

## Parcial 1:

modelo :  $t_n = \phi(x_n) w + \eta_n ; \quad n = 1, \dots, N.$

- $t_n \in \mathbb{R} \rightarrow$  Respuesta
- $x_n \in \mathbb{R}^p \rightarrow$  Entrada.
- $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^Q \rightarrow$  Características
- $w \in \mathbb{R}^Q \rightarrow$  Parámetros a Estimar
- $\eta_n \rightarrow$  ruido independiente  $\eta_n \sim \mathcal{N}(0, \sigma^2)$
- Matriz de diseño  $\Phi \in \mathbb{R}^{N \times Q}$  filas  $\phi(x_n)^T$ .
- Vector observados

## 1) Mínimos Cuadrados:

minimizar Error / Pérdida:

Forma Vectorial



$$J(w) = \frac{1}{2} \|t - \Phi w\|_2^2 = \frac{1}{2} (t - \Phi w)^T (t - \Phi w)$$

$$\begin{aligned} (t - \Phi w)^T (t - \Phi w) &= t^T t - t^T \Phi w - \Phi^T w^T t + \Phi^T w^T \Phi w \\ &= t^T t - 2 w^T \Phi^T t + w^T \Phi^T \Phi w \end{aligned}$$

\*  $t^T \Phi w \rightarrow$  Escalar / lo transpuesto es igual al esc.

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

$$J(w) = \frac{1}{2} (t^T t - 2 w^T \Phi^T t + w^T \Phi^T \Phi w)$$

$$= \frac{1}{2} \underbrace{t^T t}_0 - w^T \underbrace{\Phi^T t}_0 + \frac{1}{2} w^T \underbrace{\Phi^T \Phi}_A w$$

\* Derivato  $J(w) / w$ .

$$A = \Phi^T \Phi$$

$$\cdot \frac{\partial}{\partial w} \frac{1}{2} (w^T A w) = \frac{1}{2} (A + A^T) w = \frac{1}{2} (2A w) = A w$$

Termino Quadrato

Termino Lineal:

$$\frac{\partial}{\partial w} (\Phi^T t) w = \Phi^T t$$

$$\cdot \frac{\partial J}{\partial w} = -\Phi^T t + A w = -\Phi^T t + \Phi \Phi^T w$$

igualamos a zero "condicion optima gradiente"

$$-\Phi^T t + \Phi^T \Phi w = 0$$

↳ minimo global  
Gradiente = 0.

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

$$w_{ols} = (\Phi^T \Phi)^{-1} \Phi^T t$$



## 2) Mínimos Cuadrados Regularizados:

$$J_{\lambda}(w) = \underbrace{\frac{1}{2} \|t - \Phi w\|_2^2}_{T_1} + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{T_2} \quad \lambda > 0.$$

$$\|w\|_2^2 = w^T w$$

↑  
Parámetro de regularización

$$\hookrightarrow \frac{\lambda}{2} w^T w$$

Derivadas:

$$T_1: -\Phi^T t + \Phi \Phi^T w \quad - \text{Calculada anteriormente.}$$

$$T_2: \lambda w$$

$$\text{Gradiente: } -\Phi^T t + \Phi \Phi^T w + \lambda w$$

→ Condiciones de mínimo

$$\nabla_w J_{\lambda}(w) = 0$$

$$-\bar{\Phi}^T t + \bar{\Phi}^T \bar{\Phi} w + \lambda w = 0$$

$$(\bar{\Phi}^T \bar{\Phi} + \lambda I) w = \bar{\Phi}^T t$$

$$w_{ridge} = (\bar{\Phi}^T \bar{\Phi} + \lambda I)^{-1} \bar{\Phi}^T t$$

### 3.) Maxima Verosimilitud.

• Verosimilitud.  $p(t|w)$

$$p(t_n | w) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(t_n - \phi(x_n)^T w)^2}{2\sigma_n^2}\right)$$

• función de densidad Normal.

