

Regresión Lineal



Regresión Lineal y Gradient Descent

Breve Recordatorio

- n número de observaciones o registros
- p número de variables o features

Ejemplo:

$n=5000$ habitantes y $p=12$ características (features)

Breve Recordatorio

- Una variable del tipo x_{ij} , $x_{i,j}$, $x_j^{(i)}$ indicaría el valor del feature j para la observación i ; es decir:

$$\begin{aligned}i &= 1, 2, \dots, n \\ j &= 1, 2, \dots, p\end{aligned}$$

En forma de matriz: \mathbf{X}_{np}

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

Breve Recordatorio

	sqft	City	bedrooms	baths	price
x_i	3392	Dublin	3	2.1	339000
	4100	pleasanton	4	3	899900
	3200	Clayton	5	4	448641
	1436	Moraga	4	3	239999
	1944	Antioch	3	2	377500
	1500	Danville	3	2.5	299900
	1700	El Dorado Hills	4	3	265000
	2507	Shingle Springs	4	3	449000
	1580	McKinleyville	3	2	439950
	1500	Marina	4	2	699888
	2705	Roseville	3	2	1250000
	1715	Rocklin	4	3	439000

b)

$$x_i = \begin{bmatrix} 3392 \\ \text{Dublin} \\ 3 \\ 2.1 \end{bmatrix}$$

c)

$$x_j = \begin{bmatrix} 3392 \\ 4100 \\ 3200 \\ 1436 \\ 1944 \\ 1500 \\ 1700 \\ 2507 \\ 1580 \\ 1500 \\ 2705 \\ 1715 \end{bmatrix}$$

d)

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_p \end{bmatrix} = \begin{bmatrix} \text{sqft} & \text{City} & \text{bedrooms} & \text{baths} \\ 3392 & \text{Dublin} & 3 & 2.1 \\ 4100 & \text{pleasanton} & 4 & 3 \\ 3200 & \text{Clayton} & 5 & 4 \\ 1436 & \text{Moraga} & 4 & 3 \\ 1944 & \text{Antioch} & 3 & 2 \\ 1500 & \text{Danville} & 3 & 2.5 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = \begin{bmatrix} 3392 & 4100 & 3200 & 1436 & \dots \\ \text{Dublin} & \text{pleasanton} & \dots & & \\ 2.1 & 3 & & 4 & \end{bmatrix}$$

Breve Recordatorio

- $\mathbf{x}^{(i)}$ o \mathbf{x}_i : Variable de entrada, input feature, donde i es el item
- $\mathbf{y}^{(i)}$ o \mathbf{y}_i : Variable de salida, target variable
- $x_{ij}, x_{i,j}, x_j^{(i)}$: Dato particular de la observación / item
- Data de entrenamiento: $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$

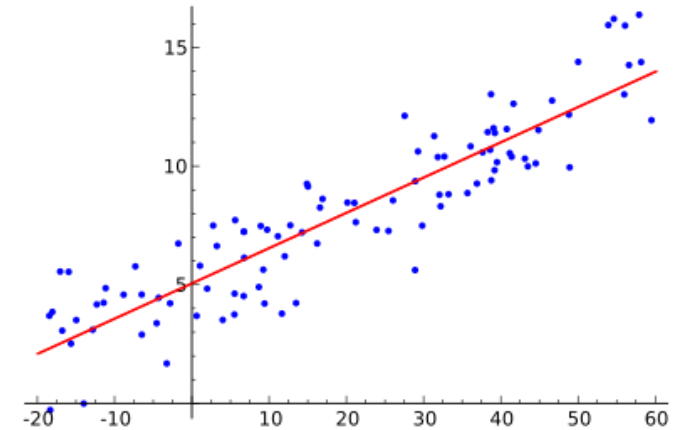
Regresión Lineal

- Es un método de aprendizaje supervisado
- Nos ayuda a predecir valores

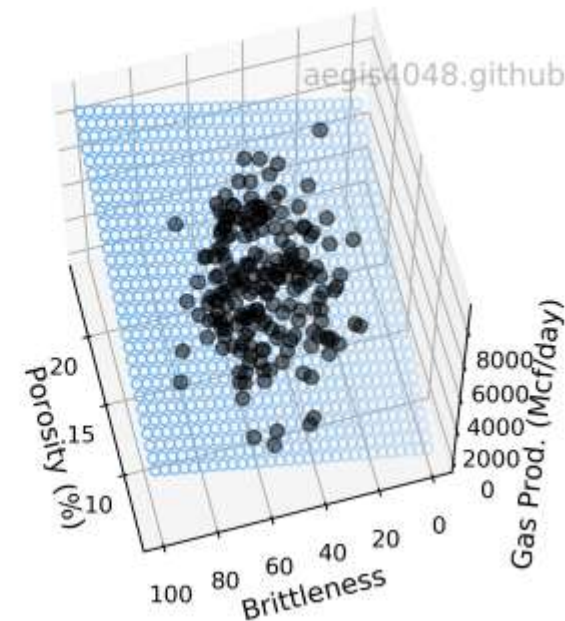
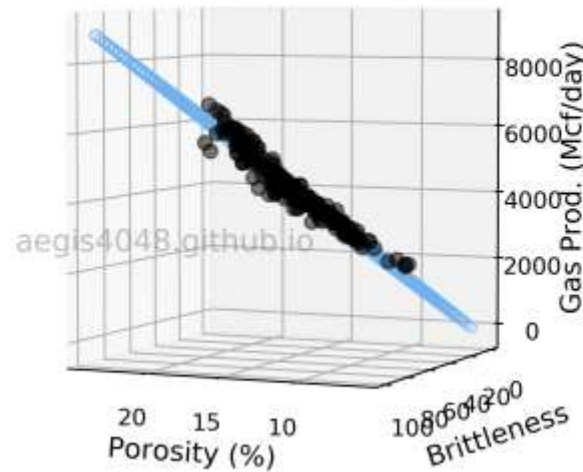
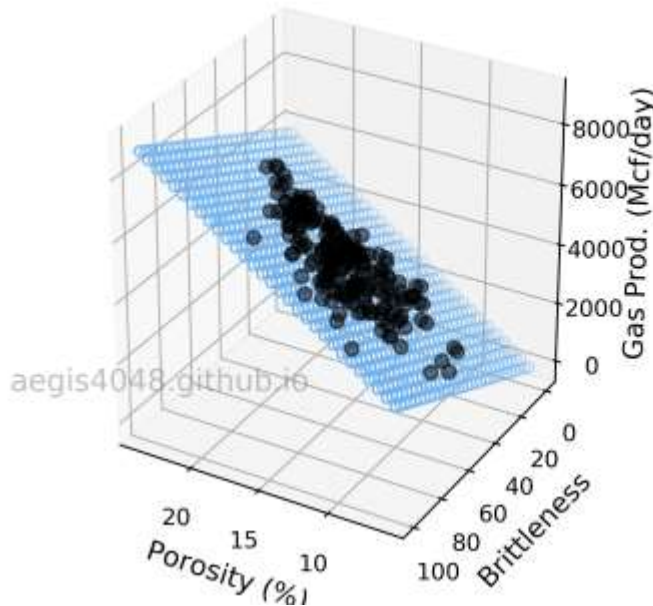


Regresión Lineal

- Relación lineal de una variable independiente con una o más variables dependientes



$$R^2 = 0.93$$



Regresión Lineal

Una sola variable independiente: $y = ax + b$

Dos variables independiente: $z = ax + by + c$

Nos quedamos sin letras.... Y estamos cambiando letras a cada rato

Generalizar

$$f(x) = w_0 + w_1x_1 + w_2x_2 \dots w_nx_n$$

Regresión Lineal en ecuaciones

$$f(x) = w_0 + w_1x_1 + w_2x_2 \dots w_nx_n$$

$$f(x) = \sum_{i=0}^n w_i x_i$$

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} ; X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Aplicando multiplicación de matrices y recordando el producto escalar

$$W = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} ; X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$f(x) = W^T X$$

Queremos encontrar los coeficientes

$$f(w, x) = W^T X$$

De tal forma que

$$f(w, x_i) \approx y_i$$

Ejemplo

$$f(w, x) = W^T X$$

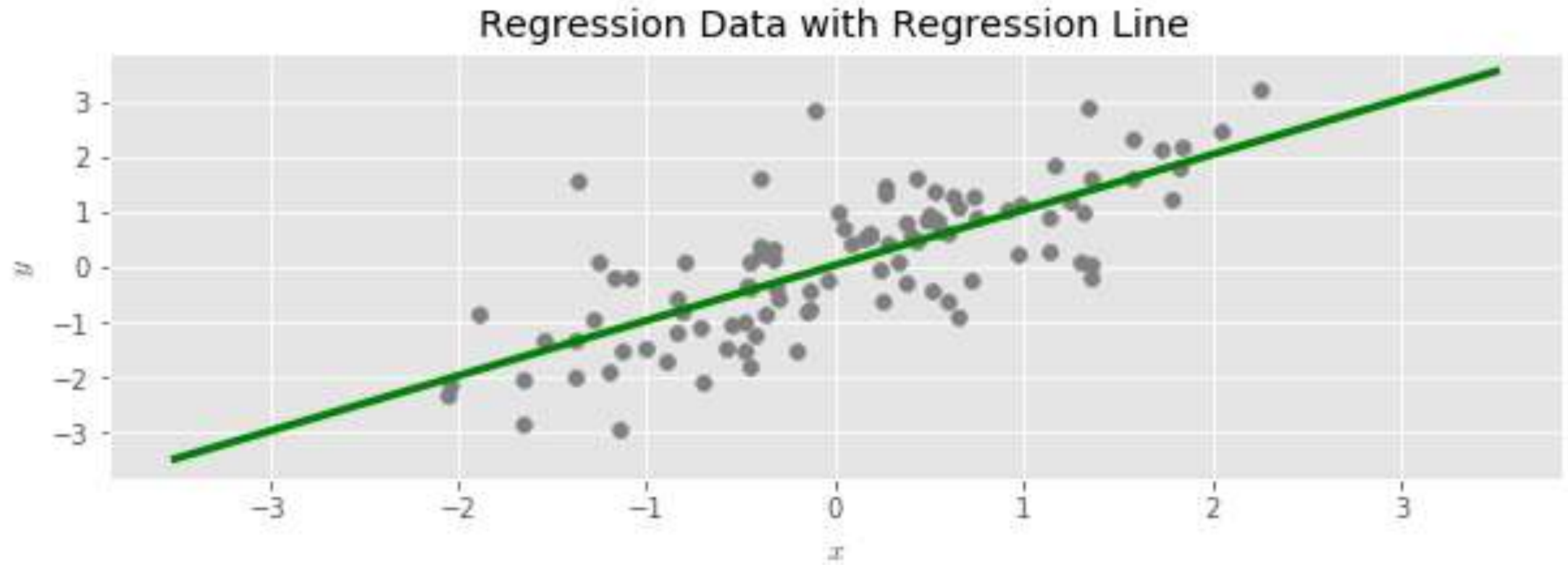
sqft	City	bedrooms	baths	price
3392	Dublin	3	2.1	339000
4100	pleasanton	4	3	899900
3200	Clayton	5	4	448641

