

Regresión lineal en línea

Raúl Martínez Noriega
Universidad Panamericana
Verano 2020

Aprendizaje supervisado

Los datos de entrenamiento consisten en pares de objetos: datos de entrada y su resultado deseado. El objetivo es encontrar una función que relacione la entrada con su resultado deseado, de tal forma que podamos hacer predicciones sobre el resultado de muestras futuras.

Básicamente, los problemas de aprendizaje supervisado están clasificados en:

- Regresión: tratamos de predecir resultados de una variable continua.
 - Regresión lineal
 - Redes neuronales
- Clasificación: predecir resultados de una variable discreta.
 - Regresión logística
 - Redes neuronales

Regresión lineal

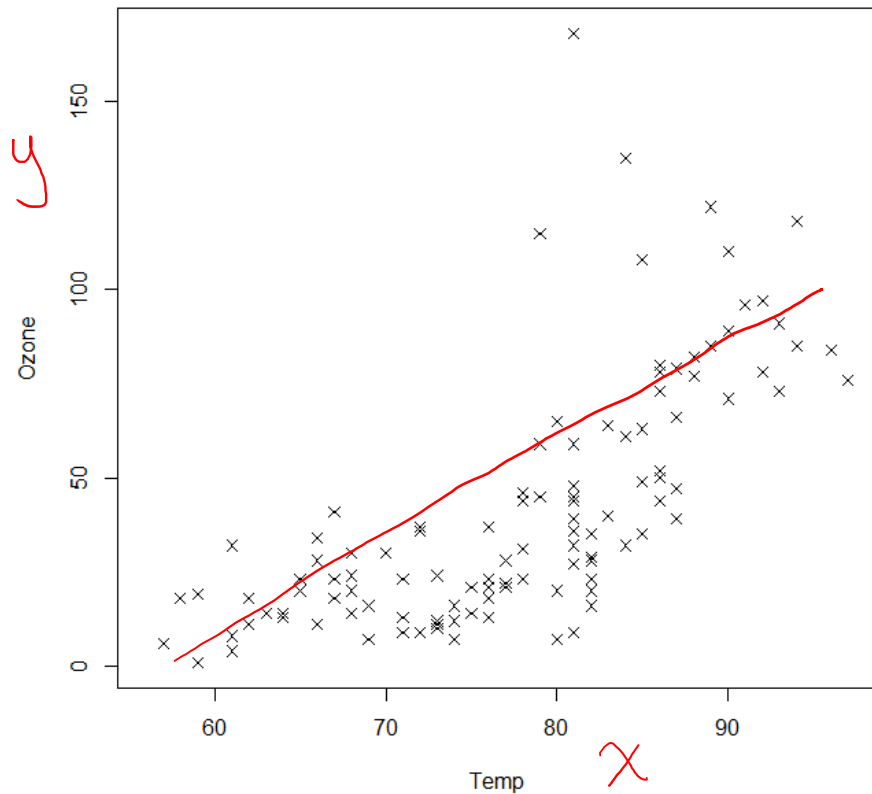
Ecuación normal

Gradient descent

Stochastic Gradient Descent (SGD)

