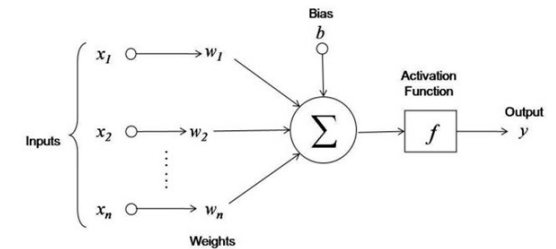
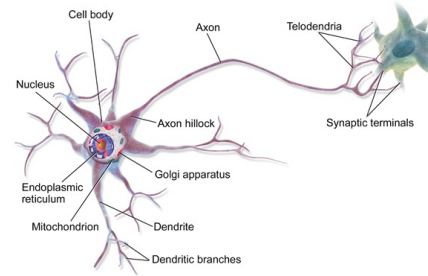
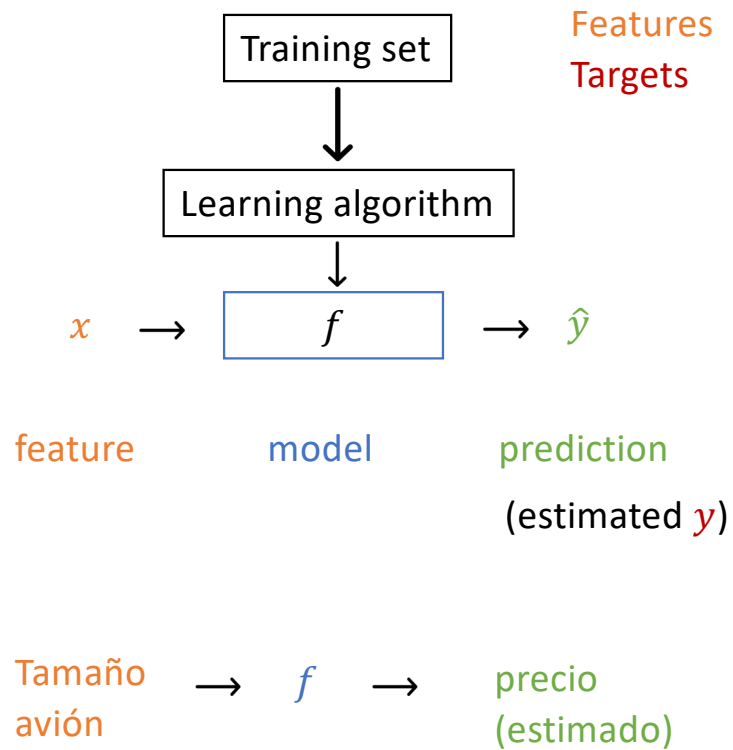


¿Cómo representamos  $f$ ?



Nuestro modelo es lineal. Sin embargo, se pueden utilizar modelos más complejos.

# Mínimos cuadrados - Machine learning



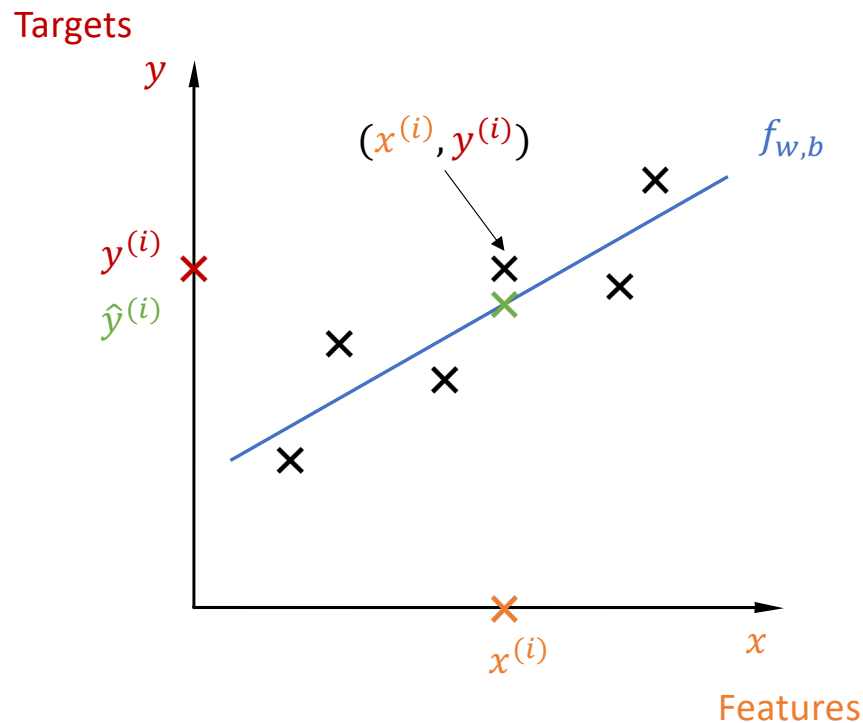
David E. Rumelhart



Geoffrey Hinton

Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533-536.

# Mínimos cuadrados - Función de coste



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b$$

Función de coste (Cost function): Error al cuadrado

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$m$  número de ejemplos en los datos de entrenamiento

Encontrar  $(w, b)$  :

$\hat{y}^{(i)}$  esté cerca de  $y^{(i)}$  para todos los  $(x^{(i)}, y^{(i)})$

# Un poco de código

Training set

Features  
Targets

Model

$$f_{w,b}(x) = wx + b$$

Lineal

Cost  
function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$m$  número de ejemplos en los  
datos de entrenamiento

```
12 # Training set
13 x = [1, 2, 3]
14 y = [1, 2, 2]
15
16 m = len(x)
17
18 # Model
19 def fwb(w,b,x):
20     return w*x+b
21
22 # Cost function
23 def J(w,b,x,y):
24     J = 0
25     for i in range(m):
26         J = J + (fwb(w,b,x[i])-y[i])**2
27     J = J / (2*m)
28     return J
```

# ¿Cómo encuentro el mínimo de $J(w, b)$ ?

La calidad de la aproximación dependerá del valor alcanzado por la función de coste  $J(w, b)$

$$\min_{w, b} J(w, b)$$

En el caso de derivación basada en cálculo, resolvimos este problema de forma analítica.

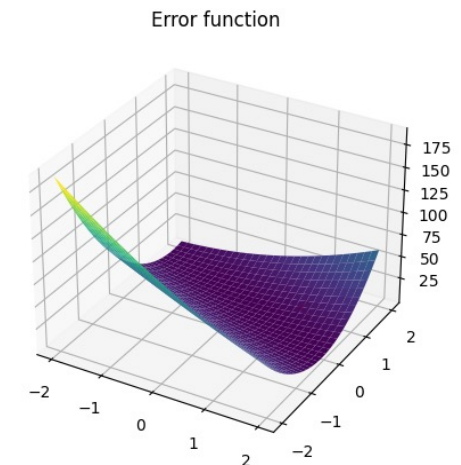
Plantear/Resolver un sistema lineal de ecuaciones puede no ser eficiente en un caso general/grande (dimensión).