$$price = w_0 + w_1 sqft + w_2 bedroom + \dots$$

Esta es nuestra situación ideal



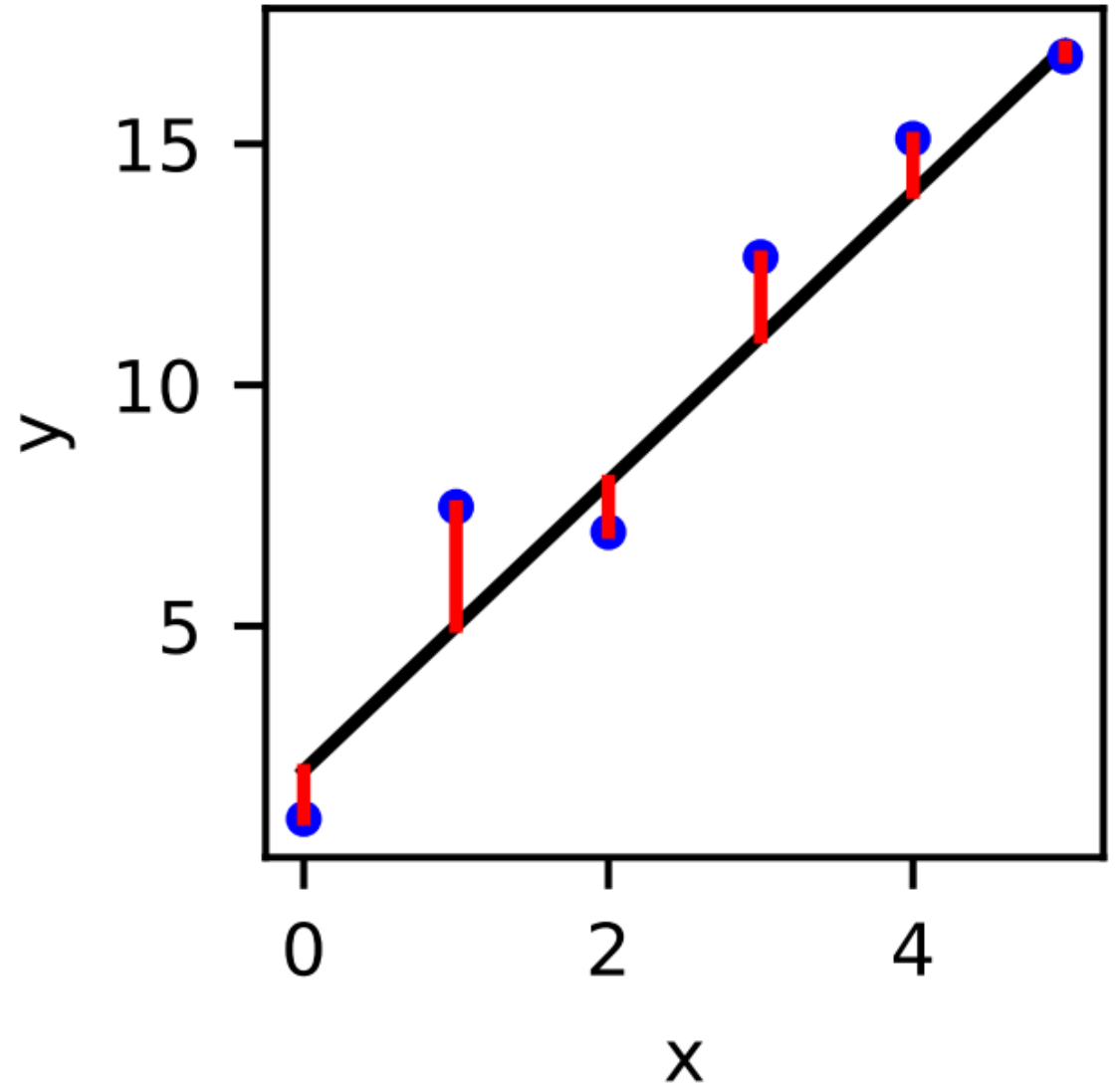
El mundo real es ruidoso



Residuales

- Los residuales, mostrado en rojo, son una métrica del fit de una función (línea negra) con respecto a un conjunto de data (azul)

$$r_i(w) = y_i - f(w, x_i)$$



Una forma simple es un problema de optimización. Definiendo una función de pérdida

$$\mathcal{L}(w)$$

Queremos encontrar los w que minimizen el “error”

$$\mathbf{w}^* = \min_w \mathcal{L}(w)$$

- Definiendo la función de pérdida en función a los Mínimos Cuadrados (Least-Squares Error) (minimizar RSS, SSR, SSE....)

$$\mathcal{L}_{LSE}(W) = \sum_{i=1}^n r_i^2 = r^T r$$

*Recordemos que todas nuestras variables son vectores

- Reemplazando términos

$$\mathcal{L}_{LSE}(W) = \sum_{i=1}^n (y_i - W^T X_i)^2$$

X = variable independiente (features)

y = variable dependiente (target)

W = coeficientes o pesos

Intuitivamente

Mientras más cercano nuestra función se asemeje a los datos, menor es el costo

Mientras menos el costo, menos desviación entre la línea (función) y nuestros datos

Mejor predicción!

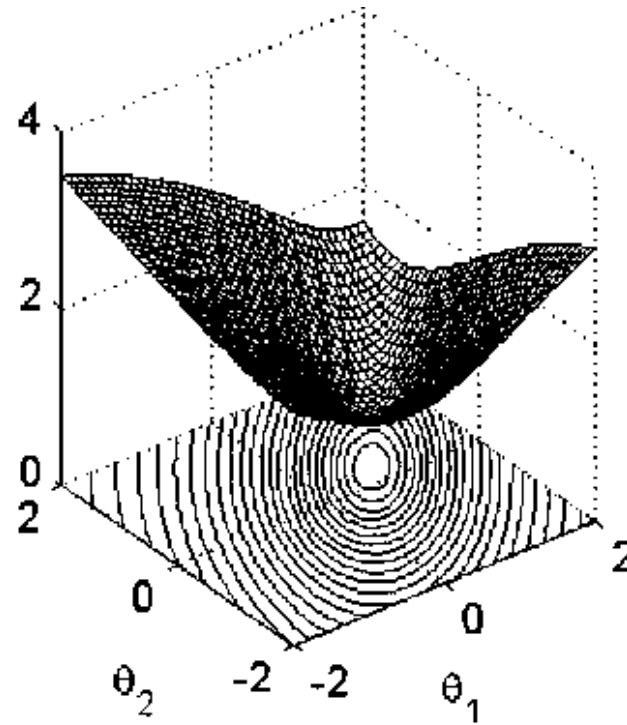
El problema

- ¿Cómo estimar los pesos W ?
- De tal forma que $f(w, x_i)$ sea cercano a y_i
- ¿Cómo minimizar la función pérdida ajustando solo W ?

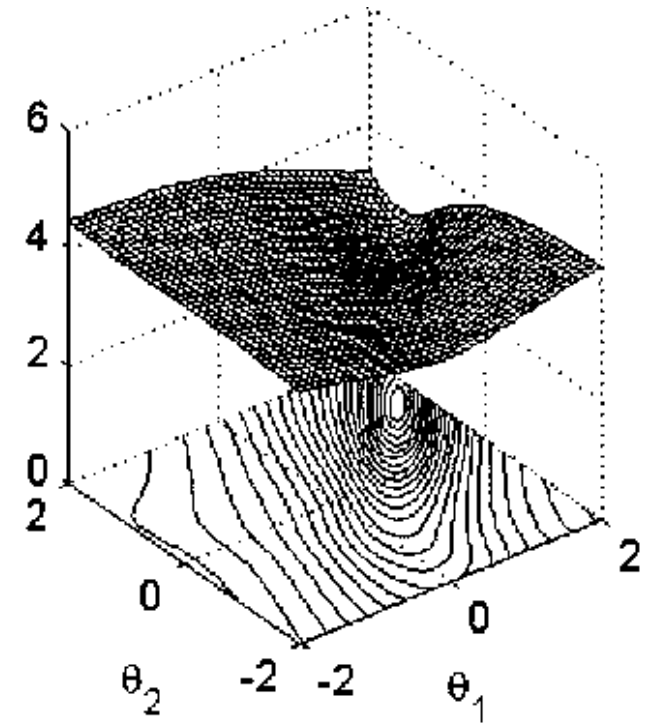


Problema

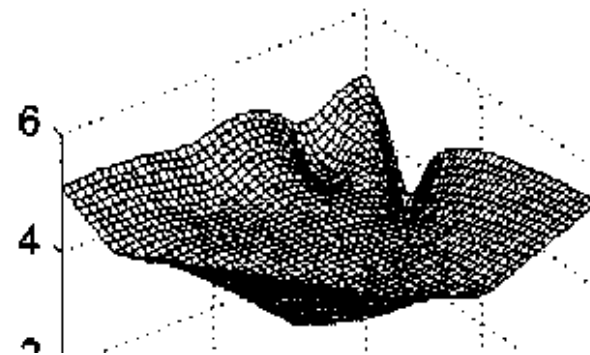
No sabemos que forma tiene
la función objetivo

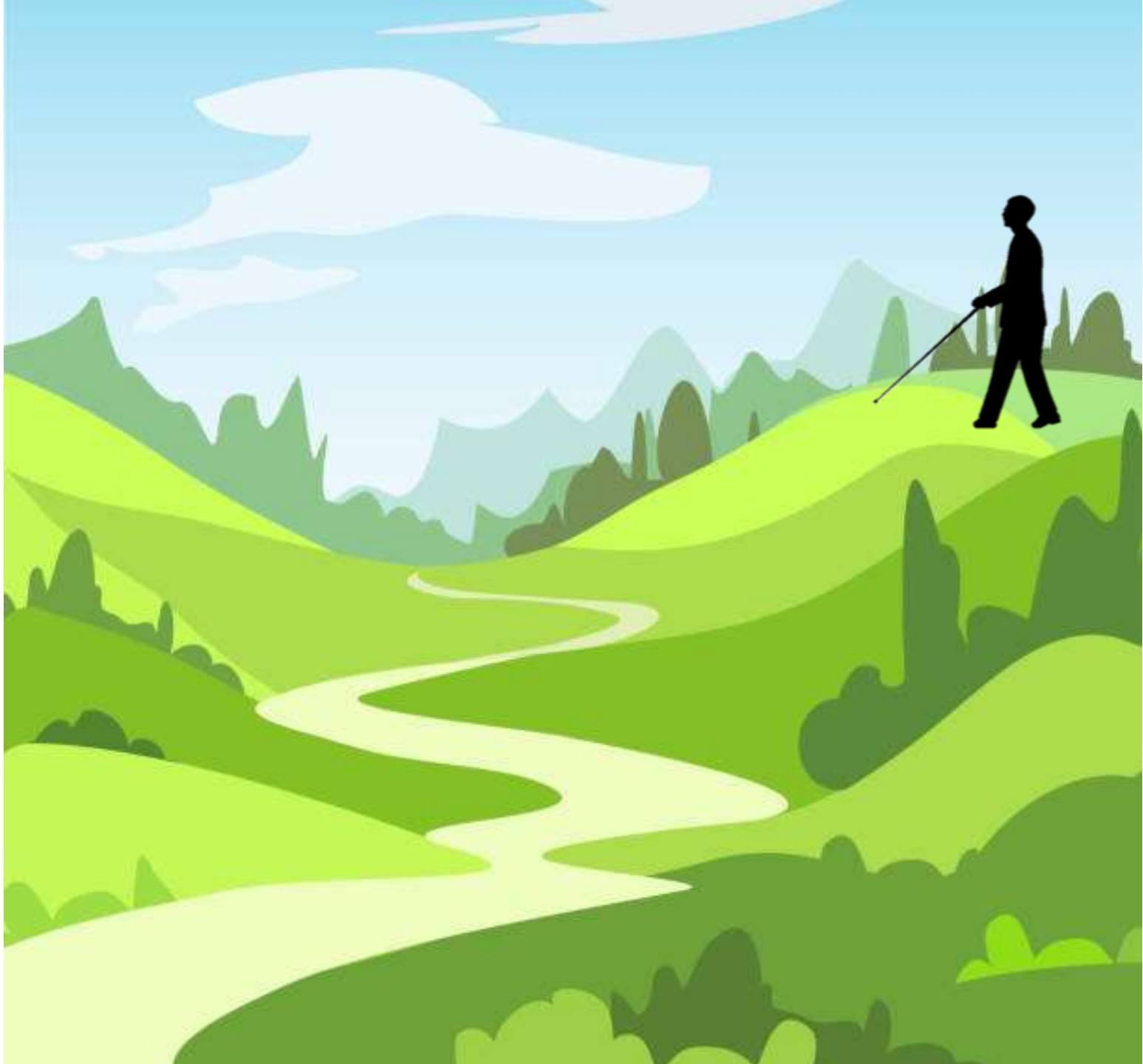


(a)



(b)