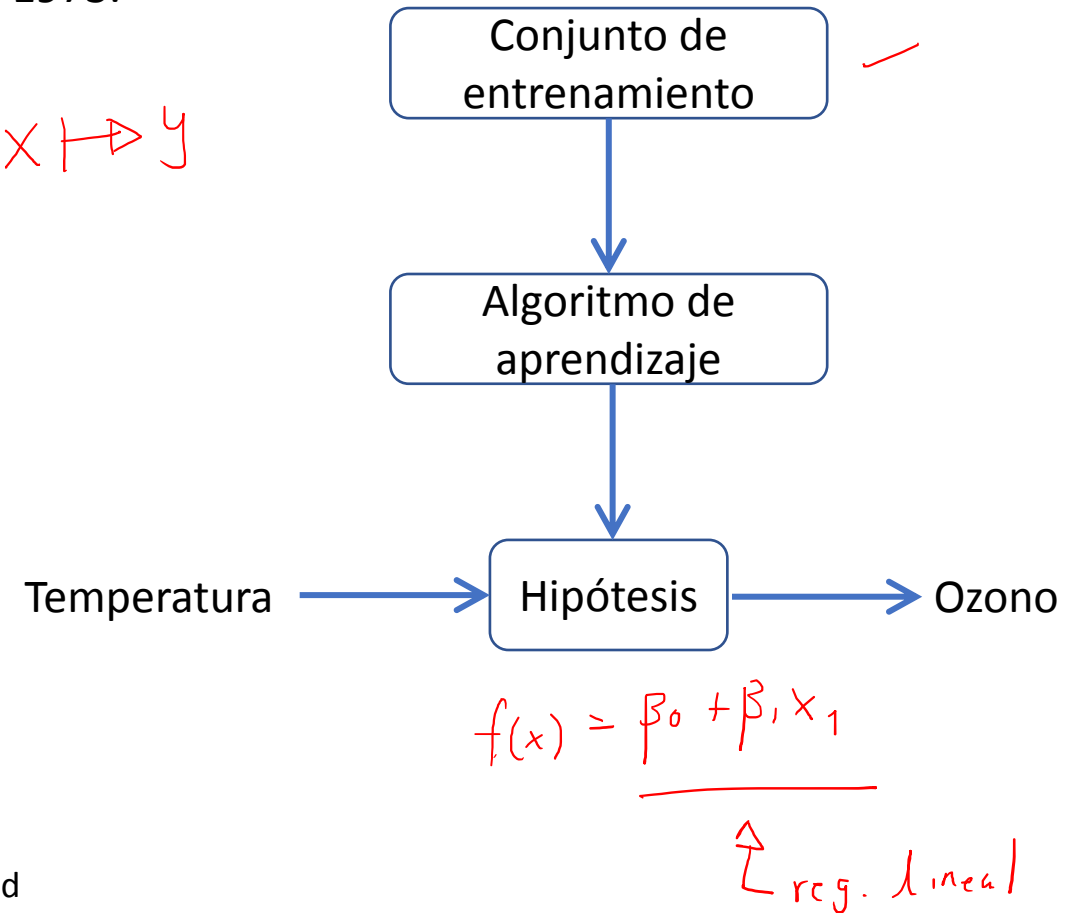
Regresión lineal

Daily air quality measurements in New York, May to September 1973.



$$f(x) := x \mapsto y$$

Ozone: Mean ozone in parts per billion from 1300 to 1500 hours at Roosevelt Island
Temp: Maximum daily temperature in degrees Fahrenheit at La Guardia Airport.



Función costo

Hipótesis:

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$$

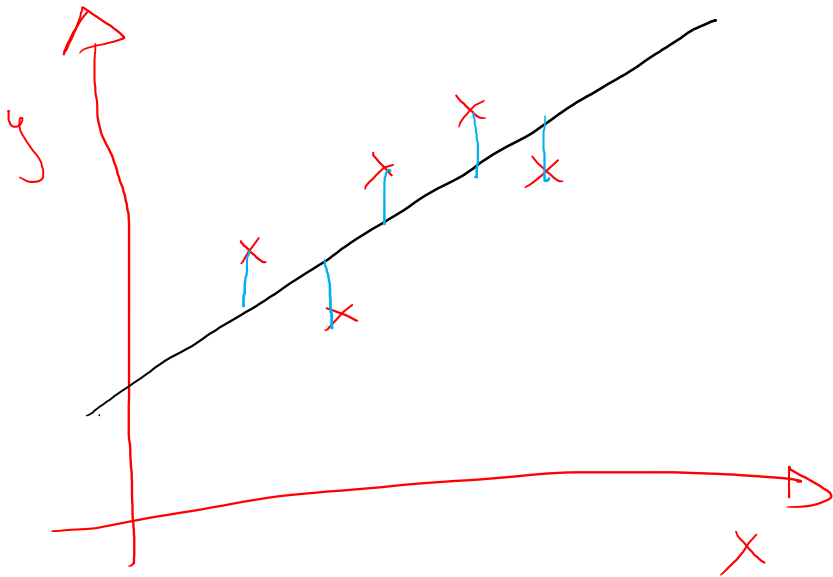
$$\hat{y} = \beta_0 + \beta_1 x$$

Función costo:

$$\begin{aligned} J(\beta_0, \beta_1) &= \sum_{i=1}^n (f(x_i) - y_i)^2 \quad \checkmark \\ &= \sum_{i=1}^n \left(\underbrace{\beta_0 + \sum_{j=1}^p X_{i,j} \beta_j}_{\hat{y}_i} - y_i \right)^2 \end{aligned}$$

minimizar $J(\beta_0, \beta_1)$

encontrar β_0 y β_1



¿Cuál es la mejor línea
que ajusta los datos?

Hipótesis: $f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$

Función costo:

$$J(\beta_0, \beta_1) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2$$

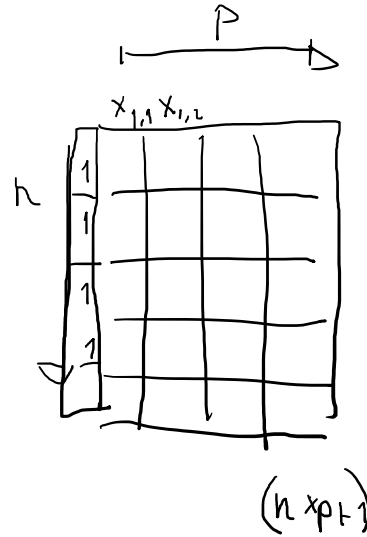
$$J(\beta_0, \beta_1) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\frac{\partial J(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