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

• independientes, Ver. Conjunta

$$p(t|w) = \prod_{n=1}^N p(t_n | w) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(t_n - \phi(x_n)^T w)^2}{2\sigma_n^2}\right)$$

$$p(t|w) = (2\pi\sigma_n^2)^{-N/2} \exp\left(-\frac{1}{2\sigma_n^2} \sum_{n=1}^N (t_n - \phi(x_n)^T w)^2\right)$$

Nota:  $\sum (t_n - \phi_n^T w)^2 = \|t - \Phi w\|_2^2$



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

• Log - Verosimilitud

$$L(w) : \log P(t|w) = -\frac{N}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \|t - \phi w\|_2^2$$

↑  
No depende de  $w$ .

$$\arg \max L(w) : \arg \max \left( -\frac{1}{2\sigma_n^2} \|t - \phi w\|_2^2 \right)$$

Nota: maximizar la Ver. similitud. = minimizar la suma de cuadrados

$$= \arg \min \frac{1}{2\sigma_n^2} \|t - \phi w\|_2^2$$

• Derivamos Log - Likelihood

$$g(w) = \frac{1}{2\sigma_n^2} \|t - \phi w\|_2^2$$

$$\text{Derivada: } \frac{\partial}{\partial w} \|t - \phi w\|_2^2 = (t - \phi w)(t - \phi w)^T$$

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

$$\nabla_{\omega} l(\omega) = \frac{1}{\sigma^2} \phi^T (t - \phi\omega)$$

\* Máximo interior / mínimo de la suma de  $^2$  ds.

$$\nabla_{\omega} l(\omega) = 0 \rightarrow \phi^T (t - \phi\omega) = 0.$$

$$\phi \phi^T \omega = \phi^T t$$

$$\hat{\omega}_{MLE} = (\phi^T \phi)^{-1} \phi^T t$$

\* Segunda derivada:

$$H(\omega) = -\frac{1}{\sigma^2} \phi^T \phi$$

Nota:

$$1. \text{ Likelihood: } p(t|\omega) = (2\pi\sigma^2)^{-N/2} \exp\left(-\frac{1}{2\sigma^2} \|t - \phi\omega\|_2^2\right)$$

$$2. \text{ log. Lik.: } l(\omega) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|t - \phi\omega\|_2^2$$

$$3. \text{ Grad: } \nabla_{\omega} l(\omega) = \frac{1}{\sigma^2} \phi^T (t - \phi\omega)$$

$$4. = 0 : \phi^T \phi \omega = \phi^T t$$

$$5. \text{ Soluc.: } (\phi^T \phi)^{-1} \phi^T t$$

$$6. \text{ Hessian: } H = -\frac{1}{\sigma^2} \phi^T \phi \rightarrow \text{máximo gl. b. b.}$$



#### 4. max. a-Posteriori:

Suponemos un Prior (a Prior)  $w$  + también gaussiano

$$P(w) = \mathcal{N}(w | 0, \sigma_w^2 I)$$

a) Likelihood  $\rightarrow p(t|w, \sigma^2)$

$$= (2\pi \sigma^2)^{N/2} \exp \left[ -\frac{1}{2\sigma^2} (t - \phi w)^T (t - \phi w) \right]$$

b) Prior.  $P(w | \sigma_w^2) = (2\pi \sigma_w^2)^{-M/2} \exp \left[ -\frac{1}{2\sigma_w^2} w^T w \right]$

• log Posterior.

$$\begin{aligned} \log p(w|t) &= \log p(t|w) + \log p(w) + \text{const.} \\ &= -\frac{1}{2\sigma^2} (t - \phi w)^T (t - \phi w) - \frac{1}{2\sigma_w^2} w^T w + \text{const.} \end{aligned}$$