Preguntas a Resolver



Que tan empinado es donde
estoy parado



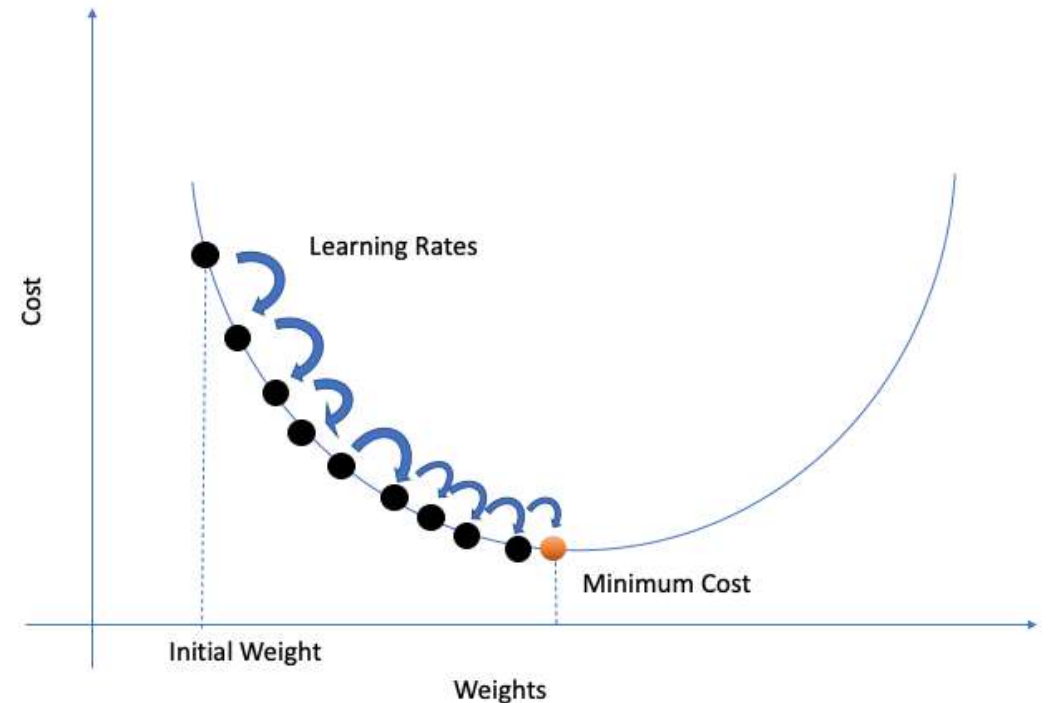
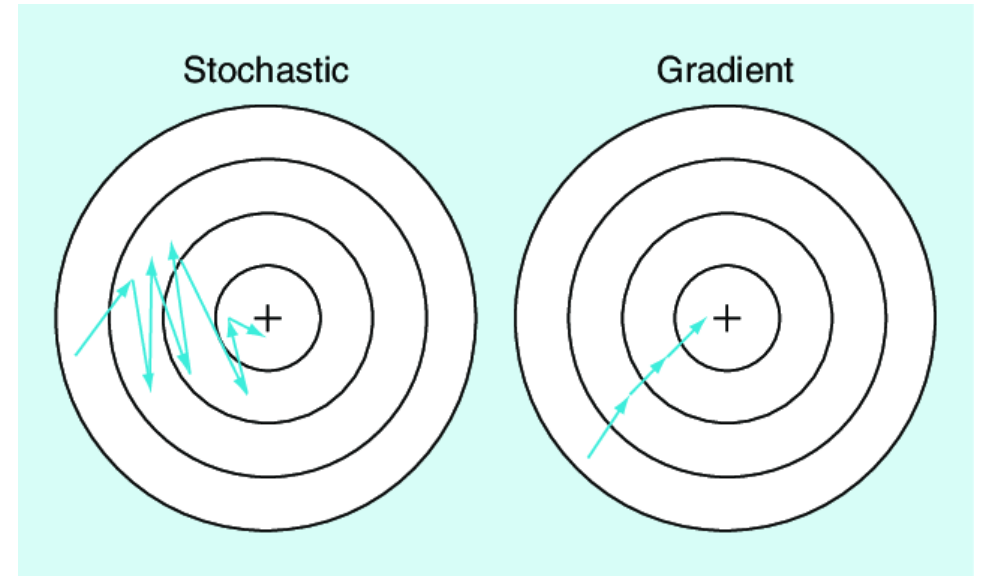
Hacia donde tengo que moverme

Puedo tomar
muchas medidas

Puedo tomar
una sola medida

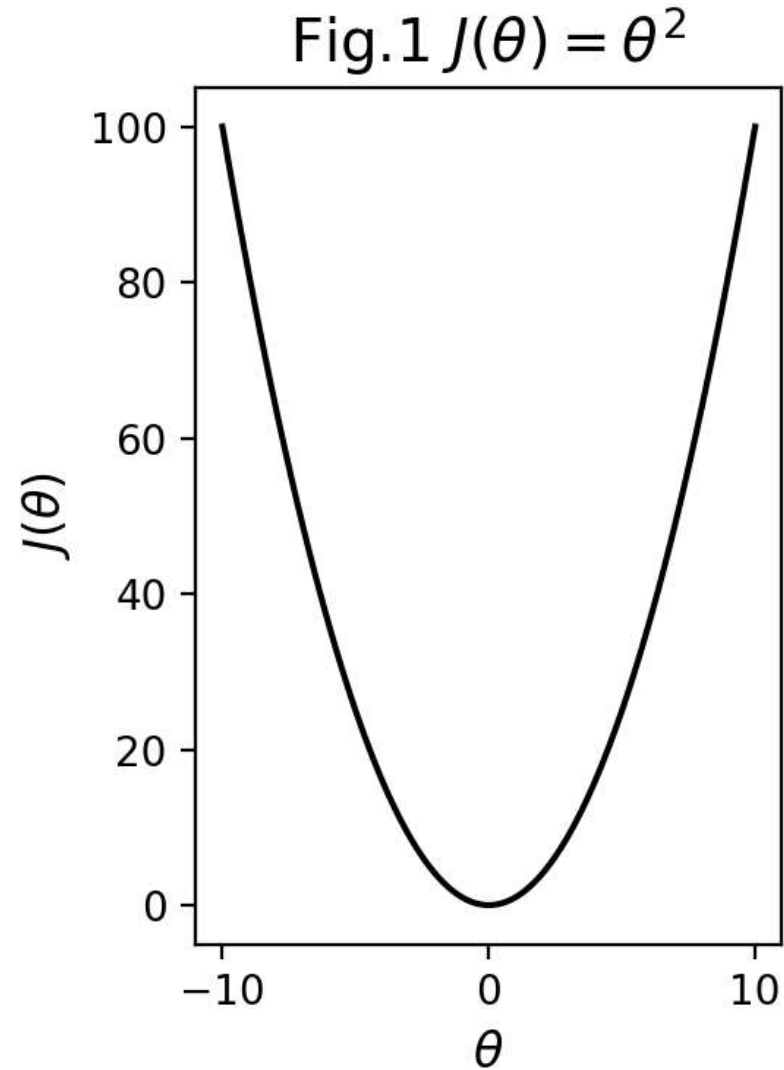
Gradient Descent

- Gradient Descent
 - La ruta más directa
 - Usa todos los data points y luego actualiza los pesos
 - Mejor cuando son pocos
- Stochastic Gradient Decent (SGD)
 - Va actualizando pesos con cada data point



Ejemplo Simple

$$J(\theta) = \theta^2$$

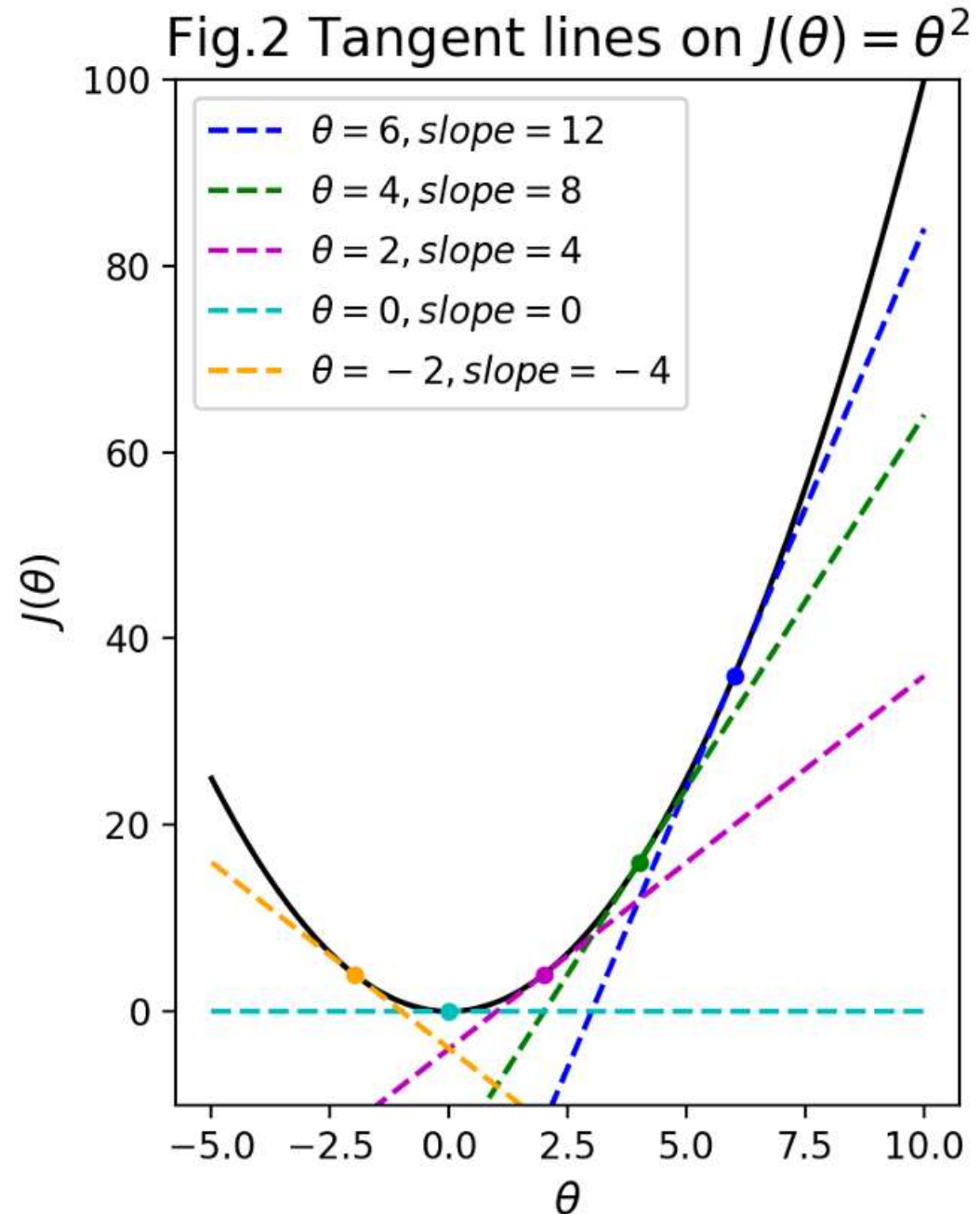


$$\min_{\theta} J(\theta)$$

- 1) ¿Cómo sabemos que $J(\theta)$ ha llegado al mínimo?
- 2) ¿Si $J(\theta)$ no está a su mínimo, como sabemos que valor de θ probar?