En machine learning es típico utilizar el algoritmo Gradient descent (gradiente descendente).

Algoritmo Gradient descent:

Empezamos con un par de valores  $(w, b)$ . Por ejemplo:  $(w = 0, b = 0)$

Cambiamos los valores de  $(w, b)$  para reducir  $J(w, b)$

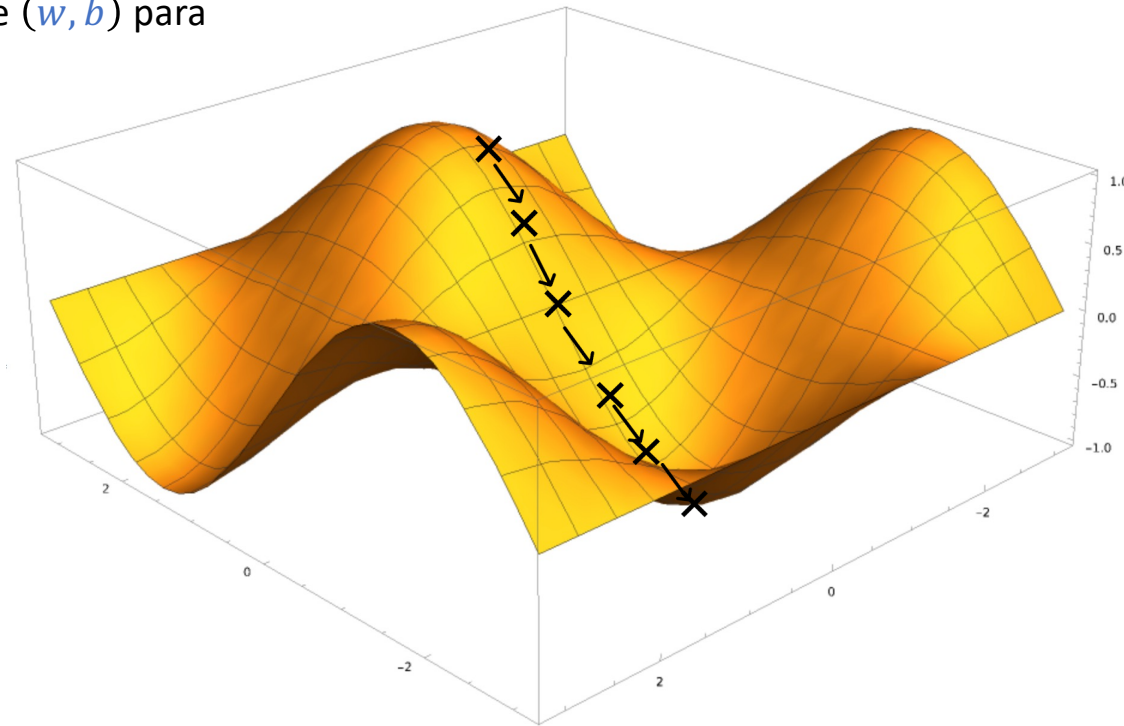


# Gradient descent - intuición

Algoritmo Gradient descent:

Empezamos con un par de valores  $(w, b)$ .

Cambiamos los valores de  $(w, b)$  para reducir  $J(w, b)$



# Gradient descent – dirección máximo crec.

Gradiente de  $J(w, b)$  me da la dirección de máximo crecimiento/decrecimiento de la función

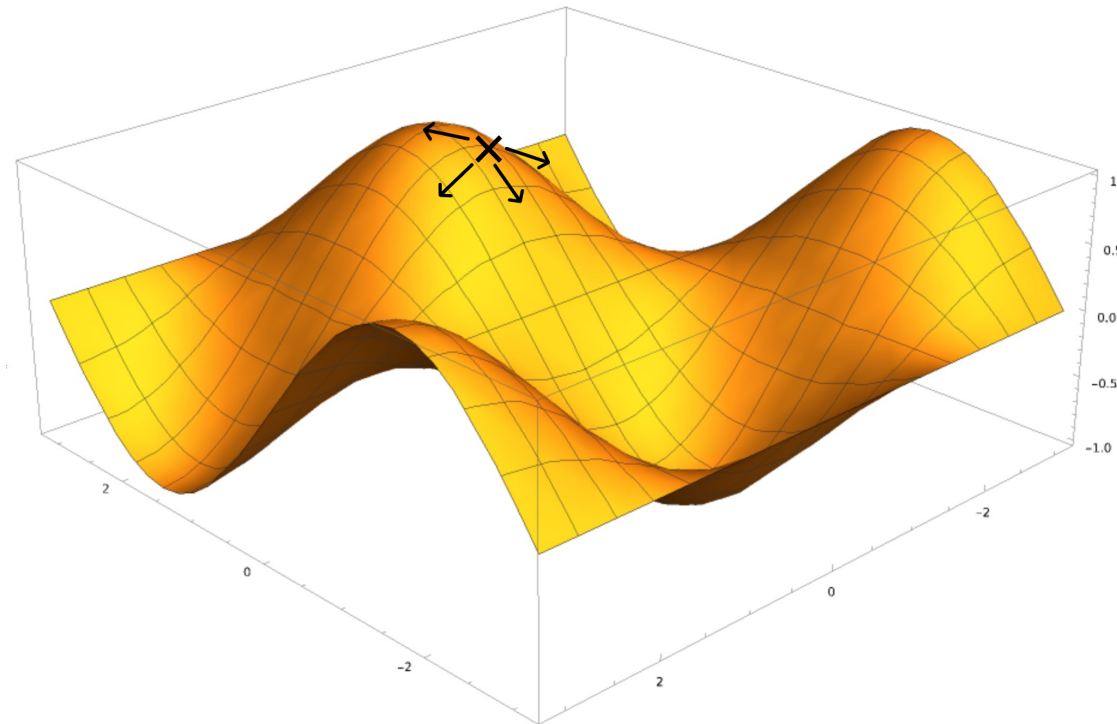
$$\nabla J(w, b) = \left( \frac{\partial J}{\partial w}, \frac{\partial J}{\partial b} \right)$$

Similitudes con el método de Newton (lo veremos más adelante)

Repetir hasta convergencia:

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$



GONZALO RUBIO CALZADO

# Gradient descent - problemas

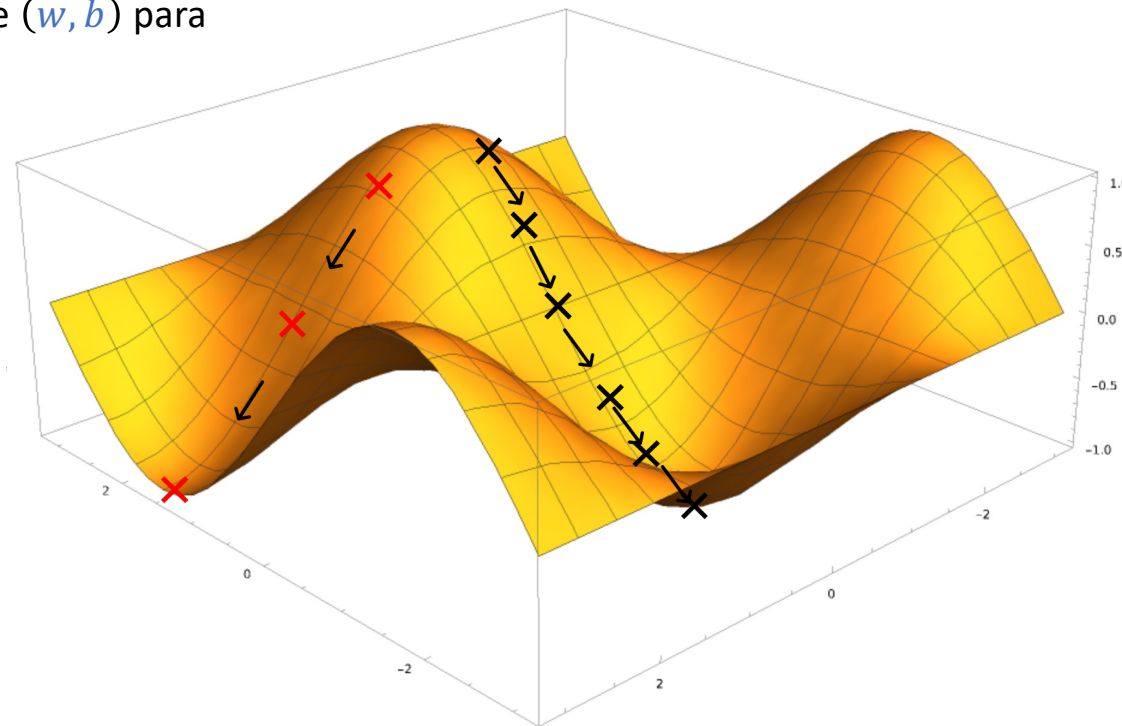
Algoritmo Gradient descent:

Empezamos con un par de valores  $(w, b)$ .

Cambiamos los valores de  $(w, b)$  para reducir  $J(w, b)$

Posibles problemas:

- Mínimo relativo (solo uno)
- Influencia condición inicial





# Gradient descent - algoritmo

$$w = w_0$$

$$b = b_0$$

Repetir hasta convergencia:

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$\alpha$  learning rate  
(tasa de aprendizaje)

Número positivo pequeño.  
Controla el tamaño de los pasos (la distancia entre las cruces).

Convergencia: cuando  $w$  y  $b$  apenas cambian entre pasos.

Parámetros  $w$  y  $b$  se actualizan a la vez.  
¡No actualizar  $J(w, b)$  con nuevo valor de  $w$  al actualizar  $b$ !

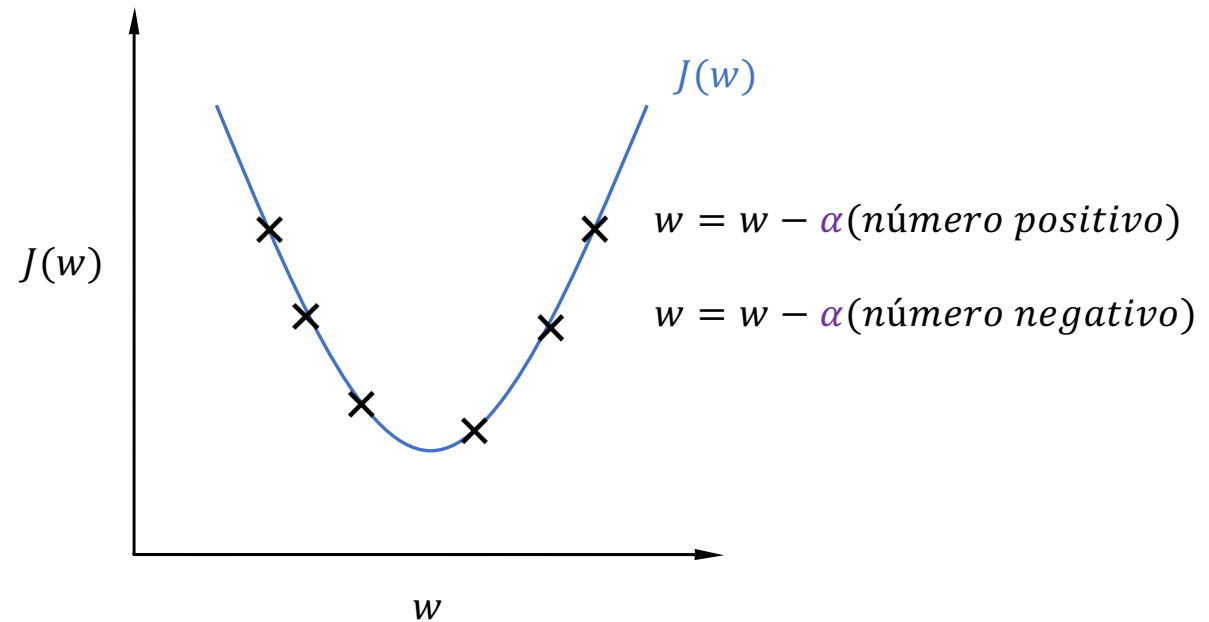
# Gradient descent – una variable

$$w = w_0$$

Repetir hasta convergencia:

$$w = w - \alpha \frac{d}{dw} J(w)$$

$\alpha$  learning rate  
(tasa de aprendizaje)



# Gradient descent – learning rate

$$w = w_0$$

Repetir hasta convergencia:

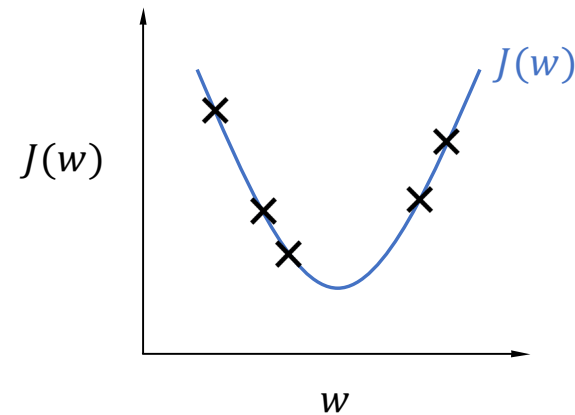
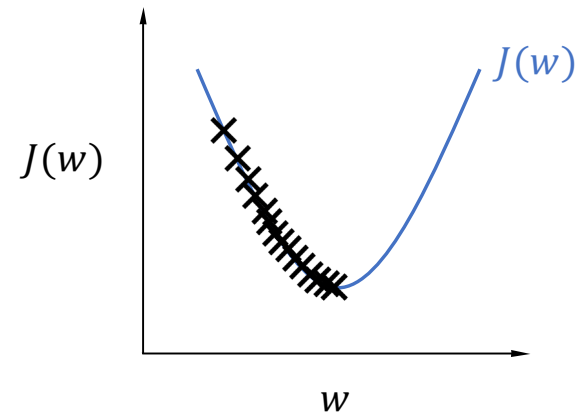
$$w = w - \alpha \frac{d}{dw} J(w)$$

Si  $\alpha$  es demasiado pequeño...