• Maximo Posterior: / minimizar su negativo

$$W_{MAP} = \arg \min \left[ \underbrace{\frac{1}{2\sigma^2} (t - \phi w)^T (t - \phi w)}_{\text{Error Cuadrático}} + \underbrace{\frac{1}{2\sigma_w^2} w^T w}_{\text{Penalización por pesos grandes, Regularización}} \right]$$

Error Cuadrático

Penalización por pesos grandes, Regularización

$$\lambda = \frac{\sigma^2}{\sigma_w^2}$$

$$J(w) = \frac{1}{2} (t - \phi w)^T (t - \phi w) + \frac{\sigma^2}{2\sigma_w^2} w^T w$$

$$J(w) = \frac{1}{2} \|t - \phi w\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

Derivamos  $J$  w.r.t  $w$ .

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

$$\bullet \nabla \frac{1}{2} w^T \phi^T \phi w = \frac{1}{2} (2 \phi^T \phi w) = \phi^T \phi w \rightarrow w^T = w \rightarrow \text{Symetric}$$

$$\bullet \nabla (-t^T \phi w) = -\phi^T t$$

$$\bullet \nabla \frac{\lambda}{2} w^T w = \frac{\lambda}{2} (2w) = \lambda w$$

Nota: Regras:  $\nabla$  matrices

$$w^T A w \quad 2Aw$$

$$C^T w \quad C$$

$$w^T w \quad 2w$$



$$\nabla J(w) = -\underbrace{\phi^T (t - \phi w)}_{\text{Grad. Error}} + \lambda w.$$

Error Pred / v real

Grad. term. no  
regularizada.  
Evita sobreajuste.

$$\nabla = 0.$$

$$\phi^T \phi w + \lambda w = \phi^T t.$$

$$(\phi^T \phi + \lambda I) w = \phi^T t$$

$$w_{\text{máx}} = (\underbrace{\phi^T \phi}_{\text{Con. Caracter.}} + \lambda I)^{-1} \phi^T t \quad ; \quad \lambda = \frac{\sigma^2}{G^2 w}$$

↓  
Prior Gaussiano.

## 5. Bayesiano Con modelo lin G.

Usamos:

$$\text{Likelihood: } p(t|w) = \mathcal{N}(t | \phi w, \sigma^2 I)$$

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

$$\text{Prior: } p(w) = \mathcal{N}(w | 0, \Sigma_p)$$

$$= (2\pi)^{-Q/2} |\Sigma_p|^{-1/2} \exp\left(-\frac{1}{2} w^T \Sigma_p^{-1} w\right)$$

2)  $P(w) = P_{prior} \times \text{likelihood}$

$$\log p(w|t) = -\frac{1}{2\sigma^2} (t - \phi w)^T (t - \phi w) - \frac{1}{2} w^T \Sigma_p^{-1} w + \text{const}$$

$$= -\frac{1}{2\sigma^2} (t^T t - 2w^T \phi^T t + w^T \phi^T \phi w) - \frac{1}{2} w^T \Sigma_p^{-1} w + \text{const}$$

$$= \frac{1}{2} w^T \left( \Sigma_p^{-1} + \frac{1}{\sigma^2} \phi^T \phi \right) w - 2 \left( \frac{1}{\sigma^2} \phi^T t \right) w + \text{const}$$

$$A := \Sigma_p^{-1} + \frac{1}{\sigma^2} \phi^T \phi \quad b := \frac{1}{\sigma^2} \phi^T t$$

$$= \frac{1}{2} w^T A w - 2 b^T w$$

Complete the Square

$$w^T A w - 2 b^T w = (w - A^{-1} b)^T A (w - A^{-1} b) - b^T A^{-1} b$$

$$\log p(w|t) = -\frac{1}{2} (w - A^{-1} b)^T A (w - A^{-1} b)$$

$$S_N: A^{-1} = \left( \Sigma_p^{-1} + \frac{1}{\sigma^2} \phi^T \phi \right)^{-1}$$



$$m_N : A^{-1} b = \frac{1}{G^2} S_w \phi^T t.$$

## Regresión Rígida por Kernel

Queremos permitir no linealidad sin trabajar explícitamente con  $w$  en  $\mathbb{R}^Q$ , usando un kernel  
 $k(x, x') = \phi(x)^T \phi(x')$