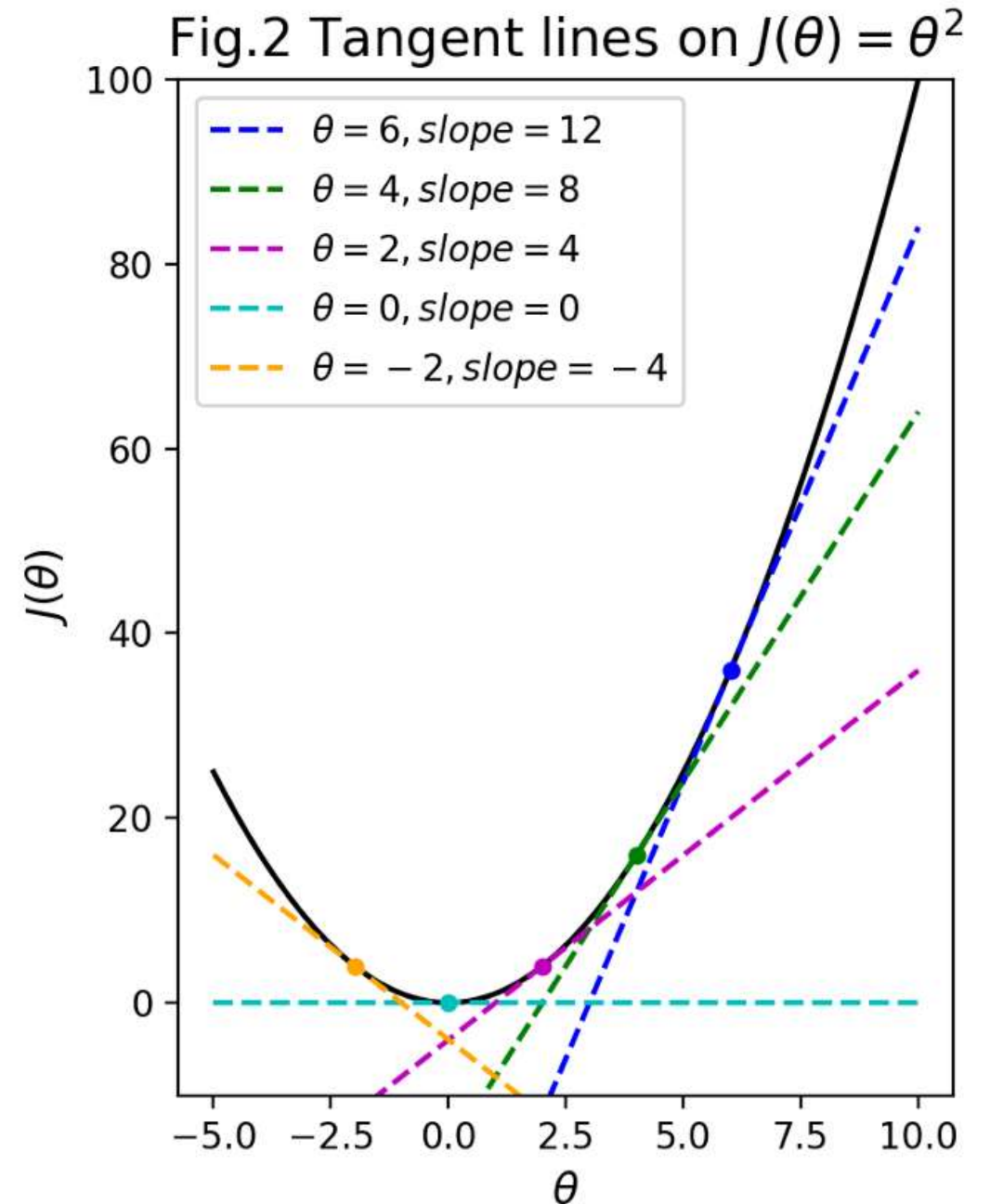
Derivadas!

- Si la pendiente es positiva, θ tiene que disminuir
- Si la pendiente es negativa, θ tiene que aumentar
- Si la pendiente es 0, felicidades! Has encontrado el mínimo (o máximo 🤖)



El Mínimo

- Si la pendiente es 0, o cambia muy poco (0.001), $J(\theta)$ ha convergido



¿Como cambiar θ ?

Regla de Actualización (Update Rule)

$$\theta = \theta - \alpha \frac{dJ}{d\theta}$$

α es conocido como “learning rate”

$\frac{\partial J}{\partial \theta}$ es la derivada de $J(\theta)$ con respecto a θ

Derivando

$$J(\theta) = \theta^2$$

$$\frac{dJ}{d\theta} = 2\theta$$

$$\theta = \theta - \alpha 2\theta$$

$$\theta = (1 - 2\alpha)\theta$$

Learning Rate

- Muy grande: El algoritmo no converge
- Muy pequeño: Demora mucho
- Es un hiperparametro
- Tradicionalmente 0.1 a 0.01
- Podría ser variable