La convergencia será lenta

Si  $\alpha$  es demasiado grande...

Fallo en la convergencia, divergencia



# Gradient descent – mínimos cuadrados

Modelo lineal

$$f_{w,b}(x) = wx + b$$

Función de coste – error al cuadrado

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Algoritmo gradient descent

$$w = w_0$$

$$b = b_0$$

Repetir hasta convergencia:

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$\frac{\partial}{\partial w} J(w, b) = \frac{2}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial w} f_{w,b}(x^{(i)})$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{2}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial b} f_{w,b}(x^{(i)})$$

# Gradient descent – mínimos cuadrados

Modelo lineal

$$f_{w,b}(x) = wx + b$$

Algoritmo gradient descent

$$w = w_0$$

$$b = b_0$$

Repetir hasta convergencia:

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

```
46 # Gradient descent
47
48 alpha = 0.001
49
50 w = 0
51 b = 0
52
53 # Number of training steps
54 for i in range(1,100000):
55     wnew = w - alpha * dJdw(w,b,x,y)
56     bnew = b - alpha * dJdb(w,b,x,y)
57     w = wnew
58     b = bnew
```

¡Parámetros  
 $w$  y  $b$  se  
actualizan a  
la vez!

# Gradient descent – mínimos cuadrados

Modelo lineal

$$f_{w,b}(x) = wx + b$$

```
18 # Model
19 def fwb(w,b,x):
20     return w*x+b
21
22 # dJ/dw
23 def dJdw(w,b,x,y):
24     dJdw = 0
25     for i in range(m):
26         dJdw = dJdw + (fwb(w,b,x[i])-y[i]) * x[i]
27     dJdw = dJdw / (m)
28     return dJdw
29
30 # dJ/db
31 def dJdb(w,b,x,y):
32     dJdb = 0
33     for i in range(m):
34         dJdb = dJdb + (fwb(w,b,x[i])-y[i])
35     dJdb = dJdb / (m)
36     return dJdb
```

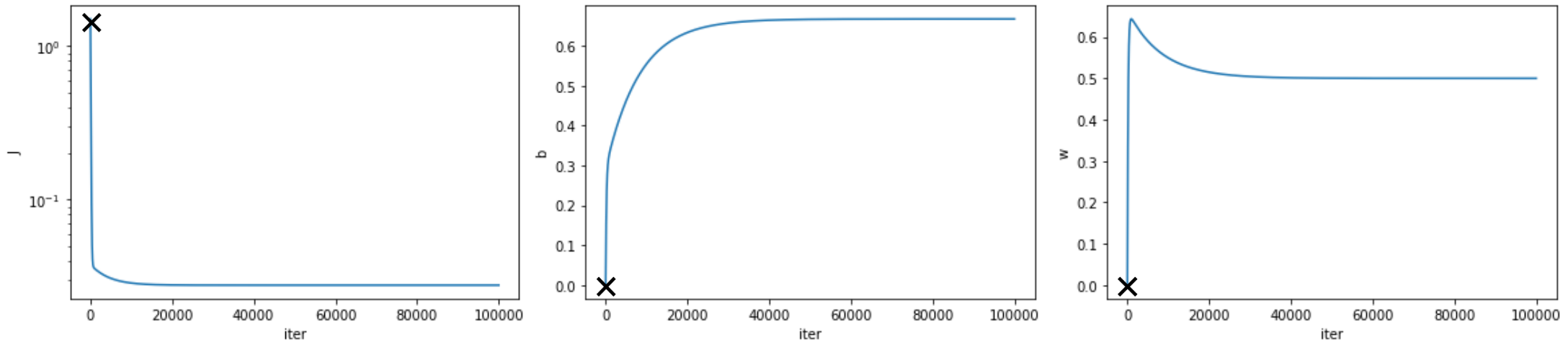
Función de coste – error al cuadrado

$$J(w,b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial w} J(w,b) = \frac{2}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

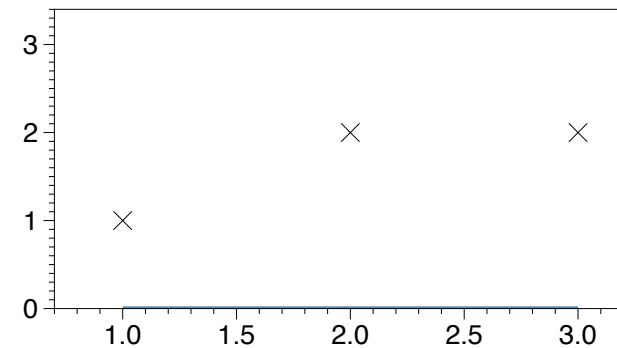
$$\frac{\partial}{\partial b} J(w,b) = \frac{2}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

# Gradient descent – mínimos cuadrados

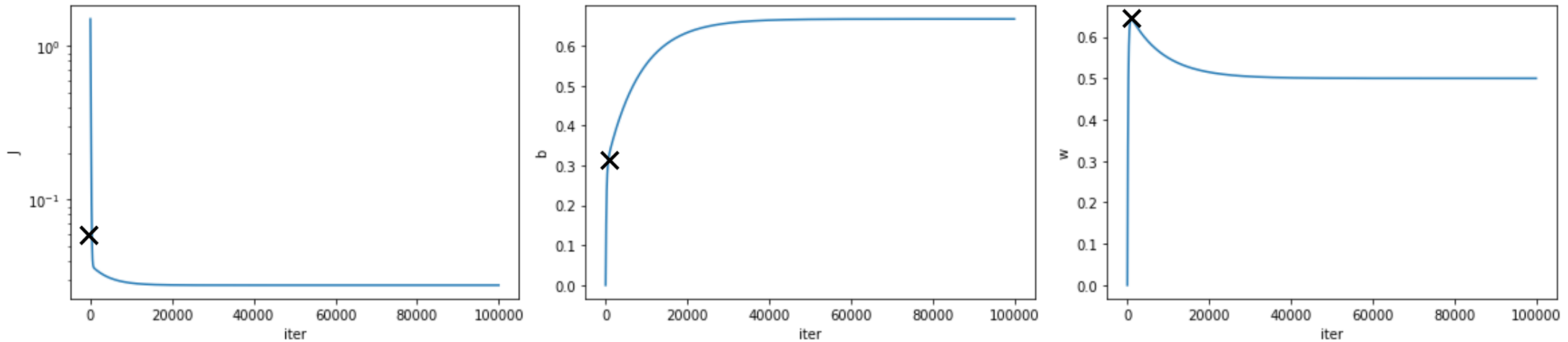


Inicio:

$$f_{w,b}(x) = wx + b = 0x + 0$$

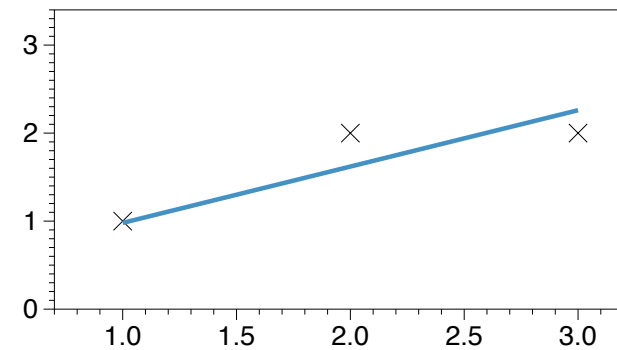


# Gradient descent – mínimos cuadrados



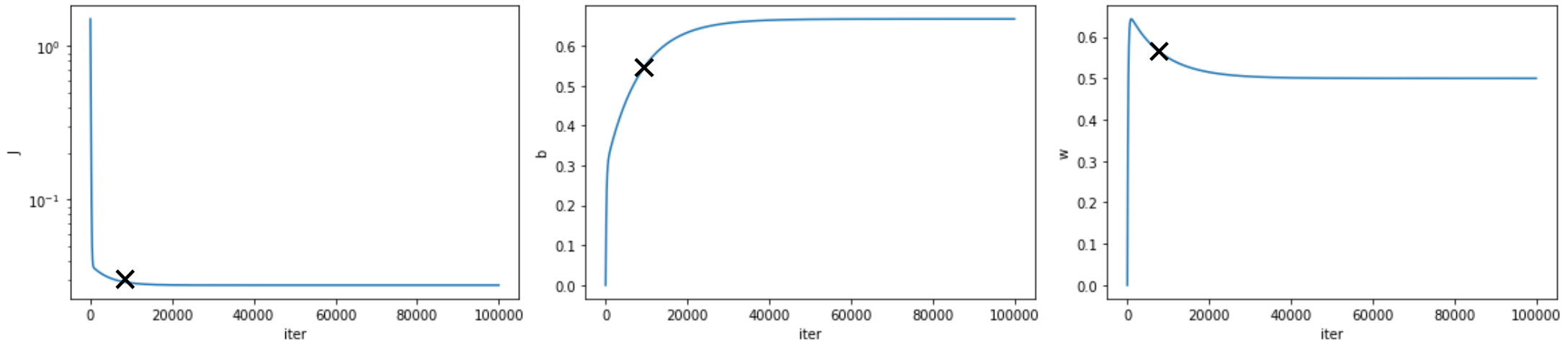
Tras 1000 iteraciones:

$$f_{w,b}(x) = wx + b = 0.64x + 0.34$$



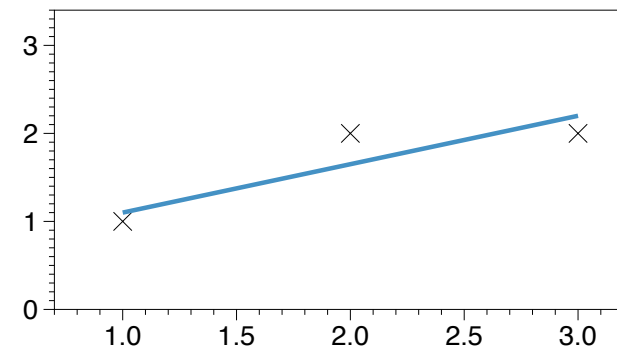


# Gradient descent – mínimos cuadrados

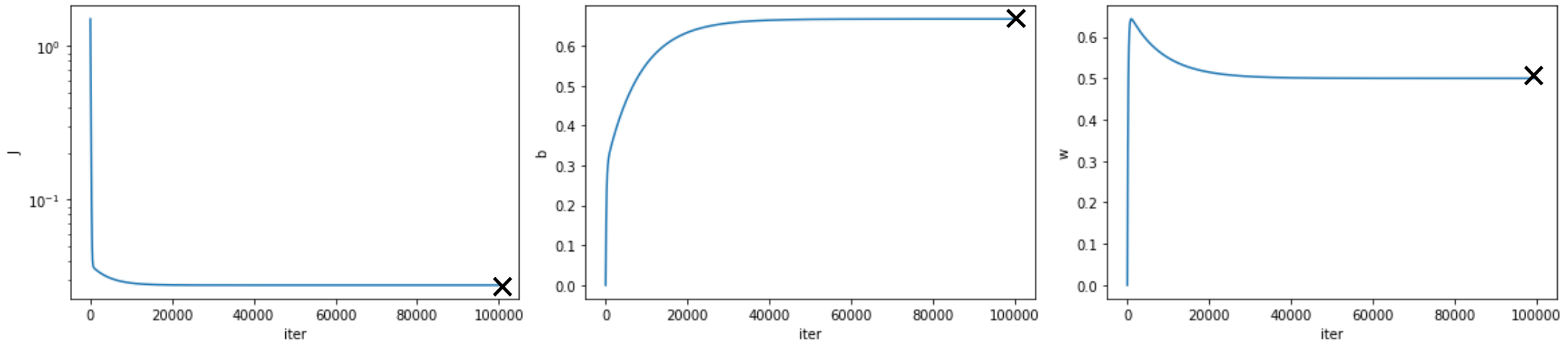


Tras 10000 iteraciones:

$$f_{w,b}(x) = wx + b = 0.55x + 0.55$$



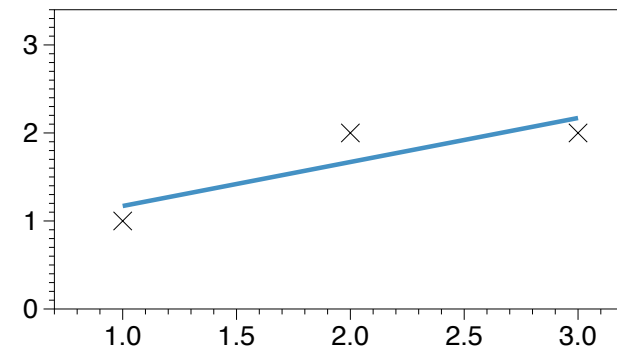
# Gradient descent – mínimos cuadrados



Tras 100000 iteraciones:

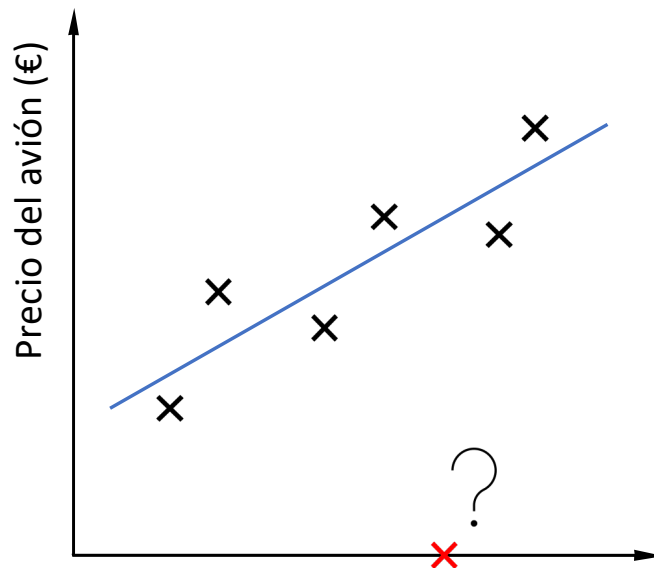
$$f_{w,b}(x) = wx + b = 0.5x + 0.67$$

¡Llego al mismo resultado!