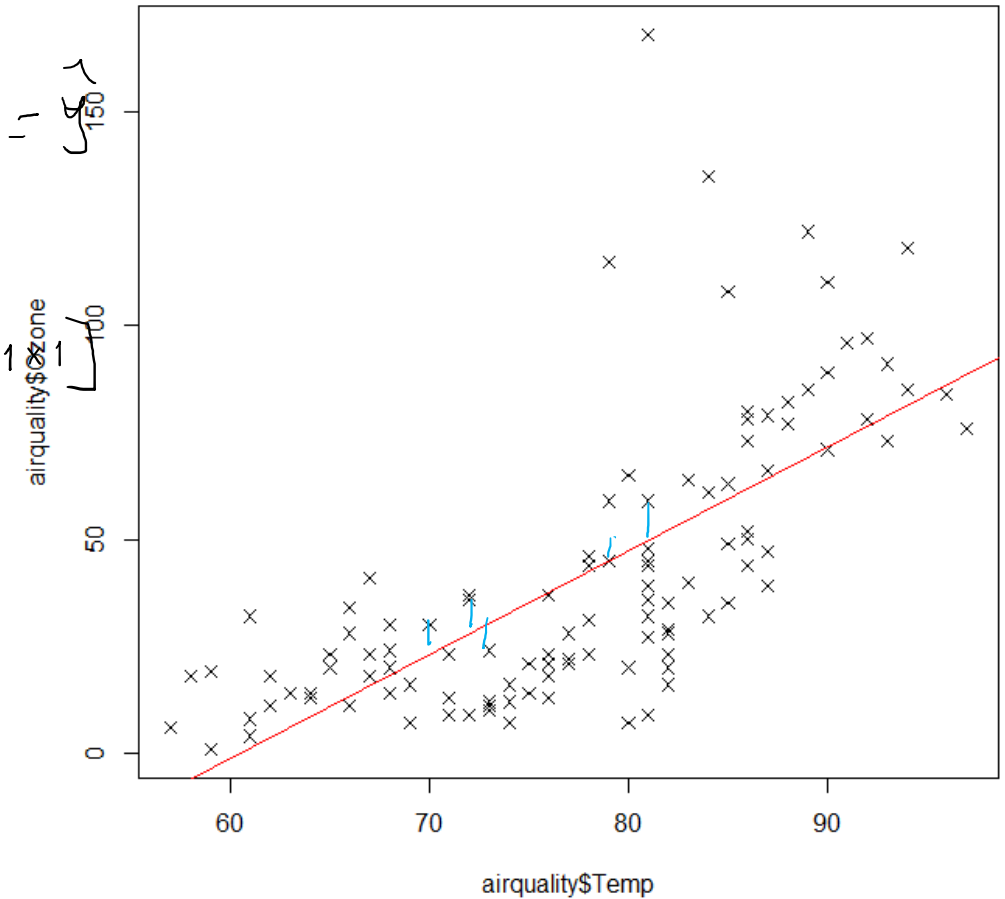
$\mathcal{O}(n^3)$



$$\mathbf{X}\boldsymbol{\beta} = \hat{\mathbf{y}}$$

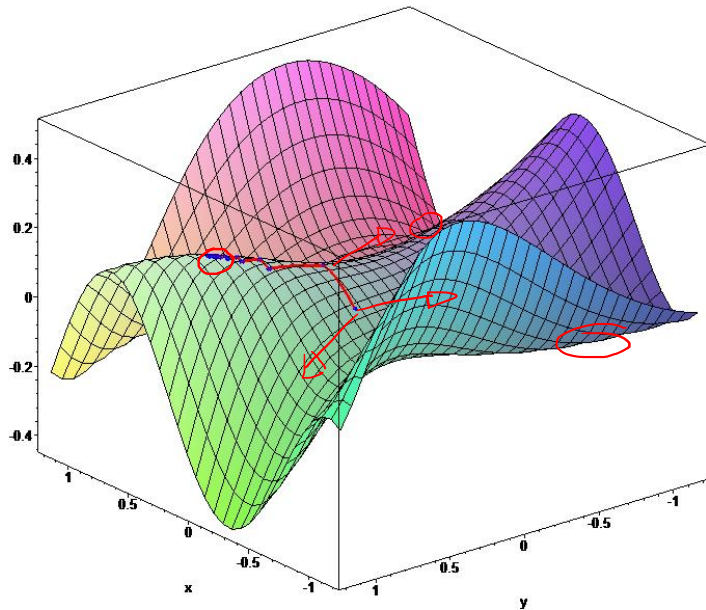
Ecuación normal

closed form solution



Gradient descent - I

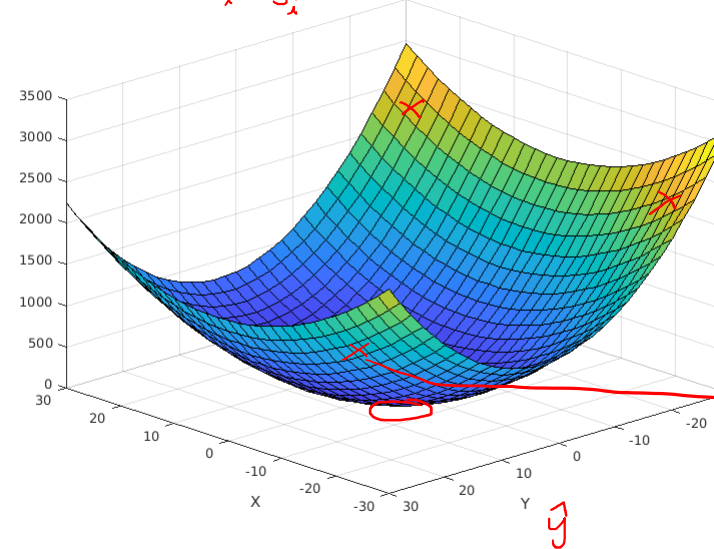
Es un algoritmo de optimización iterativo que nos ayuda a encontrar un mínimo local de una función diferenciable.



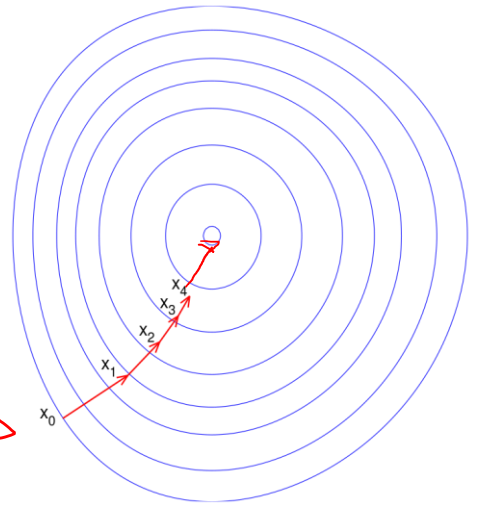
$$\text{error} = (y_i - \hat{y}_i)$$