Fig.3a Univariate gradient descent

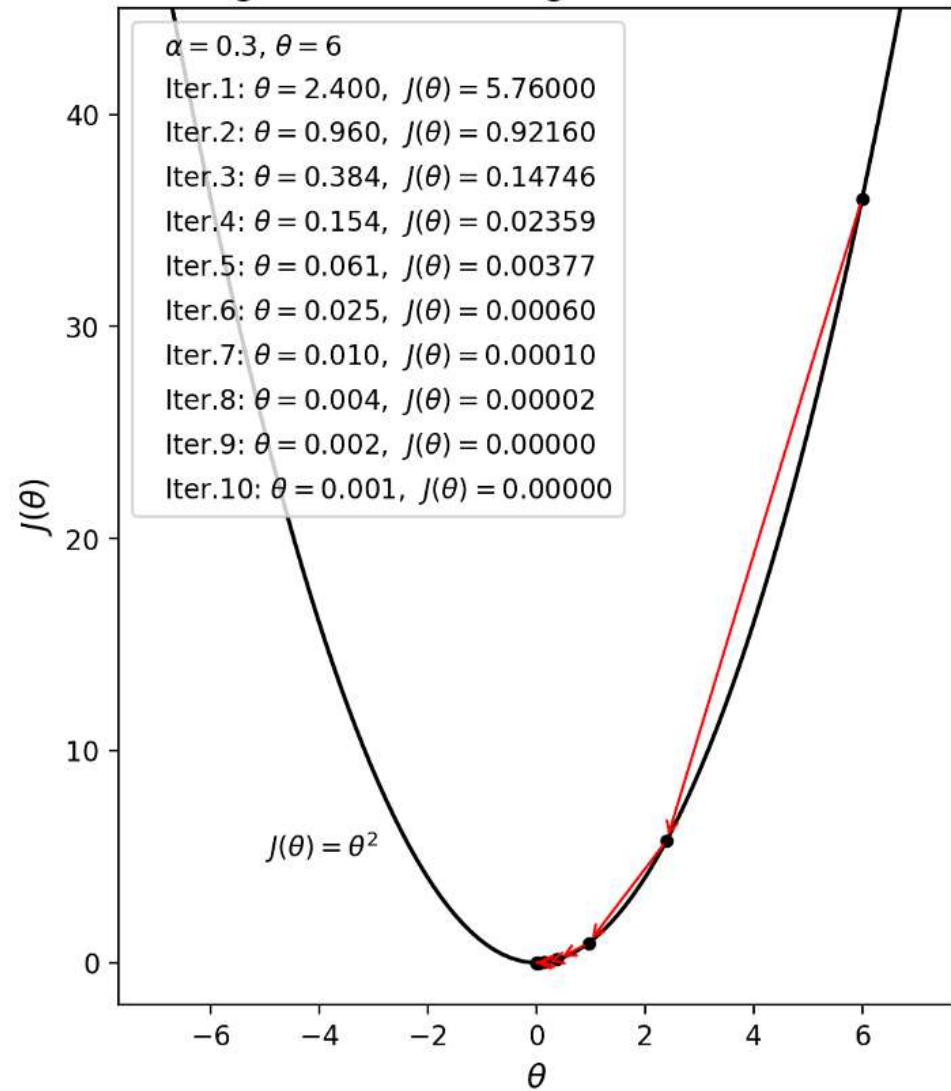


Fig.3b Monitoring convergence

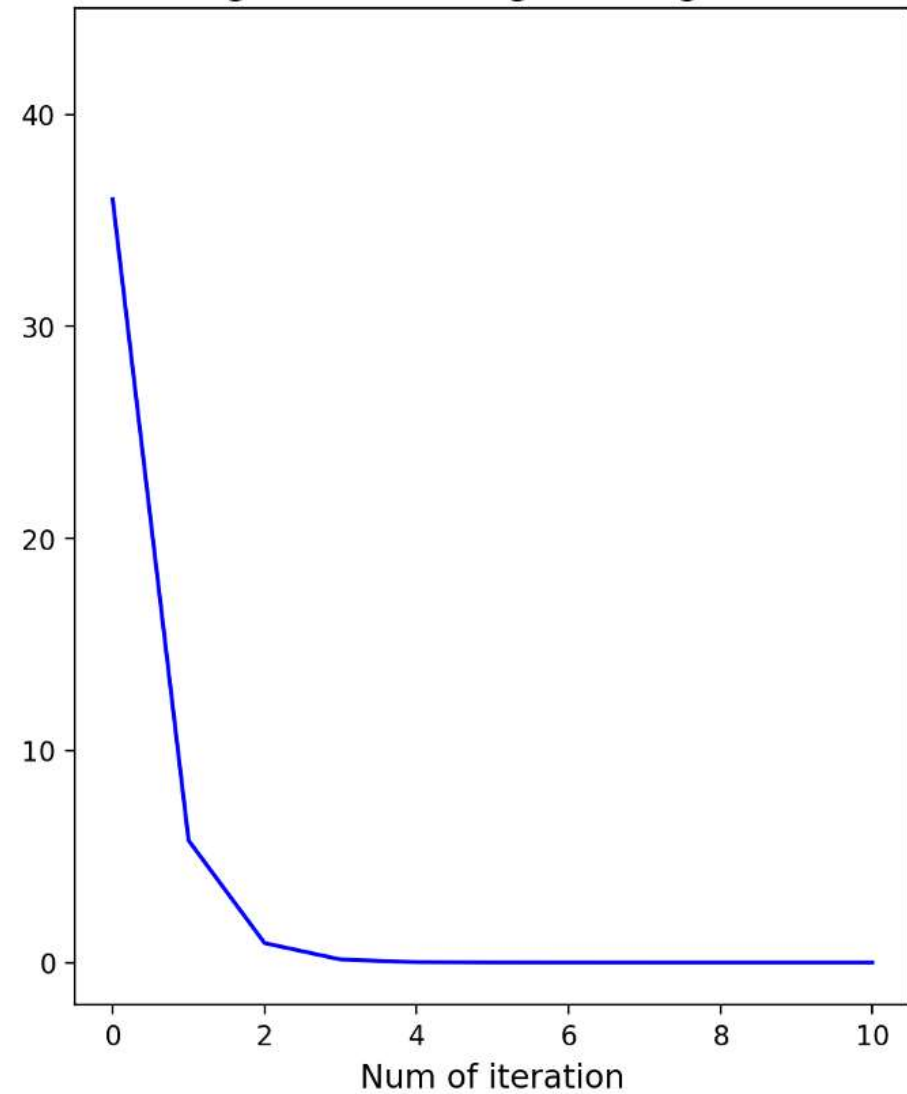


Fig.4a Learning rate - too large

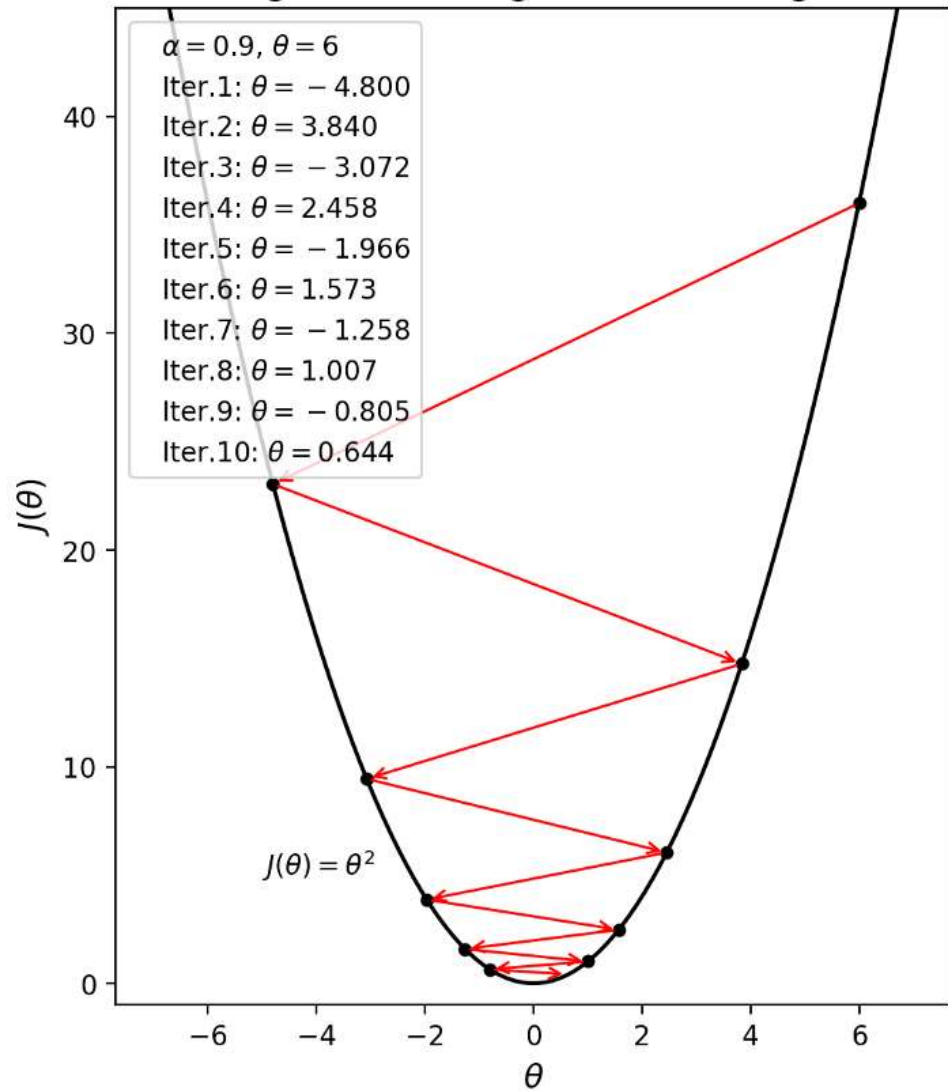


Fig.4b Monitoring convergence

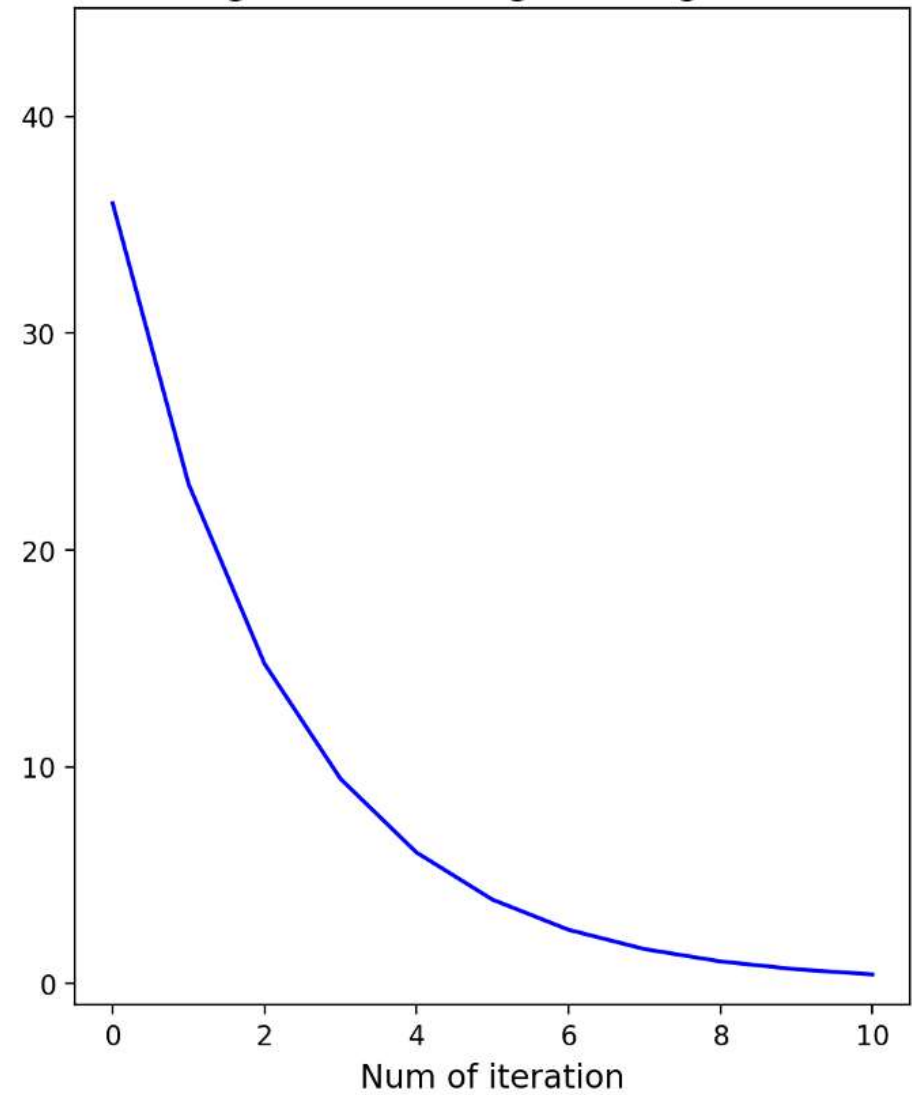


Fig.5a Learning rate - too small

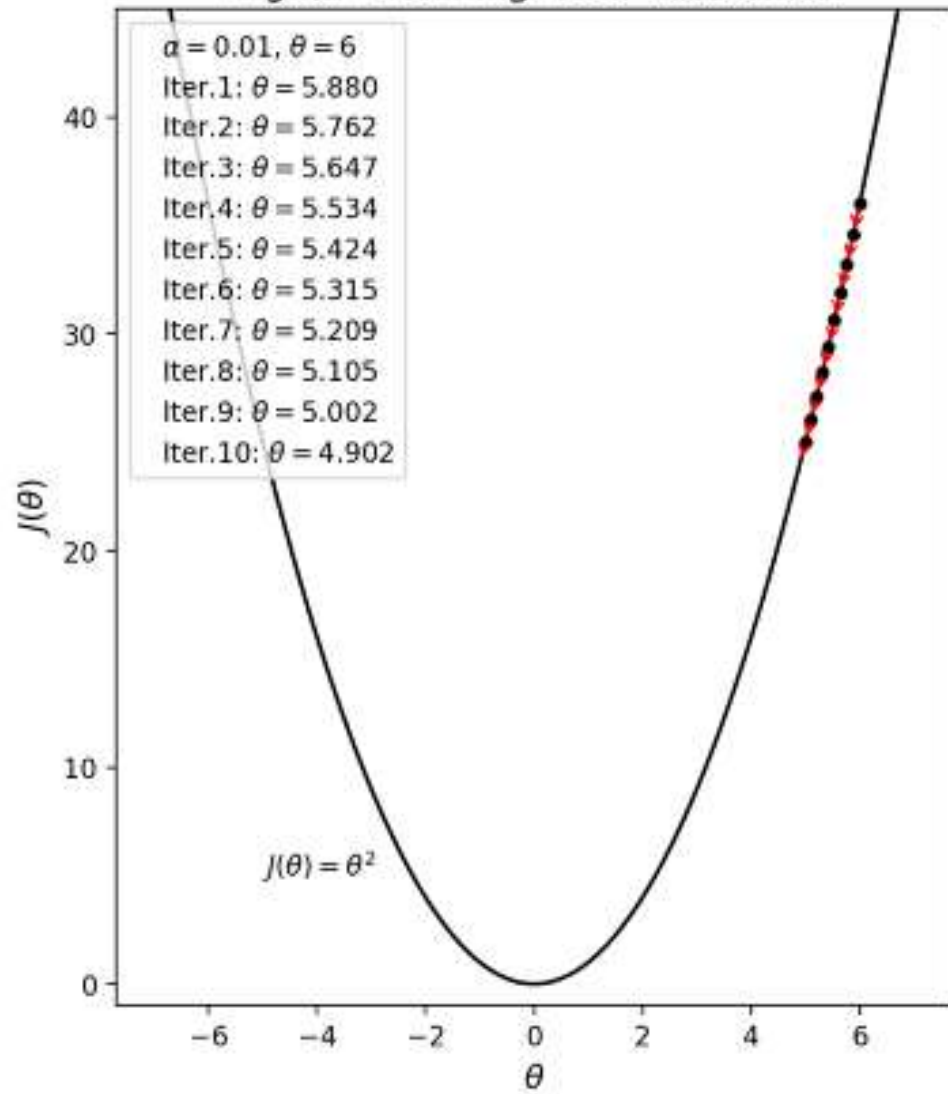
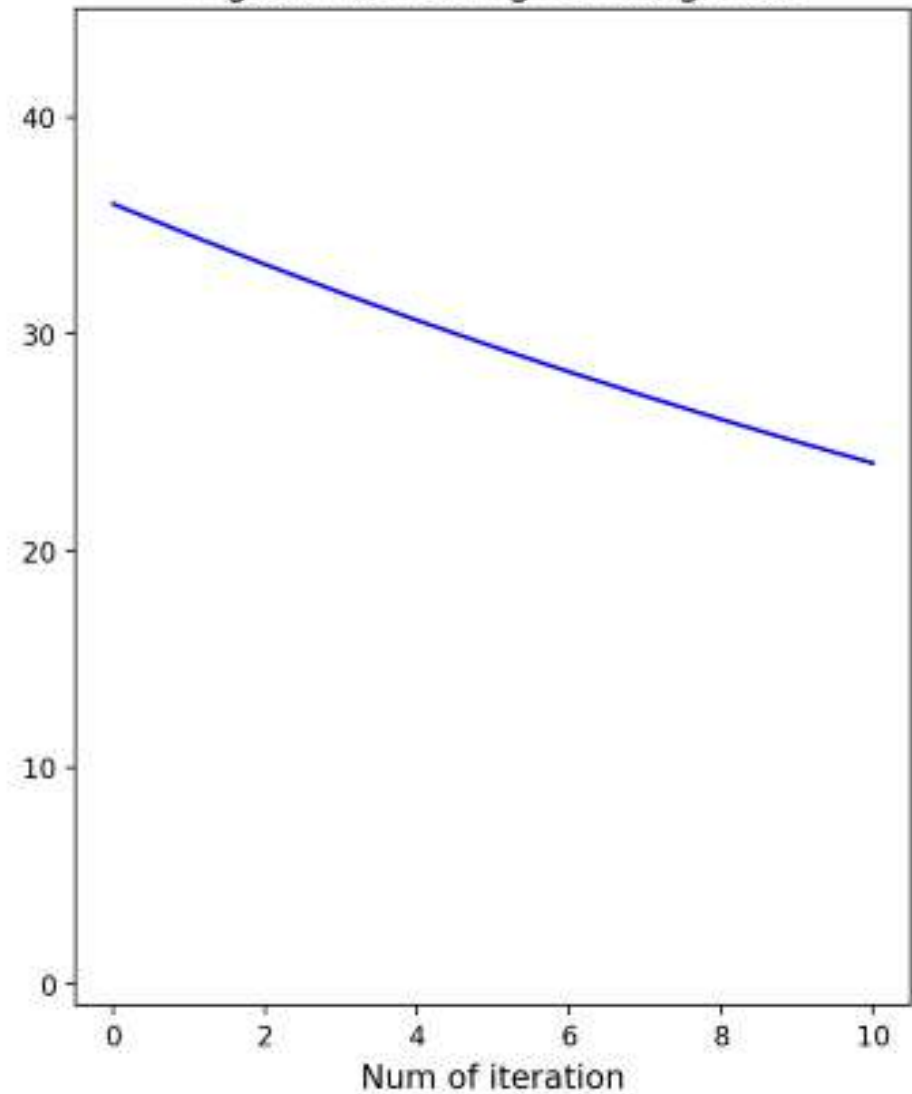


