

# PARCIAL 1 TAM

①

Teniendo el modelo de regresión

$$t_n = \phi(x_n)w^T + \eta_n,$$

con  $\{t_n \in \mathbb{R}, x_n \in \mathbb{R}^p\}_{n=1}^N$   
 target                      vector de entrada (características)

$w \in \mathbb{R}^Q \rightarrow$  parámetros (vector de pesos)

$\phi: \mathbb{R}^p \rightarrow \mathbb{R}^Q, Q \geq p$

$\eta_n \sim N(\eta_n | 0, \sigma_n^2) \rightarrow$  ruido gaussiano

Así tenemos la representación matricial del modelo como:

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$$

$\Downarrow$   
matriz de diseño  
 $\in \mathbb{R}^{N \times Q}$

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$$

$\Downarrow$   
vector de salida  
 $\in \mathbb{R}^N$

$$\eta = \begin{bmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_n \end{bmatrix}$$

$\Downarrow$   
ruido

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_Q \end{bmatrix}$$

$\Downarrow$   
vector de parámetros  
 $\in \mathbb{R}^Q$

$$t = \Phi w + \eta$$

Solución por  
mínimos  
cuadrados

$$J(w) = \sum_{n=1}^N (t_n - \phi(x_n)^T w)^2$$

$$= (t - \Phi w)^T (t - \Phi w)$$

$$= t^T t - t^T \Phi w - (\Phi w)^T t + (\Phi w)^T (\Phi w)$$

$$= t^T t - 2t^T \Phi w + w^T \Phi^T \Phi w$$

derivamos

$$\frac{dJ(w)}{dw} = -2\Phi^T t + 2\Phi^T \Phi w = 0$$

$$\Phi^T \Phi w = \Phi^T t$$

asumiendo que  $\Phi^T \Phi$  es invertible  $\Rightarrow (\Phi^T \Phi)^{-1} \Phi^T t$

## Minimos Cuadrados Regularizados

②

$$J(\omega) = \sum_{n=1}^N (t_n - \Phi(x_n)^T \omega)^2 + \lambda R(\omega)$$

Con  $\lambda \geq 0$

↑  
hiperparametro  
controla el grado  
de regularización

$R(\omega) \rightarrow$  función de penalización

$$\begin{cases} \text{Ridge } (\ell_2): R(\omega) = \|\omega\|_2^2 = \sum_{j=1}^Q \omega_j^2 \\ \text{Lasso } (\ell_1): R(\omega) = \|\omega\|_1 = \sum_{j=1}^Q |\omega_j| \end{cases}$$

Or Ridge

$$J(\omega) = (t - \Phi\omega)^T (t - \Phi\omega) + \lambda \omega^T \omega$$

$$J(\omega) = t^T t - 2t^T \Phi \omega + \omega^T \Phi^T \Phi \omega + \lambda \omega^T \omega$$

Con gradiente

$$\nabla_{\omega} J_{\text{ridge}} = -2\Phi^T t + 2\Phi^T \Phi \omega + 2\lambda \omega$$

Así tenemos

$$-2\Phi^T t + 2\Phi^T \Phi \omega + 2\lambda \omega = 0$$

$$-\Phi^T t + (\Phi^T \Phi + \lambda I) \omega = 0$$

$\Phi^T t$

$Q \times Q$   
Identidad

entonces:

$$\boxed{(\Phi^T \Phi + \lambda I)^{-1} \Phi^T t}$$

$\lambda \geq 0$

# MAXIMA VEROSIMILITUD (MLE)

modelo base:  $t_n = \Phi(x_n)^T \omega + \eta_n$  ;  $\eta_n \sim N(0, \sigma_n^2)$

función de verosimilitud:

$$P(t_n | \omega, \sigma_n^2) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(t_n - \Phi(x_n)^T \omega)^2}{2\sigma_n^2}\right)$$

Aplicamos log para simplificar el producto

$$\begin{aligned} \ln(\omega) &= \ln P(t_n | \omega, \sigma_n^2) \\ &= \sum_{n=1}^N \ln \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(t_n - \Phi(x_n)^T \omega)^2}{2\sigma_n^2}\right) \right) \\ &= \sum_{n=1}^N \left( -\frac{1}{2} \ln(2\pi\sigma_n^2) - \frac{(t_n - \Phi(x_n)^T \omega)^2}{2\sigma_n^2} \right) \\ &= -\frac{N}{2} \ln(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N (t_n - \Phi(x_n)^T \omega)^2 \end{aligned}$$

problema de optimización:

Maximizamos  $\ln(\omega)$  con respecto a  $\omega$ ,  $-\frac{N}{2} \ln(2\pi\sigma_n^2) \rightarrow$  constante  
entonces es equivalente a

$$\min_{\omega} J(\omega) = \min_{\omega} \sum_{n=1}^N (t_n - \Phi(x_n)^T \omega)^2$$

entonces podemos usar  $J(\omega) = \|t - \Phi\omega\|_2^2$

$$= (t - \Phi\omega)^T (t - \Phi\omega) = -2\Phi^T t + 2\Phi^T \Phi \omega = 0$$

$$= (\Phi^T \Phi)^{-1} \Phi^T t \rightarrow \text{Igual a mínimos cuadrados}$$

# MAXIMO A-POSTERIOR (MAP)

$$\hat{w}_{MAP} = \arg \max_w P(w|t)$$

Asumiendo independencia y distribución gaussiana

$$P(t|w) \propto P(t_n | \Phi(x_n)^T w, \sigma_n^2)$$

$$= (2\pi\sigma_n^2)^{-\frac{N}{2}} \exp\left(-\frac{1}{2\sigma_n^2} \|t - \Phi w\|_2^2\right)$$

Suponiendo prior gaussiano

$$P(w) = (2\pi\sigma_w^2)^{-\frac{Q}{2}} \exp\left(-\frac{1}{2\sigma_w^2} w^T w\right)$$

no normalizados

$$= \exp\left(-\frac{1}{2\sigma_n^2} \|t - \Phi w\|_2^2 - \frac{1}{2\sigma_w^2} w^T w\right)$$

Problema de optimización

$$\hat{w}_{MAP} = \arg \min_w \left\{ \frac{1}{\sigma_n^2} \|t - \Phi w\|_2^2 + \frac{1}{\sigma_w^2} w^T w \right\}$$

Como  $\lambda = \frac{\sigma_n^2}{\sigma_w^2} > 0$  entonces

$$\hat{w}_{MAP} = \arg \min_w \{ \|t - \Phi w\|_2^2 + \lambda w^T w \}$$

Así:

$$J(w) = \|t - \Phi w\|_2^2 + \lambda w^T w$$

$$J(w) = t^T t - 2t^T \Phi w + w^T \Phi^T \Phi w + \lambda w^T w \rightarrow \text{mínimos cuadrados regulados}$$

entonces

$$\boxed{(\Phi^T \Phi + \lambda I)^{-1} \Phi^T t}$$



# BAYESIANO PARA UN MODELO LINEAL GAUSSIANO

3

para parametro  $w$ , con distribucion gaussiana  $p(w) = N(w|0, \Sigma_w)$

Tenemos:

$$\text{modelo} = \phi(x_n)w^T + \eta_n$$

matriz  
de covarianza

$$\Sigma_w = \alpha^{-1} I$$

$$p(t|x, w, \sigma_n^2) = N(t|\phi w, \sigma_n^2 I)$$

encontramos la distribucion posterior de  $w$

$$p(w|t, x, \sigma_n^2) \propto p(t|x, w, \sigma_n^2) p(w)$$

asi que maximizamos el logaritmo de la probabilidad posterior

$$\log p(w|t) = \log p(t|x, w, \sigma_n^2) + \log(p(w)) +$$

$$= -\frac{1}{2\sigma_n^2} \|t - \phi w\|^2 - \frac{1}{2} w^T \Sigma_w^{-1} w$$

asi minimizamos

$$J(w) = \frac{1}{2\sigma_n^2} \|t - \phi w\|^2 + \frac{1}{2} w^T \Sigma_w^{-1} w$$

convina verosimilitud y la regularizacion con

$$S_N = (S_m^{-1} + \beta \phi^T \phi)^{-1}$$

$$M_N = S_N (S_m^{-1} \mu_m + \beta \phi^T t)$$

$$\text{prior} \rightarrow \mu_m = 0 \text{ y } S_0 = \alpha^{-1} I$$

$$S_N = (\alpha I + \beta \phi^T \phi)^{-1}$$

$$M_N = \beta S_N \phi^T t = (\alpha \beta^{-1} I + \phi^T \phi)^{-1} \phi^T t$$

siendo  $\alpha \beta^{-1} = \alpha \sigma_n^2$  se vuelve maximo o posteriori

# REGRESION RIGIDA KERNEL

minimizamos el error cuadrático

$$J(\omega) = \sum_{n=1}^N (t_n - \phi(x_n) \omega^T)^2 + \lambda \|\omega\|^2$$

$$\omega = \sum_{n=1}^N \alpha_n \phi(x_n)$$

sustituimos para  $x$

$$f(x) = \sum_{n=1}^N \alpha_n \underbrace{\phi(x_n)^T \phi(x)}_{\rightarrow \text{Kernel}}$$

así problema de optimización

$$J(\alpha) = \|t - K\alpha\|^2 + \lambda \alpha^T K \alpha$$

matriz de kernel

minimizamos  $J(\alpha)$

$$\alpha = (K + \lambda I)^{-1} t$$

$$f(x_*) = \kappa(x_*, x) (K + \lambda I)^{-1} t$$

$\sigma_n^2 \rightarrow$  no está  
explícita en la  
predicción