$$\hat{y}_i = y_i$$

$$f(x,y) = (x-2)^2 + 2(y-3)^2$$



$$\beta_0$$



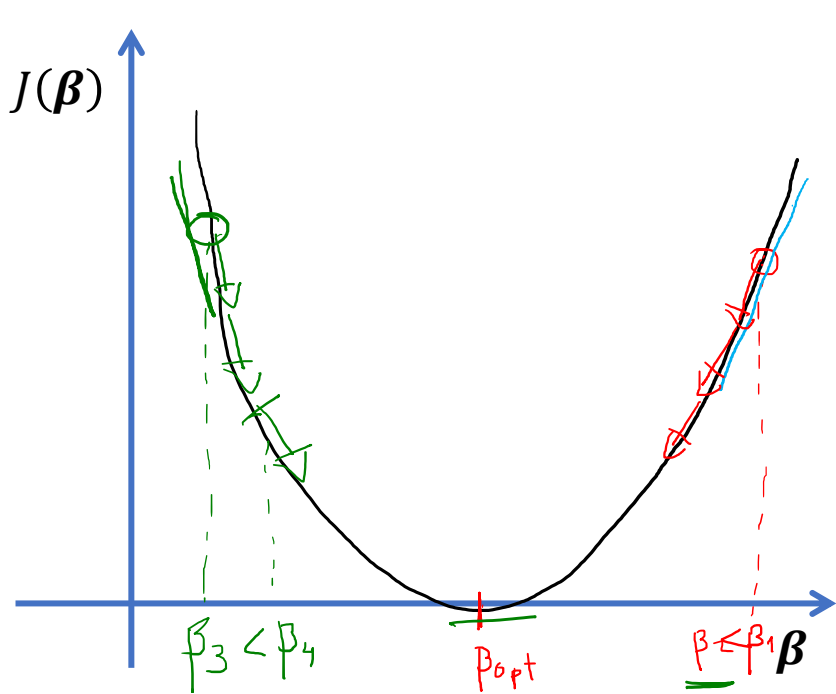
$$J(\beta_0, \beta_1) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \overbrace{X_j \beta_j}^g \right)^2$$

Función costo:

$$J(\beta_0, \beta_1) = \frac{1}{2n} \sum_{i=1}^n \left(\beta_0 + \sum_{j=1}^p X_{i,j} \beta_j - y_i \right)^2$$

Gradient descent - II

Por simplicidad, asumamos una regresión lineal univariable:



Gradient descent:

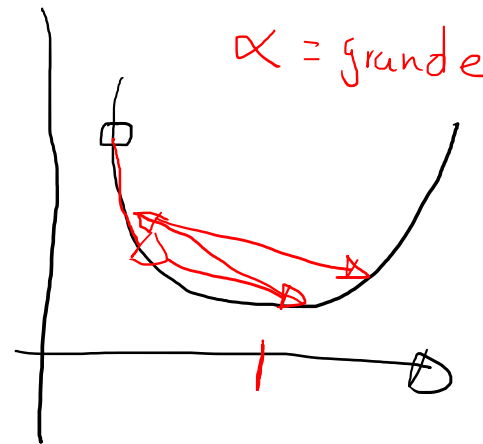
Repetir hasta converger{

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\beta)$$

↑ paso

$$f(x) = x\beta$$

$$\beta_2 = \beta_1 - \alpha \frac{\partial J(\beta)}{\partial \beta}$$



$$\begin{aligned} \frac{\partial}{\partial \beta_j} J(\beta) &= \frac{\partial}{\partial \beta_j} \frac{1}{2n} \sum_{i=1}^n (\beta_0 + x_i \beta_1 - y_i)^2 \\ &= \frac{2}{2n} \sum_{i=1}^n (\beta_0 + x_i \beta_1 - y_i) \frac{\partial}{\partial \beta_j} (\beta_0 + x_i \beta_1 - y_i) \end{aligned}$$

$\beta_0 = 1$
 $\beta_1 = x_i$

Por tanto las derivadas parciales son:

$$\frac{\partial}{\partial \beta_0} J(\beta) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)$$

$$\frac{\partial}{\partial \beta_1} J(\beta) = \frac{1}{n} \sum_{i=1}^n x_i * (\hat{y}_i - y_i)$$

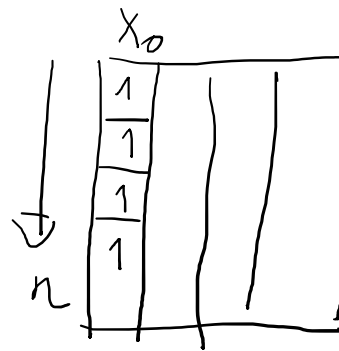
Gradient descent - III

Linear regression

Update rules:

$$\beta_0 = \beta_0 - \alpha \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)$$

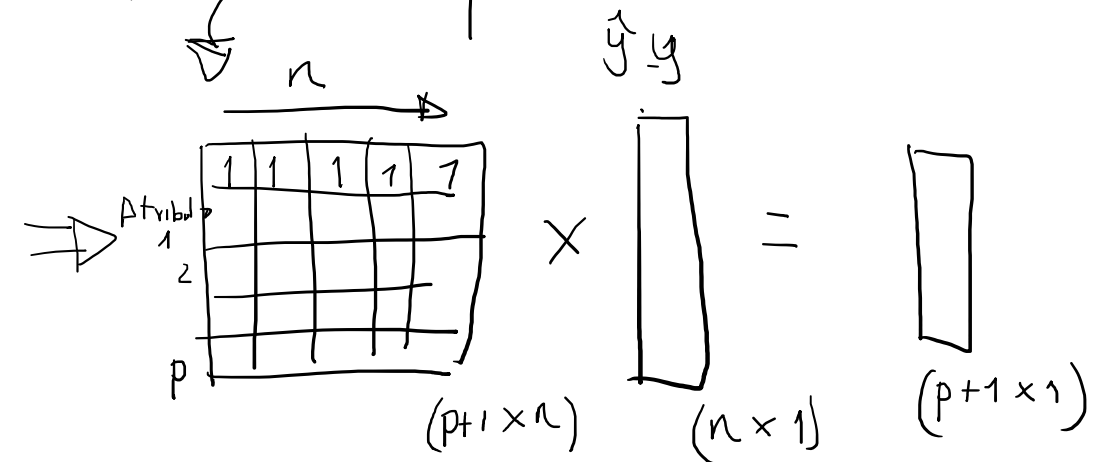
$$\beta_1 = \beta_1 - \alpha \frac{1}{n} \sum_{i=1}^n x_i (\hat{y}_i - y_i)$$



Multivariate linear regression

Update rule:

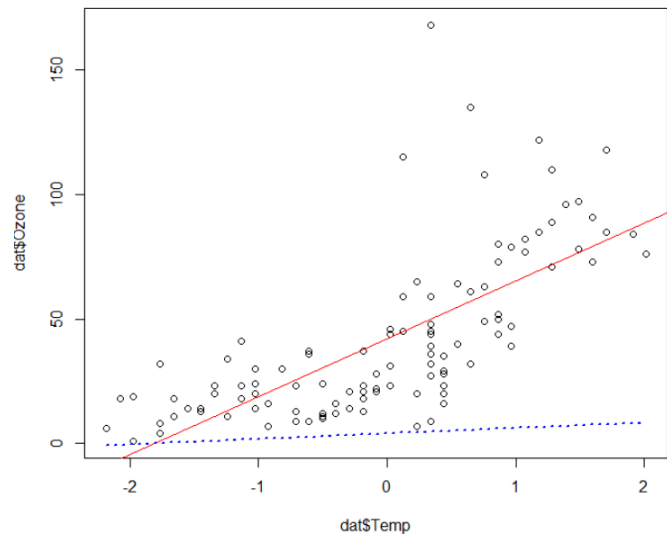
$$\beta = \beta - \alpha \frac{1}{n} X^T (\hat{y} - y)$$