Fig.5b Monitoring convergence



Agreguemos
un poco de
complejidad



Dos variables independientes

$$J(\theta_0, \theta_1) = \theta_0^2 + \theta_1^2$$

Fig.1a 3D plot for $J(\theta_1, \theta_2) = (\theta_1^2 + \theta_2^2)$

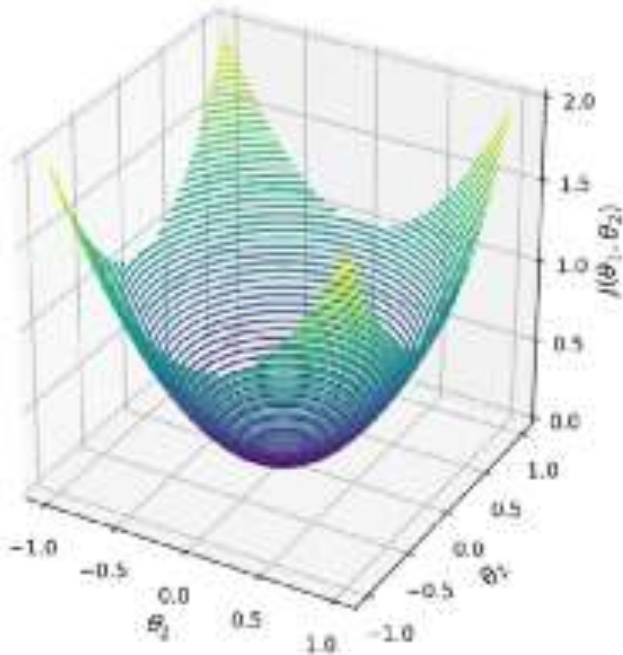


Fig.1b 3D plot - rotated

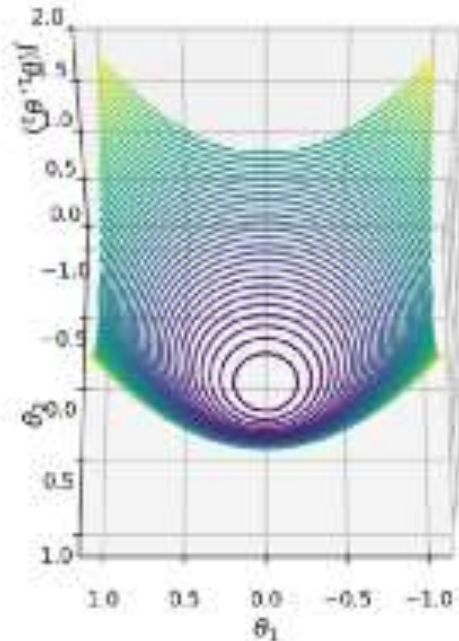
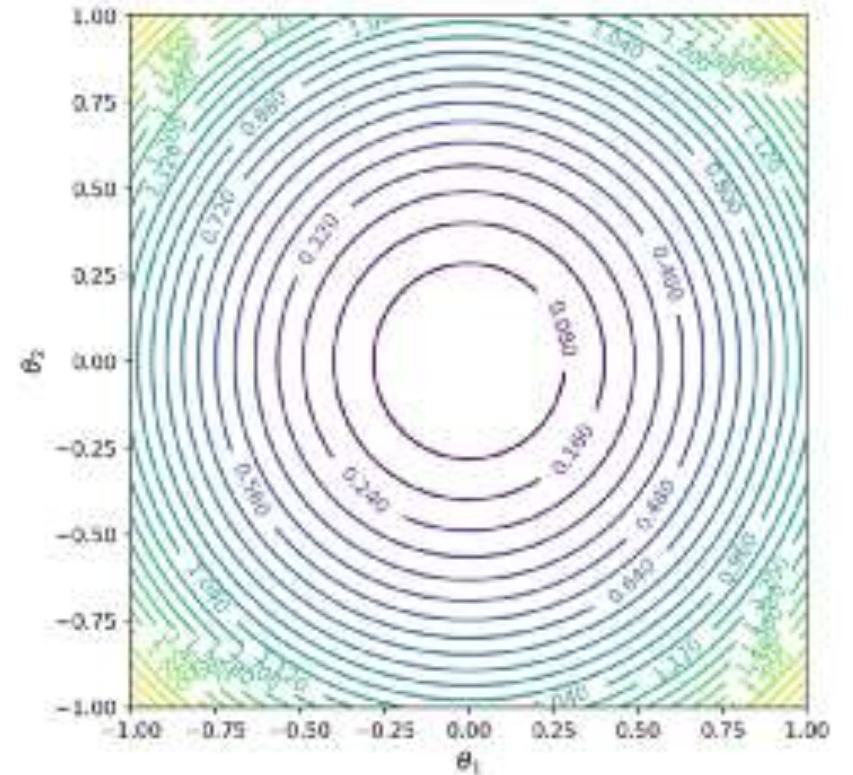


Fig.1c Contour plot for $J(\theta_1, \theta_2) = (\theta_1^2 + \theta_2^2)$



Update Rule

- Tenemos que derivar a cada una de las dos variables

$$\theta_0 = \theta_0 + \alpha \frac{\partial J}{\partial \theta_0}$$

$$\theta_1 = \theta_1 + \alpha \frac{\partial J}{\partial \theta_1}$$

Razón: No podemos calcularlo en simultáneo, por eso se hace uno por uno

Derivando

$$J(\theta_0, \theta_1) = \theta_0^2 + \theta_1^2$$

$$\begin{aligned}\frac{\partial J}{\partial \theta_0} &= 2\theta_0 \\ \theta_0 &= \theta_0 + \alpha \frac{\partial J}{\partial \theta_0}\end{aligned}$$

$$\theta_0 = (1 - 2\alpha)\theta_0$$

$$\begin{aligned}\frac{\partial J}{\partial \theta_1} &= 2\theta_1 \\ \theta_1 &= \theta_1 + \alpha \frac{\partial J}{\partial \theta_1}\end{aligned}$$

$$\theta_1 = (1 - 2\alpha)\theta_1$$

Fig.3a Multivariate gradient descent

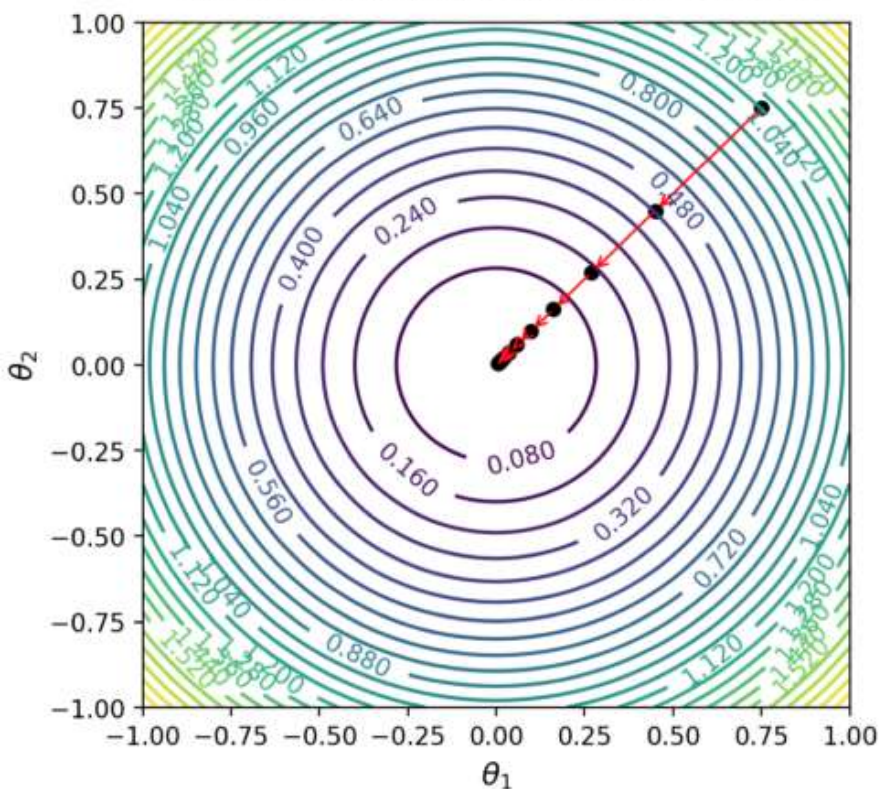
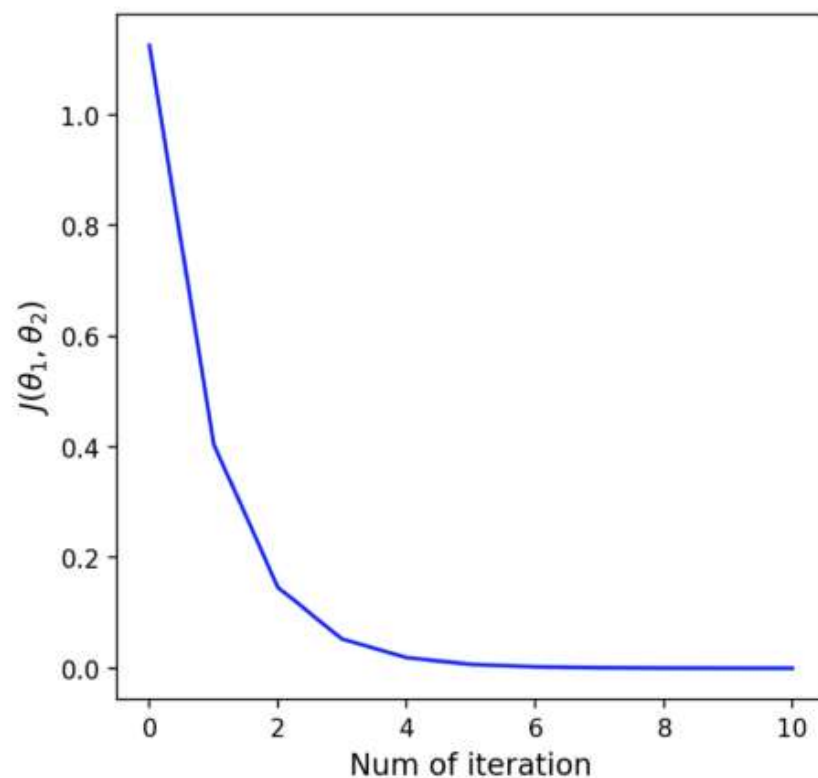


Fig.3b Monitoring convergence



Iter.1: $[\theta_1, \theta_2] = [0.450, 0.450]$, $J(\theta) = 0.40500$
 Iter.2: $[\theta_1, \theta_2] = [0.270, 0.270]$, $J(\theta) = 0.14580$
 Iter.3: $[\theta_1, \theta_2] = [0.162, 0.162]$, $J(\theta) = 0.05249$
 Iter.4: $[\theta_1, \theta_2] = [0.097, 0.097]$, $J(\theta) = 0.01890$
 Iter.5: $[\theta_1, \theta_2] = [0.058, 0.058]$, $J(\theta) = 0.00680$
 Iter.6: $[\theta_1, \theta_2] = [0.035, 0.035]$, $J(\theta) = 0.00245$
 Iter.7: $[\theta_1, \theta_2] = [0.021, 0.021]$, $J(\theta) = 0.00088$
 Iter.8: $[\theta_1, \theta_2] = [0.013, 0.013]$, $J(\theta) = 0.00032$
 Iter.9: $[\theta_1, \theta_2] = [0.008, 0.008]$, $J(\theta) = 0.00011$
 Iter.10: $[\theta_1, \theta_2] = [0.005, 0.005]$, $J(\theta) = 0.00004$

Aún más
complejidad!!





Loss function is back!

$$\mathcal{L}_{LSE}(W) = \sum_{i=0}^n (y_i - W^T X_i)^2$$

Diagram illustrating the components of the Least Squares Error (LSE) loss function:

- $(w_0, w_1 \dots w_d)$ (Weights)
- $\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}$ (Target values)
- $\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$ (Weight vector)
- $\begin{bmatrix} x_{0,0} & x_{0,1} & \dots & x_{0,d} \\ x_{1,0} & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,0} & x_{n,1} & \dots & x_{n,d} \end{bmatrix}$ (Feature matrix)

La Derivada

$$\frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^n x_j^{(i)} (w_j x_j^{(i)} - y^{(i)})$$

Lo igualamos a 0, y se pone feo rápido...

Si quieren> OLS, Ecuación Normal (one-step learning).

O hacemos gradient descent...