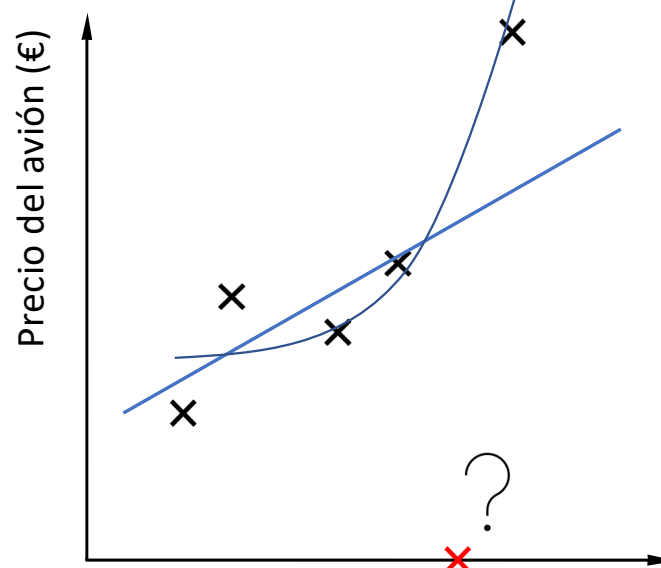
# Conclusiones

La técnica de mínimos cuadrados permite realizar predicciones basadas en datos



Tamaño del avión  
(max número de pasajeros)

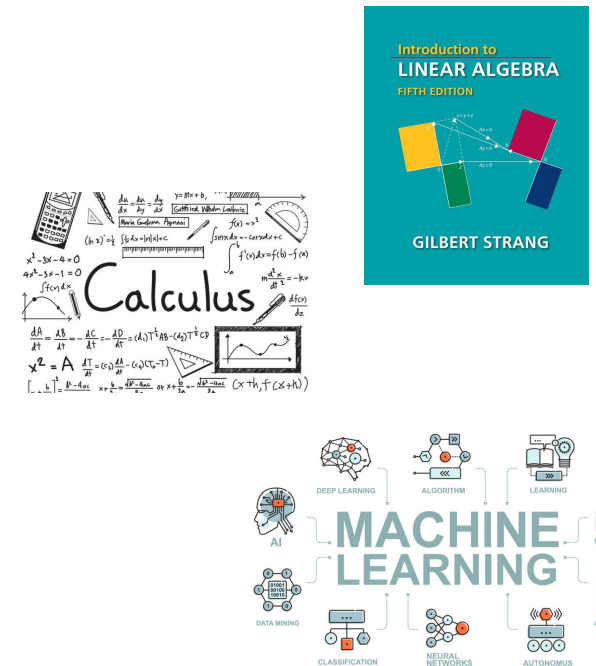
El modelo puede ser lineal o no lineal



Tamaño del avión (max número de pasajeros)

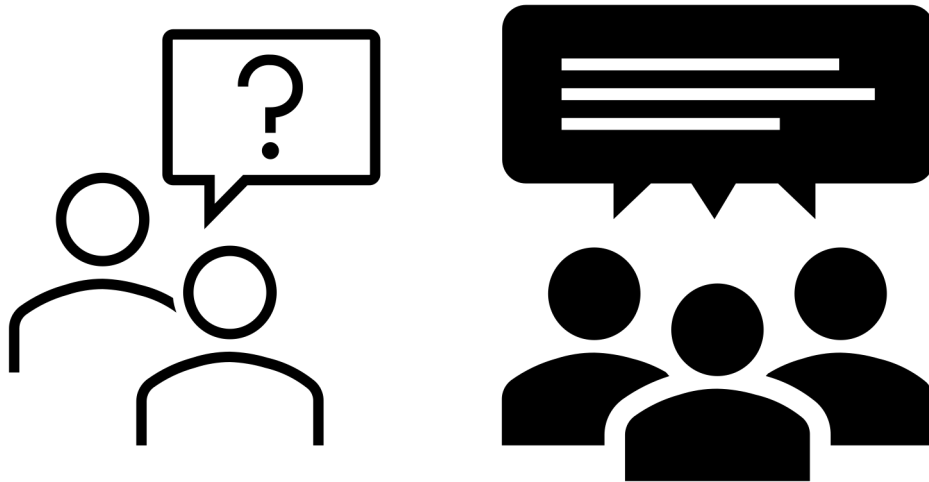
GONZALO RUBIO CALZADO

Puedo llegar al mismo resultado final por distintos caminos





Gracias por su atención. ¿Preguntas?

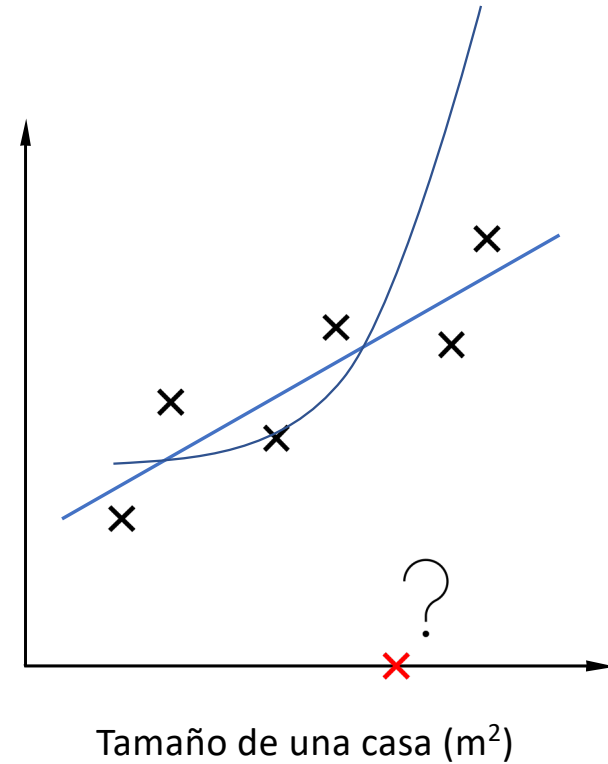


MATERIAL ADICIONAL

# Predicción basada en datos

- ¿Y si quiero aproximar con algo que no sea una recta?
  - Camino 1: Álgebra lineal
  - Camino 2: Cálculo (de varias variables)

Precio de la casa (€)