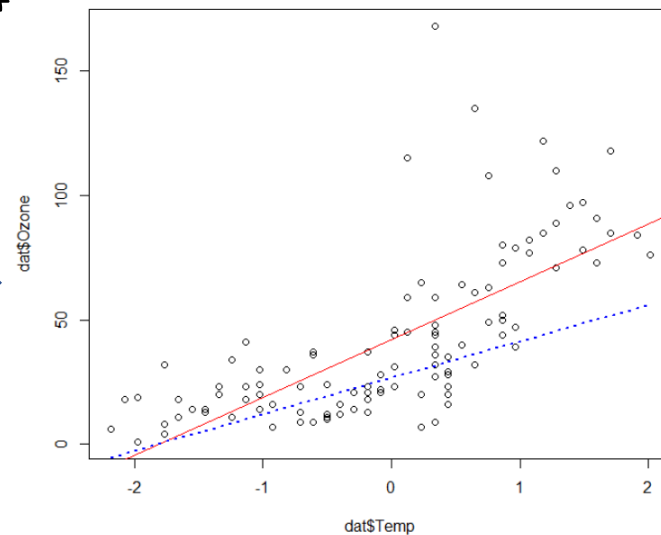
Ejemplo

todo el data set

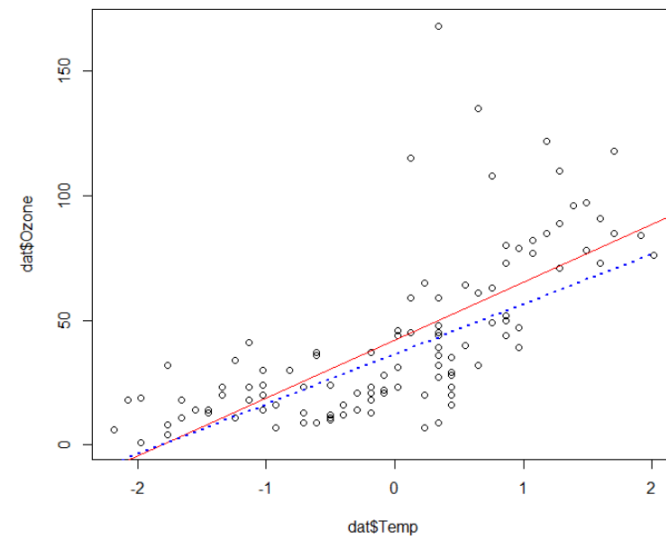
iteraciones: 10 = epoch



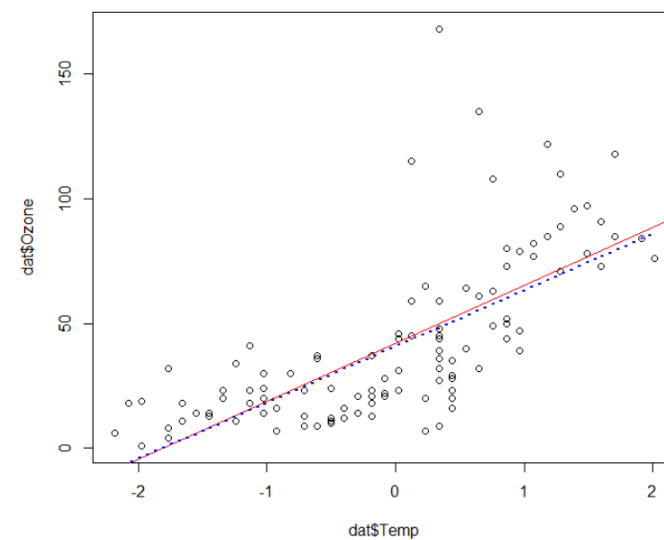
iteraciones: 100



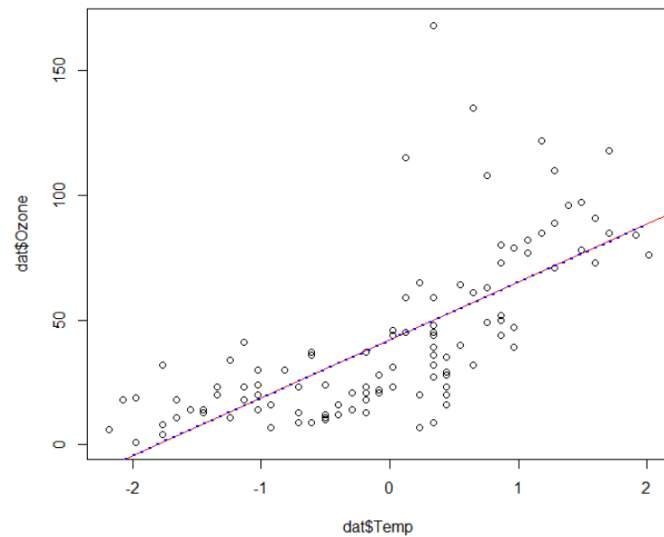
iteraciones: 200



iteraciones: 350

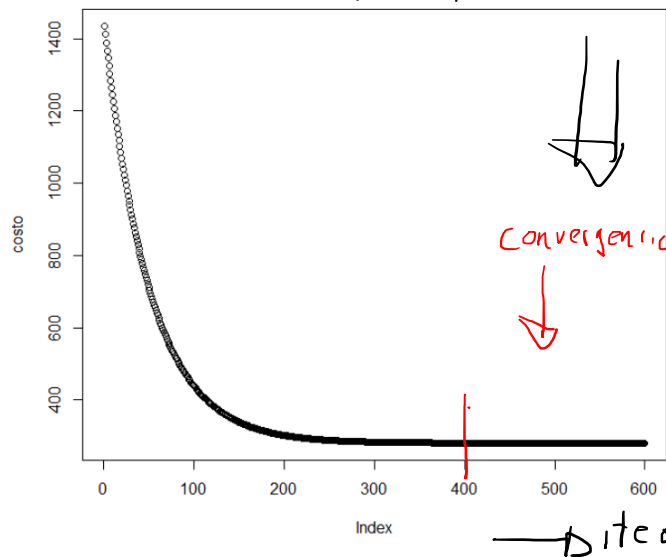


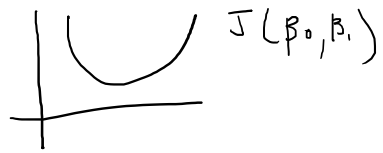
iteraciones: 600



$$J_i - J_{i+1} < 10^{-3}$$

$$J(\beta_1, \beta_0)$$





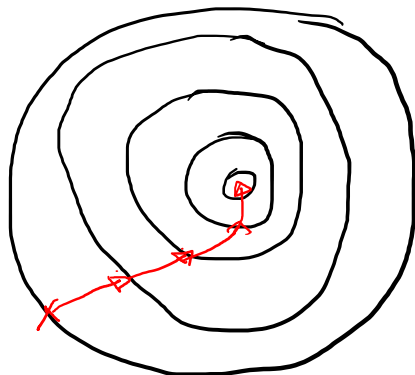