El Gradient Descent

$$w_j \leftarrow w_j - \alpha \sum_{i=1}^n x_j^{(i)} (w_j x_j^{(i)} - y^{(i)})$$

Lo podemos poner más bonito y sin tantos índices.

Vectorizar el problema

- X es una matriz de $N \times D$
 - N es el número de observaciones / ítems / data points
 - D es el número de features / características / dimensiones
- W es un vector columna de D pesos
- Y es un vector columna de N elementos target
- \hat{Y} es el vector columna predicho de X con W

$$\hat{Y} = XW$$

Vectorizando

La función de pérdida

$$\mathcal{L} = (Y - XW)^2$$

El Gradient Descent update rule

$$W \leftarrow W - \alpha X^T (XW - Y)$$

Y por si lo quieren resolver en un solo paso complicado:

Ecuación Normal

$$W = (X^T X)^{-1} X^T \hat{Y}$$

Recordando Gradient Descent

- 1) Calcular la gradiente (derivada) de la función de pérdida en el punto actual
- 2) Mover en la dirección opuesta (update rule)
- 3) Repetir hasta converger, o el cansancio, o el final de los tiempos...

Gradient Descent

$$W \leftarrow W - \alpha X^T (XW - Y)$$

Ventajas

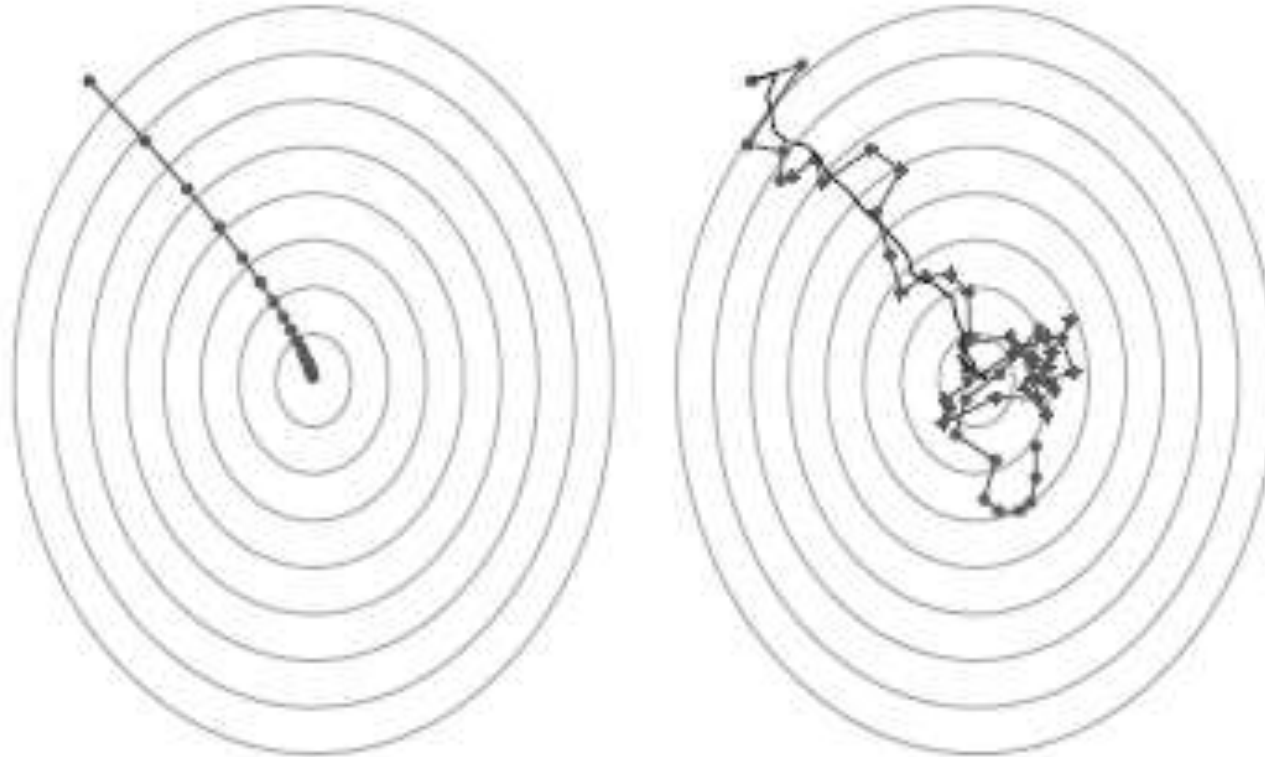
- Funciona para todas las funciones de optimización
- Eficiente en tiempos de computación

Desventaja

- No llega al valor exacto, pero es suficiente
- En cada iteración tiene que usar todos los data points.

Stochastic Gradient Descent (SDG)

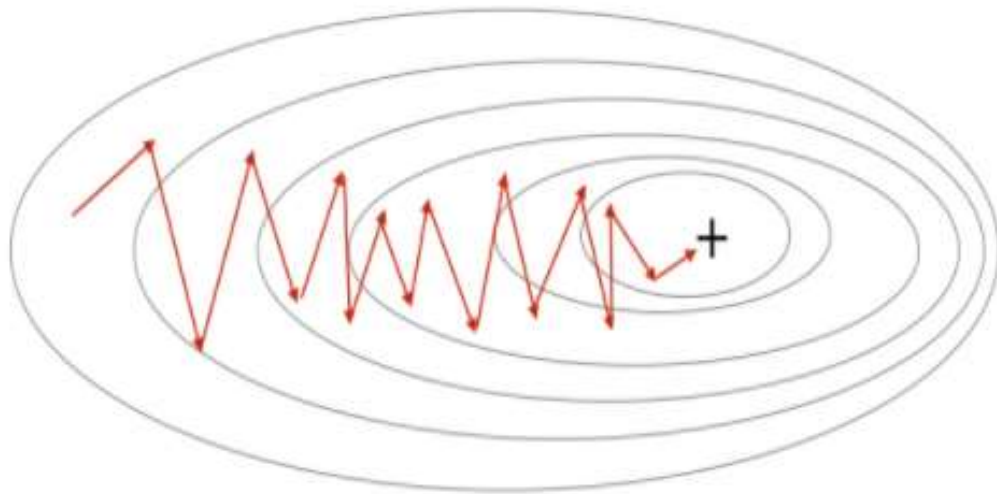
En lugar de tomar toda la matriz X , se toma aleatoriamente un solo ítem, se calcula la gradiente y se actualiza



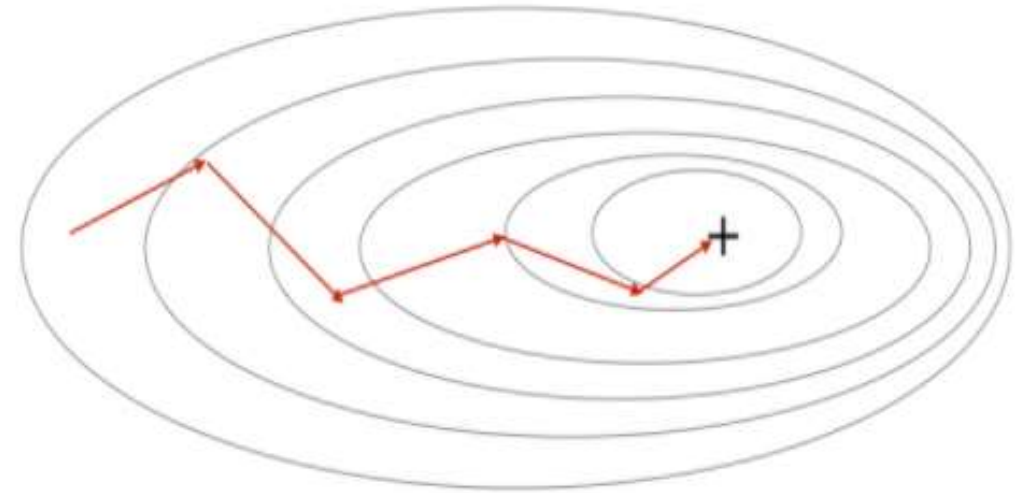
Mini-Batch Gradient Descent

En lugar de tomar toda la matriz X , se toma unos cuantos ítems en cada iteración, se calcula la gradiente y se actualiza

Stochastic Gradient Descent

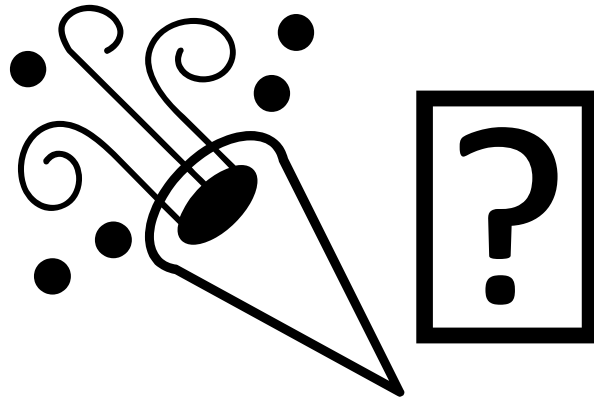


Mini-Batch Gradient Descent



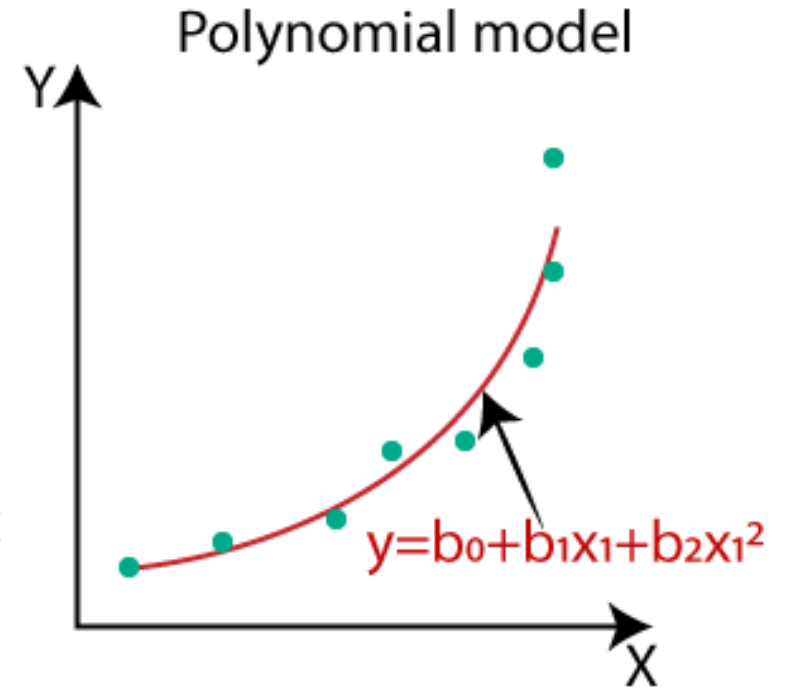
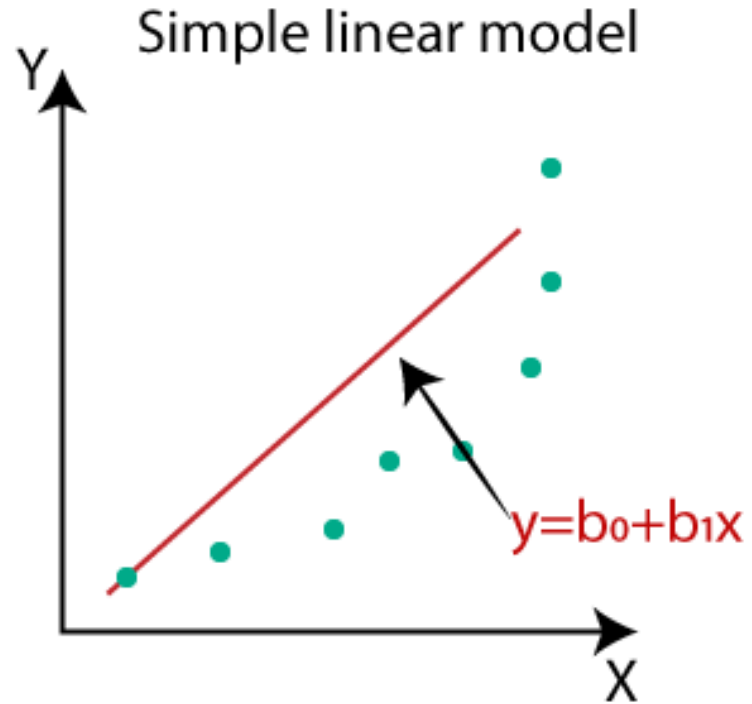
Finalizando

Cuando la gradiente converge en un valor cercano a 0, quiere decir que hemos llegado a un mínimo en la función pérdida, y que nuestros pesos se ajustan bien a la data.