# Más allá de la recta. Álgebra lineal

$$(1,1), (2,2), (3,2)$$

$$1 = C + D + E$$

$$2 = C + 2D + 4E$$

$$2 = C + 3D + 9E$$

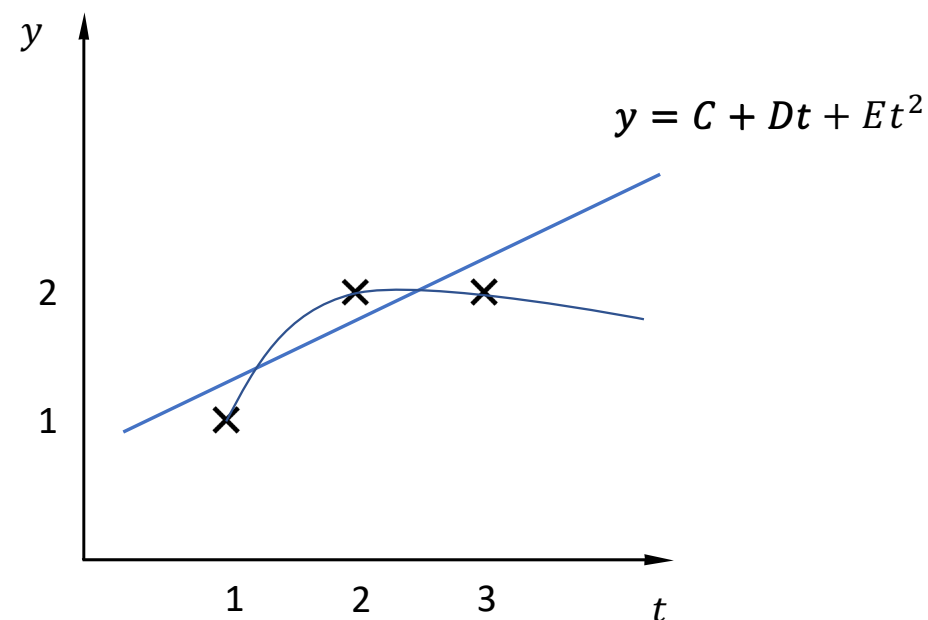
Este sistema tiene mejor pinta

$$Ax = b$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix} \begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}$$

En este caso, la función puede pasar por todos los puntos (Sistema tiene solución)

$$\begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} -1 \\ 5/2 \\ -1/2 \end{bmatrix}$$



Este caso particular se llama interpolación. Se verá en detalle más adelante.

# Más allá de la recta. Álgebra lineal

$$(1,1), (2,2), (3,2), (4,2)$$

$$1 = C + D + E$$

$$2 = C + 2D + 4E$$

$$2 = C + 3D + 9E$$

$$2 = C + 4D + 16E$$

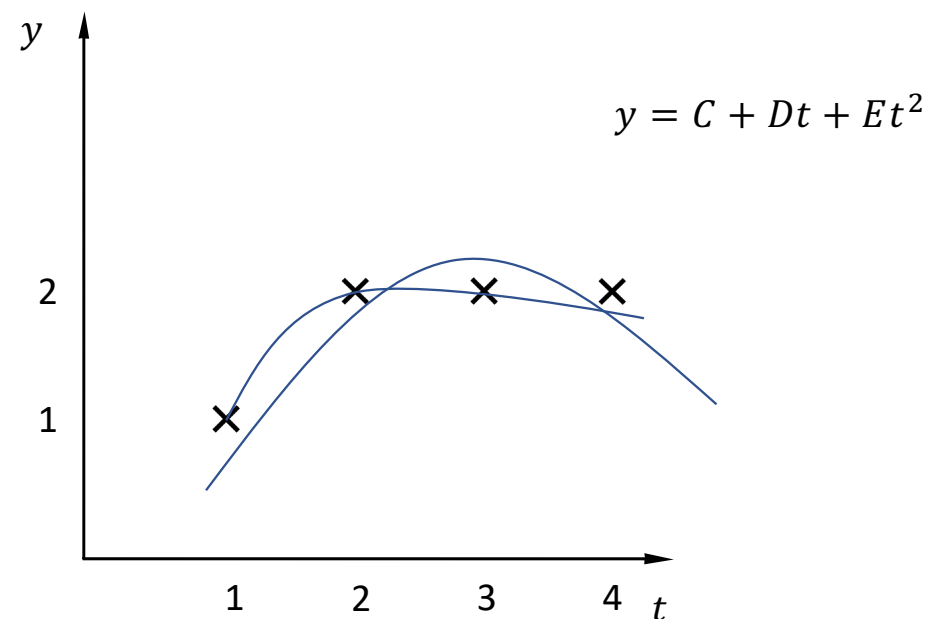
Volvemos a tener problemas

$$A\mathbf{x} = \mathbf{b}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} \begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

Y proponemos la misma solución

$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$$

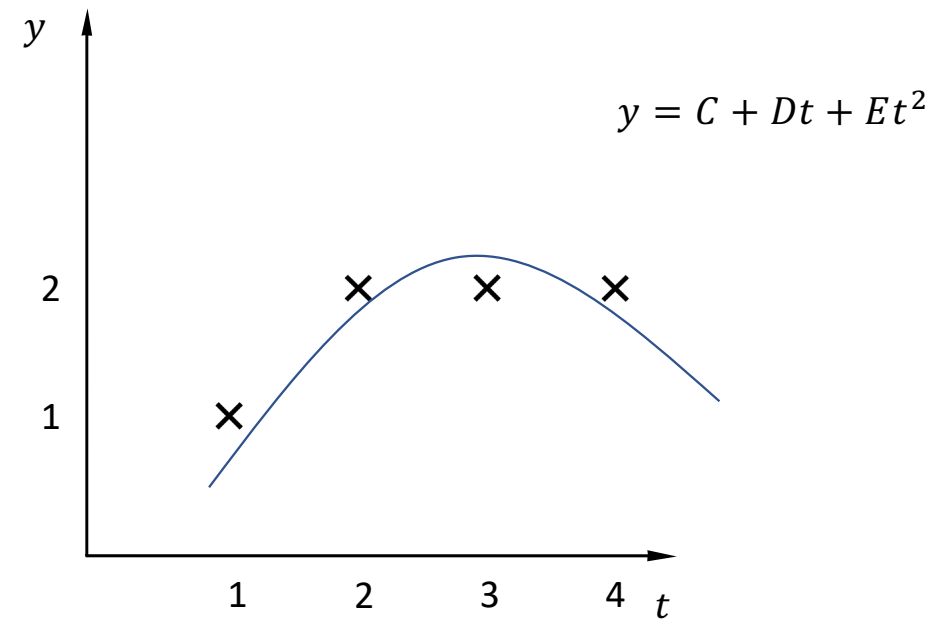




# Más allá de la recta. Álgebra lineal. Generalización

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} \begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ 1 & t_3 & t_3^2 \\ 1 & t_4 & t_4^2 \end{bmatrix} \begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