Stochastic Gradient descent ^(SGD)

Gradient descent:

Repetir hasta converger{

$$\checkmark \beta = \beta - \alpha \frac{1}{n} X^T (\hat{y} - y) \checkmark$$

}



Stochastic Gradient descent:

Repetir hasta converger{

$$\checkmark \beta = \beta - \alpha x_i (\hat{y}_i - y_i)$$

}

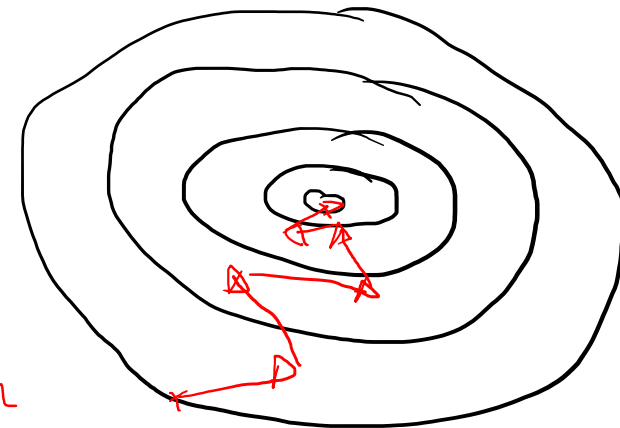
¹ muestra
↓
✓ rápido

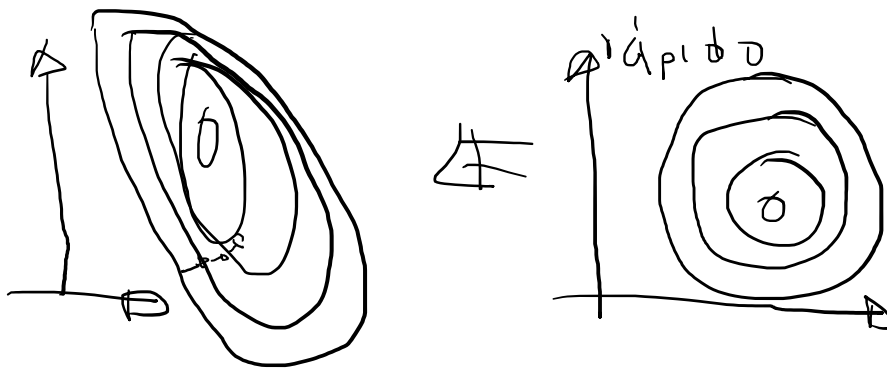
K muestras
Batch - SGD

$$\beta = \beta - \frac{\alpha}{K} X^T (\hat{y}_i - y_i)$$

X^* → batch

d,m ⇒ K × (p+1)