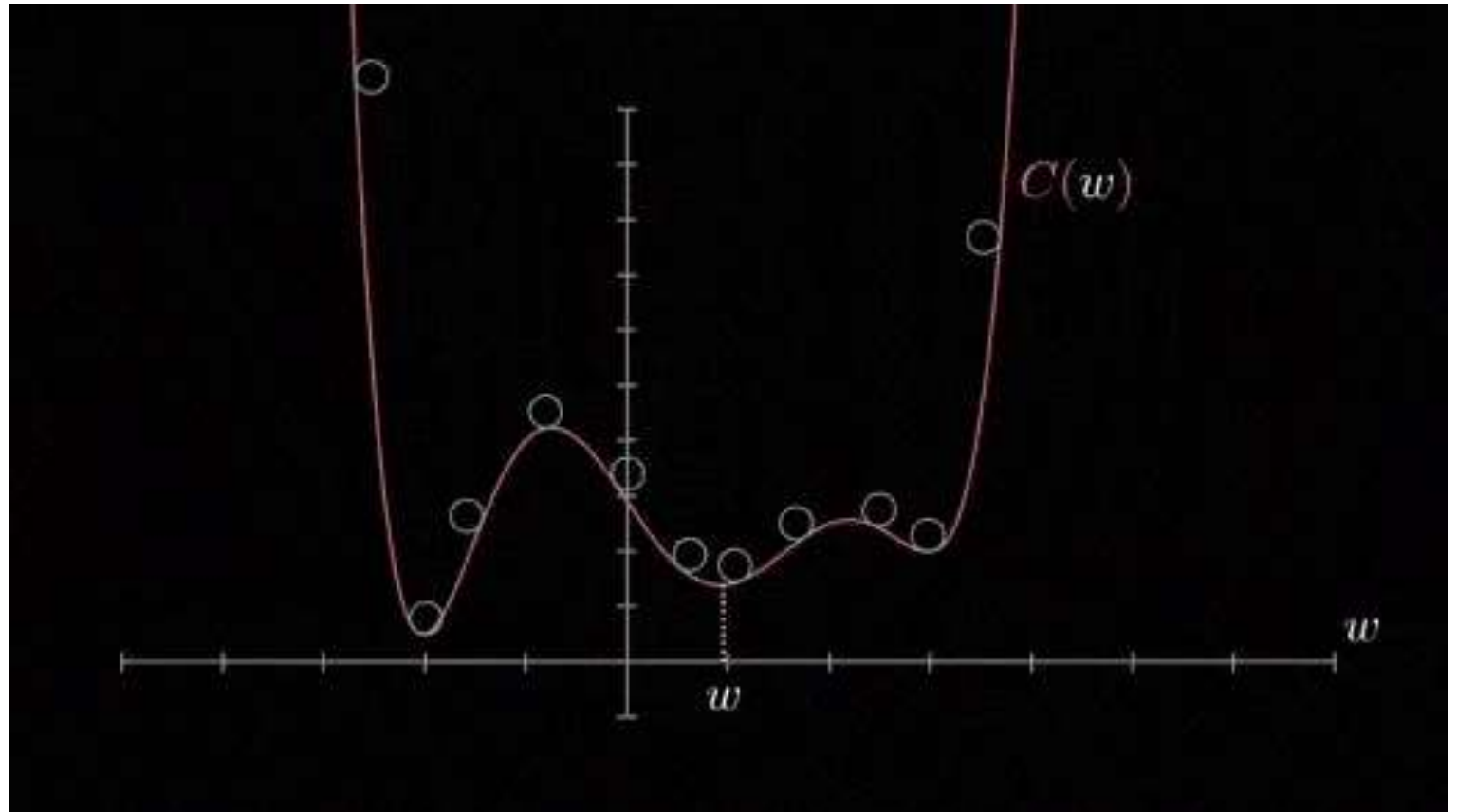




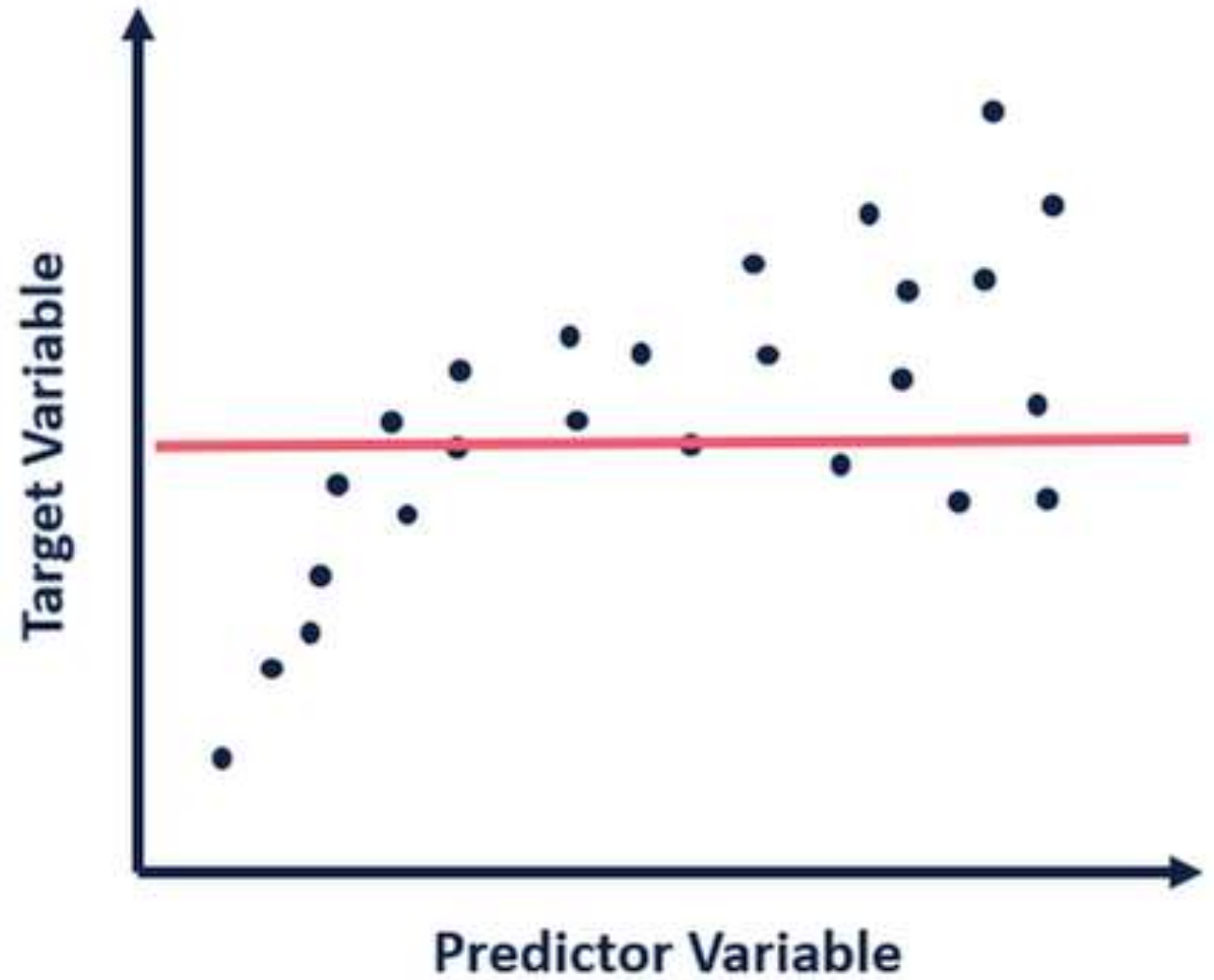
Regresión Polinomial



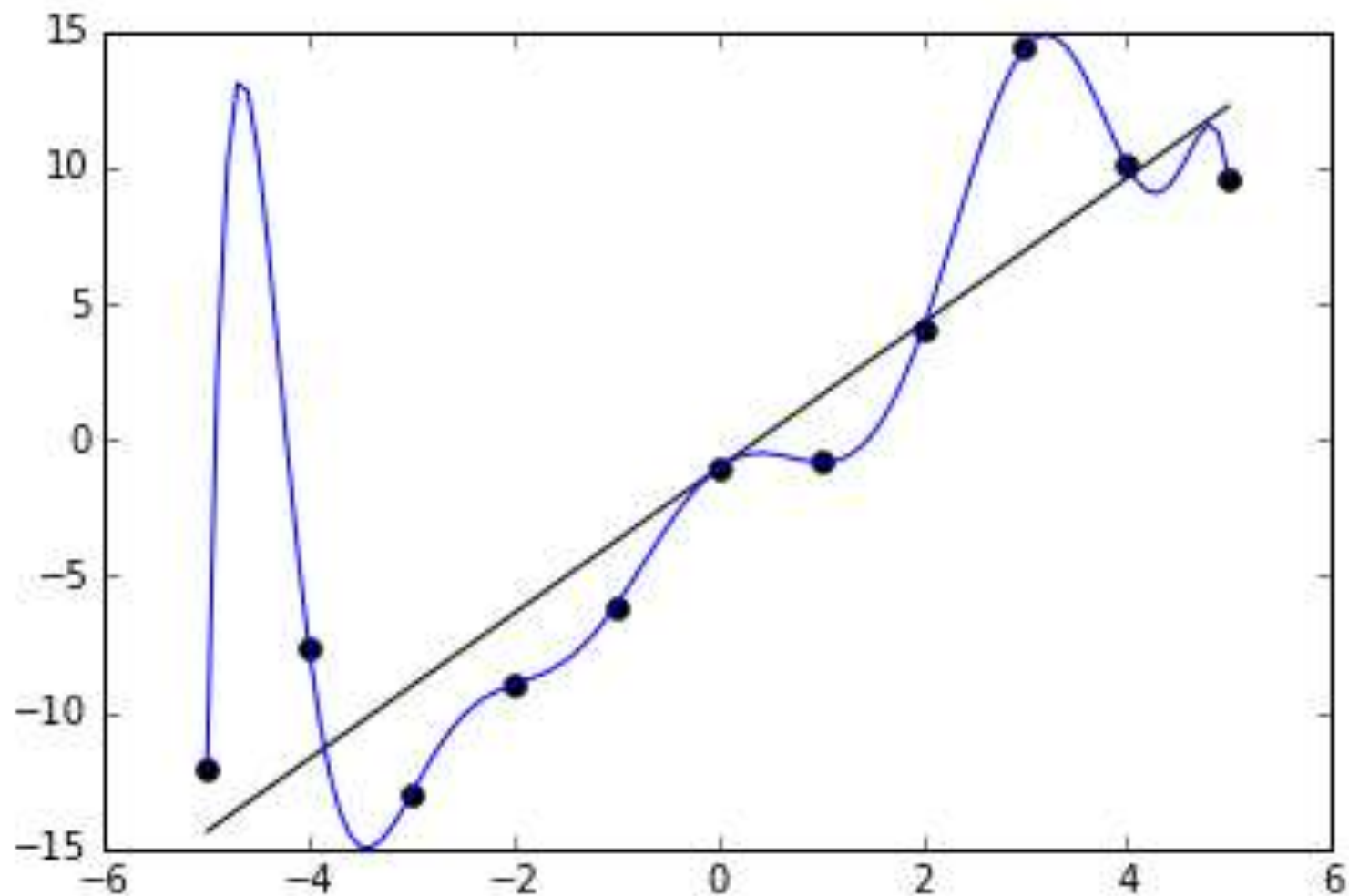
Mínimo Local y Mínimo Global



Underfitting



Overfitting



Regularización