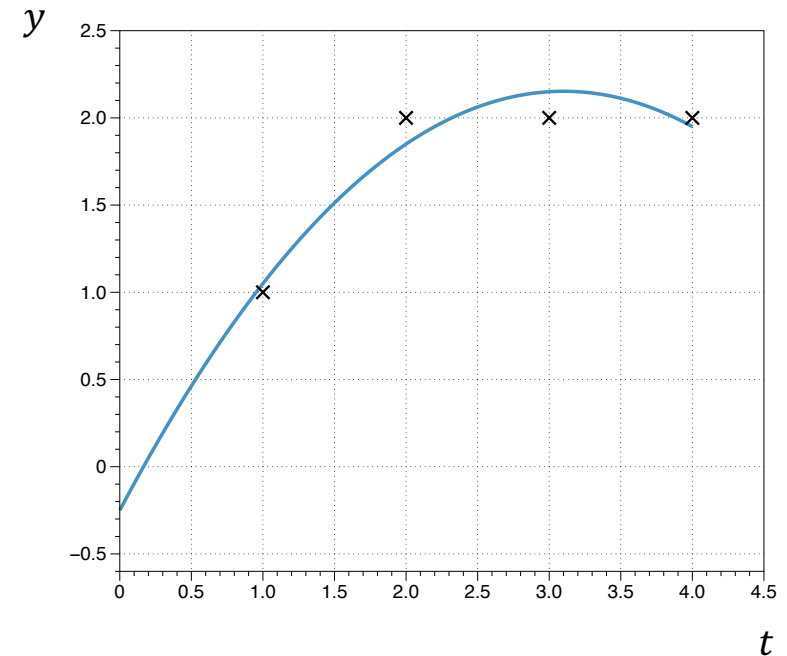


# Más allá de la recta. Álgebra lineal

$$A^T A \hat{x} = A^T b$$

```
8  import numpy as np
9
10
11  A = np.matrix([[1, 1, 1], [1, 2, 4], [1, 3, 9], [1, 4, 16]])
12  b = np.matrix([1,2,2,2])
13  b = b.transpose()
14
15  AT = A.transpose()
16
17  Ahat = np.matmul(AT,A)
18  bhat = np.matmul(AT,b)
19
20  x = np.linalg.solve(Ahat, bhat)
21
```

$$\begin{bmatrix} C \\ D \\ E \end{bmatrix} = \begin{bmatrix} -0.25 \\ 1.55 \\ -0.25 \end{bmatrix}$$



# Gradient descent – learning rate fijo

$$w = w_0$$

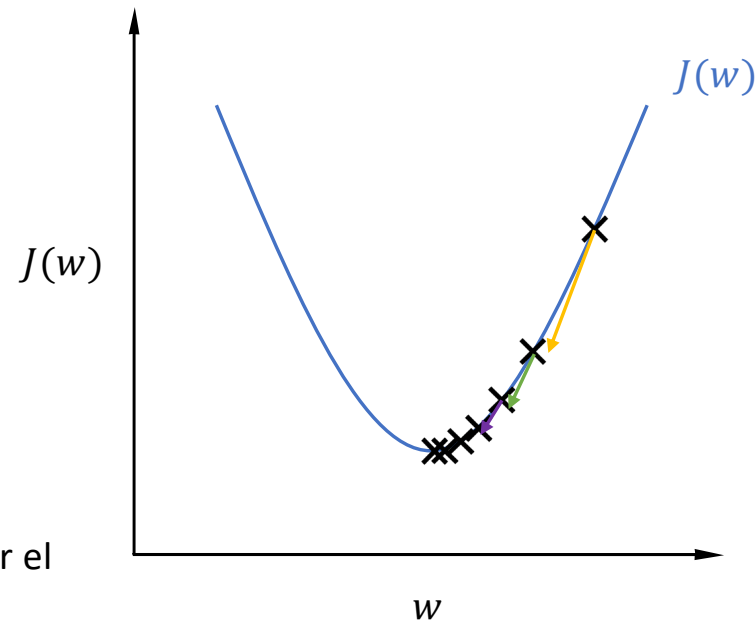
Repetir hasta convergencia:

$$w = w - \alpha \frac{d}{dw} J(w)$$

Cerca de un mínimo local:

- La derivada se hace más pequeña
- Los pasos se hacen más pequeños

Algoritmo alcanza el mínimo sin necesidad de variar el learning rate  $\alpha$



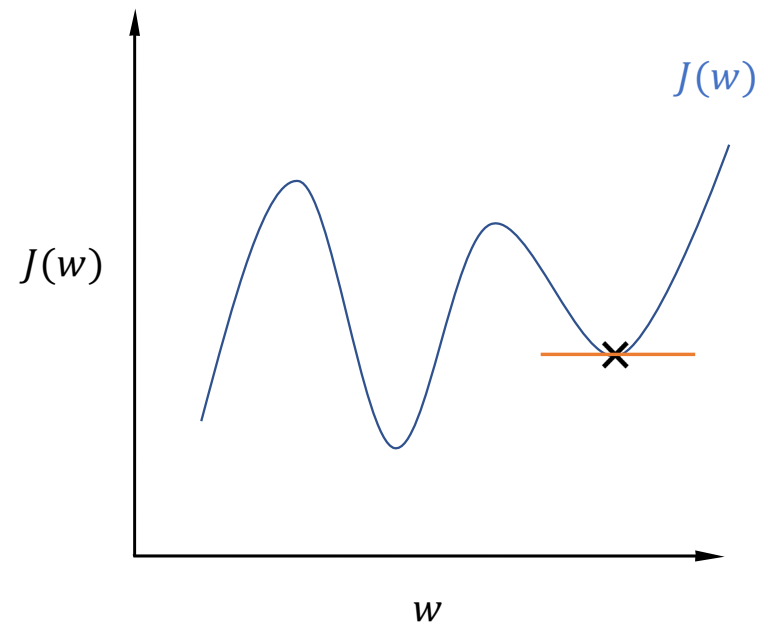
# Gradient descent – learning rate

$$w = w_0$$

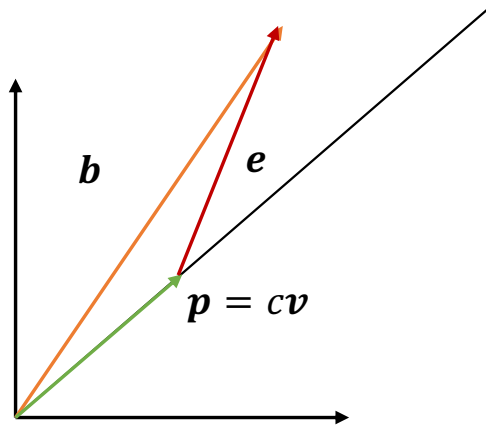
Repetir hasta convergencia:

$$w = w - \alpha \frac{d}{dw} J(w)$$

Si ya estoy en un mínimo local, el algoritmo deja de evolucionar



# ¿EQUIVALENCIA ALGEBRA Y CÁLCULO?



El objetivo en ambos casos es minimizar la distancia entre  $b$  y  $p$ :

Algebra:

La distancia será mínima si  $e$  es ortogonal a  $p$

$$(b - p) \cdot p = 0$$

$$(b - cv) \cdot cv = 0$$

$$c = \frac{b \cdot v}{v \cdot v}$$

Cálculo:

Minimizo la función distancia  $\|b - p\|^2 = (b - cv) \cdot (b - cv)$

$$b \cdot b + c^2 v \cdot v - 2cb \cdot v$$

$$\frac{d\|b - p\|^2}{dc} = 2cv \cdot v - 2b \cdot v = 0$$

$$c = \frac{b \cdot v}{v \cdot v}$$