Ridge en espacio de características

$$\min_w \|t - \Phi w\|^2 + \lambda \|w\|^2$$

Solución dual (teorema de representación):

La solución  $w$  pertenece al subespacio generado por las filas de  $\Phi$ , es decir, existe  $\alpha \in \mathbb{R}^N$  tal que  $w = \Phi^T \alpha$ .

$$t - \Phi w = t - \Phi \Phi^T \alpha = t - K \alpha$$

$K = \Phi \Phi^T \in \mathbb{R}^{N \times N}$  es la matriz de kernel con  $K_{ij} = k(x_i, x_j)$

El problema en  $\alpha$  es:

$$\min_{\alpha} \|t - K \alpha\|^2 + \lambda \alpha^T \Phi \Phi^T \alpha = \|t - K \alpha\|^2 + \lambda \alpha^T K \alpha$$

Derivando e igualando a cero:

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

La predicción para un nuevo  $x_*$  usa:

$$\hat{t}_* = \phi(x_*)^T w = \phi(x_*)^T \Phi^T \alpha = k_*^T \alpha$$

Donde  $k_* = (k(x_*, x_1), \dots, k(x_*, x_N))^T$

$$\text{Por lo tanto: } \hat{t}_* = k_*^T (K + \lambda I)^{-1} t$$

Interpretación: kernel ridge permite operar con  $k$  directamente sin calcular  $\Phi$ . Es equivalente a Ridge en el espacio de características.



## Regresión Por proceso Gaussiano

Un proceso Gaussiano (GP) pone un prior directamente sobre funciones  $f(\cdot)$

$$f(\cdot) \sim \text{GP}(0, k(\cdot, \cdot))$$

$$t_n = f(x_n) + \eta_n \text{ con } \eta_n \sim \mathcal{N}(0, \sigma_\eta^2)$$

Prior conjunto sobre valores observados y no observados

Sea  $f = [f(x_1), \dots, f(x_N)]^T$  entonces

$$P(f) = \mathcal{N}(0, K), \quad K_{ij} = k(x_i, x_j)$$

$$t = f + \eta \text{ con } \eta \sim \mathcal{N}(0, \sigma_\eta^2 I) \text{ por lo tanto.}$$

$$P(t) = \mathcal{N}(0, K + \sigma_\eta^2 I)$$

Predictiva para  $x_*$

La distribución conjunta de  $\begin{pmatrix} t \\ f_* \end{pmatrix}$  es gaussiana, condicionando se obtiene la predictiva nominal.

$$P(f_* | x_*, t) = \mathcal{N}(\mu_*, \sigma_*^2)$$

$$\mu_* = K_* (K + \sigma_\eta^2 I)^{-1} t$$

$$\sigma_*^2 = K(x_*, x_*) - K_* (K + \sigma_\eta^2 I)^{-1} K_*$$

Si deseamos la predictiva sobre  $t_* = f_* + \eta_*$  añadiremos  $\sigma_\eta^2$  a la varianza.

Relación con kernel ridge y modelo Bayesiano lineal

$$\text{la media predictiva de GPR: } \mu_* = K_*^T (K + \sigma_\eta^2 I)^{-1} t$$

Es idéntica a la predicción de kernel ridge con  $\lambda = \sigma_\eta^2$

GP y kernel ridge Comparten la misma media predictiva cuando se toma hiperparámetros coincidentes. La diferencia principal es que GP entrega una Varianza Predictiva natural.

El GP puede verse como el límite no paramétrico del modelo lineal bayesiano cuando el número de Características  $Q$  tiende a infinito apropiadamente y el prior induce el kernel  $k$ .