Consejos prácticos - I

- Revisar que las variables tengan una escala similar

- Usualmente entre $[-1, 1]$, e.g. normalizar

- Convergencia más rápida

$$\frac{x - \mu}{\sigma}$$

$[-1, 1]$ ✓

$[0, 3]$ ✓

$[-1, 2]$ ✓

atributos

X_1	$[-10,000, -10,000]$
X_2	$[0.5, 2]$
X_3	$[-3, 10]$

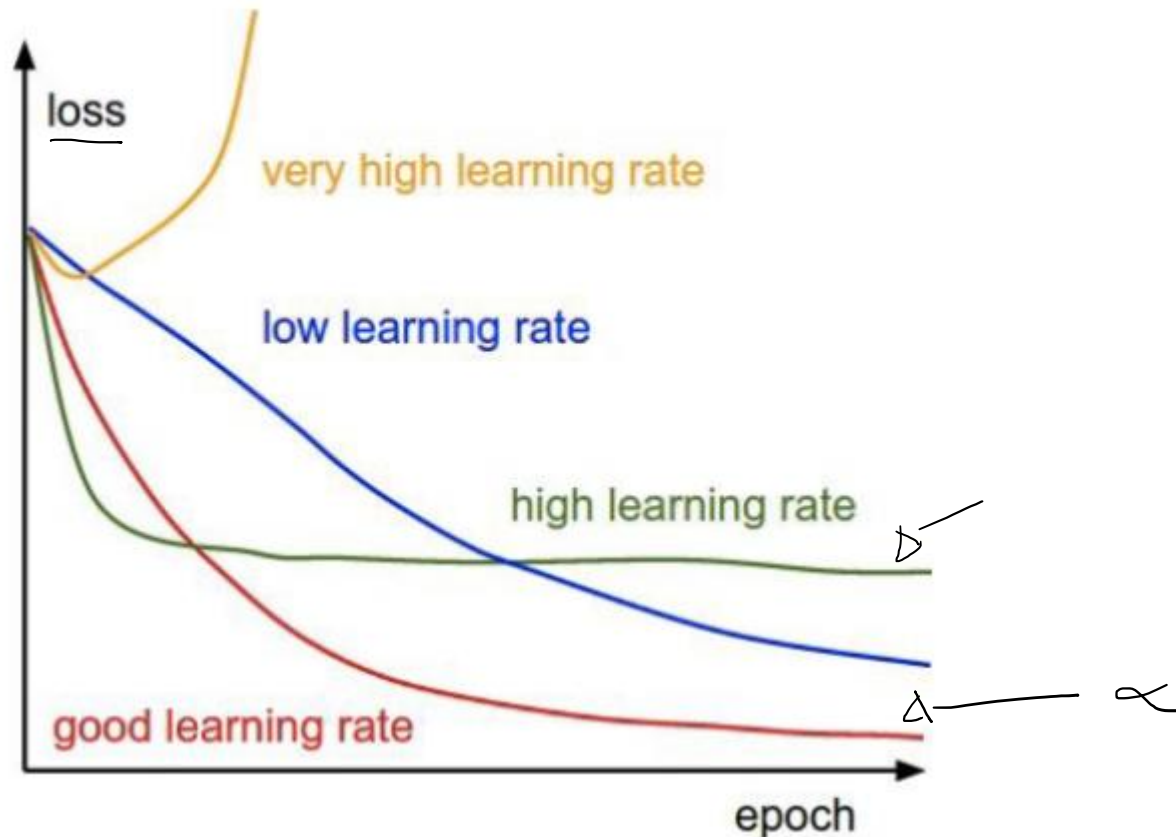
- ¿Qué valor de α escoger?

• .001, , 0.01, 0.03, 0.1, , 1 ...

(Handwritten arrows indicate a sequence from 0.01 to 1, and a checkmark is under 1.)

Consejos prácticos - II

- ¿Cómo asegurarnos que gradient descent esta trabajando correctamente?



Podemos detener el algoritmo automáticamente:

Decidimos que el algoritmo ha convergido si de una iteración a otra $J(\boldsymbol{\beta})$ tiene un cambio menos a 10^{-3} .