

# Regresión lineal multivariada

⑦

Datos con  $n$  variables  $m$  records

dimensiones

Dataframe

	$y$	$x_1$	$x_2$	...	$x_n$
1	$y_1$	$x_1^1$	$x_2^1$	...	$x_n^1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$m$	$y_m$	$x_1^m$	$x_2^m$	...	$x_n^m$

	$y$	$x_1$	$x_2$	...	$x_d$
1	$y^{(1)}$	$x_1^{(1)}$	$x_2^{(1)}$	...	$x_d^{(1)}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$y^{(n)}$	$x_1^{(n)}$	$x_2^{(n)}$	...	$x_d^{(n)}$

$(x^{(n)})^T$

En notación matricial

target  $\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}_{n \times 1}$

$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(n)})^T \end{bmatrix}_{n \times d+1}$

fila datos entre  
columnas

$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}_{d+1 \times 1}$

Modelo  $h_{\theta}(x^{(i)}) = (x^{(i)})^T \vec{\theta}$   
 $= \vec{\theta}^T x^{(i)}$

ej  $(x_0^{(i)}, x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)})$   
 $(\theta_0, \dots, \theta_d) \begin{pmatrix} x_0^{(i)} = 1 \\ x_1^{(i)} \\ \vdots \\ x_d^{(i)} \end{pmatrix}$

$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$   
 $= \frac{1}{2} \sum_{i=1}^n (x^{(i)T} \vec{\theta} - y^{(i)})^2$   
 $= \frac{1}{2} (\vec{X}\vec{\theta} - \vec{y})^T (\vec{X}\vec{\theta} - \vec{y})$

se reescribe

$J(\theta) = \frac{1}{2} (\vec{X}\vec{\theta} - \vec{y})^T (\vec{X}\vec{\theta} - \vec{y})$   
 $= \frac{1}{2} (\vec{\theta}^T \vec{X}^T \vec{X} \vec{\theta} - \vec{y}^T \vec{X} \vec{\theta} - \vec{y}^T \vec{X} \vec{\theta} + \vec{y}^T \vec{y})$   
 $= \frac{1}{2} (\vec{\theta}^T (\vec{X}^T \vec{X}) \vec{\theta} - 2 \vec{y}^T (\vec{X} \vec{\theta}) + \vec{y}^T \vec{y})$

\* use  $a^T b \leq b^T a$

Minimizamos  $J$  con respecto a  $\theta$  con

②

$$\frac{\partial J(\theta)}{\partial \theta_j} = 0$$

$$M^T = M$$

termino  $y^T y$  no contiene  $\theta_j$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{2} \left( \frac{\partial}{\partial \theta_j} \left( \vec{\theta}^T \underbrace{(X^T X)}_M \vec{\theta} - 2 \underbrace{y^T X}_N \vec{\theta} \right) \right)$$

expresando en coordenadas

$$\frac{\partial}{\partial \theta_j} \left\{ \sum_{ik} \theta_i M_{ik} \theta_k - \sum_e N_e \theta_e \right\} = \frac{\partial}{\partial \theta_j} \left\{ \sum_k \theta_i M_{ik} \theta_k - 2 N_j \theta_j \right\}$$

$$= \frac{\partial}{\partial \theta_j} \left\{ \sum_{k \neq j} \theta_i M_{ik} \theta_k + M_{jj} \theta_j^2 - 2 N_j \theta_j \right\}$$

$$= \sum_{k \neq j} M_{jk} \theta_k + 2 M_{jj} \theta_j + \sum_{i \neq j} \theta_i \underbrace{M_{ij}}_{M_{ji}} - 2 N_j$$

$$= 2 \sum_{k \neq j} M_{jk} \theta_k + 2 M_{jj} \theta_j - 2 N_j = 2 (M \vec{\theta} - N)_j$$

$$= 2 (X^T X \vec{\theta} - y^T X)_j \quad (\Rightarrow) \frac{\partial J}{\partial \theta_j} = (X^T X \vec{\theta} - y^T X)_j$$

En forma de vector

$$\nabla_{\theta} J = X^T X \vec{\theta} - y^T X$$

Iguando a cero

$$X^T X \vec{\theta} = y^T X = X^T y$$

Despejando

$$\vec{\theta} = \underbrace{(X^T X)^{-1}}_{X^+} X^T y$$

$X^+$  Moore-Penrose inverse

Dados datos  $X$  y  $y$  por multiplicación matricial obtenemos  $\theta$  parametros que minimizan coste de la regresión lineal ... sin iteraciones



## Punto de vista probabilístico sobre ML

Datos tienen incertidumbre inherente x siendo

- 1) Datos aleatorios (quantum) o errores aleatorios (clásico)
- 2) Variables no observadas que influyen en el mecanismo
- 3) Discretización en la toma de datos

Utilizamos conceptos de proba y estadística. v.a

• Independencia  $P(X=x, Y=y) = P(X=x) P(Y=y)$

• Def proba condicional  $P(Y=y | X=x) = \frac{P(Y=y, X=x)}{P(X=x)}$

• Regla de producto  $P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} | x^{(1)}, \dots, x^{(i-1)})$

• Regla de Bayes  $P(x | y) = \frac{P(y | x) P(x)}{P(y)}$   $P(x)$  ← prior knowledge w/o evidence

↑  
updated knowledge

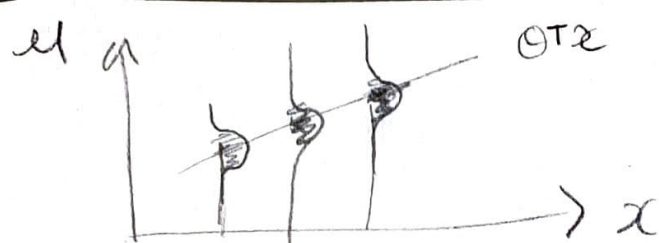
• Regla de suma  $P(y) = \sum_x P(y | x) P(x)$

## Problema ML regresión

$$\begin{cases} y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \\ \epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2) \otimes \end{cases}$$

Hay una relación lineal y la incertidumbre contenida en  $\epsilon^{(i)}$

\* quiere decir que  $P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\epsilon^{(i)})^2}{2\sigma^2}}$



en cada medida  
 todos los valores de  $y$  son  
 posibles siguiendo distribu-  
 ción de proba dado  $\theta$  como  
 parámetro

$$p(y^{(i)} | x^{(i)}; \theta) \quad \text{verosimilitud}$$

**Función de  $\theta$**  := Likelihood function

$$L(\theta) = p(\vec{y} | X; \theta) \quad \text{para escoger los mejores parámetros (donde máximo de } L \text{)}$$

$$L(\theta) = p(\vec{y} | X; \theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \quad \text{por ind.}$$

$$= \prod_{i=1}^n \mathcal{N}(y_i | x_i^T \theta, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

Buscamos  $\theta_{MLE} = \arg \max_{\theta} p(\vec{y} | X; \theta)$

$$= \arg \max_{\theta} \log p(\vec{y} | X; \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$$

$$= \arg \max_{\theta} \left\{ n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 \right\}$$

$$\theta_{MLE} = \arg \max_{\theta} \left\{ \underbrace{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2}_{\text{coste}} \right\}$$

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

$\propto$  distancia cuadrática

**Regresión lineal = MLE con ruido gaussiano**