

## Formulas General del problema de optimización

$$\min_W \left( \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - \phi(X_n) W^T)^2}_{\text{Función de pérdida}} + \underbrace{\frac{\lambda}{2} \|W\|^2}_{\text{término de regularización } L_2} \right)$$

- La función de pérdida representa el error cuadrático medio entre las predicciones y los valores reales.
- El término de regularización es la norma  $L_2$  de los pesos, que penaliza los valores grandes de  $W$ .
- $\lambda$  controla la fuerza de penalización.

### 1. Mínimos cuadrados ordinarios (OLS)

Se desea minimizar el error entre las predicciones y los valores reales.

$$\min_W \sum_{n=1}^N (t_n - \phi(X_n) W^T)^2$$

- Se define la matriz de diseño  $\Phi (N \times Q)$ , donde cada fila es  $\phi(X_n)$ .
- El vector de observaciones es  $t = [t_1, \dots, t_N]^T$ .
- El error cuadrático se reescribe como:

$$J(W) = (t - \Phi W)^T (t - \Phi W)$$

Para encontrar el  $W$  que minimiza  $J(W)$ , necesitamos calcular la derivada de  $J(W)$  con respecto a  $W$  e igualarla a cero.

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

Simplificamos los términos. dado que  $\mathbf{t}^T \Phi \mathbf{W}$  es un escalar, es igual a su traspuesta  $(\mathbf{t}^T \Phi \mathbf{W})^T = \mathbf{W}^T \Phi^T \mathbf{t}$ .

$$J(\mathbf{W}) = \mathbf{t}^T \mathbf{t} - 2 \mathbf{W}^T \Phi^T \mathbf{t} + \mathbf{W}^T \Phi^T \Phi \mathbf{W}$$

Calculamos la derivada de  $J(\mathbf{W})$  con respecto a  $\mathbf{W}$

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = -2 \Phi^T \mathbf{t} + 2 \Phi^T \Phi \mathbf{W}$$

Igualemos a cero

$$\begin{aligned} -2 \Phi^T \mathbf{t} + 2 \Phi^T \Phi \mathbf{W} &= 0 \\ \Phi^T \Phi \mathbf{W} &= \Phi^T \mathbf{t} \end{aligned}$$

Si la matriz  $\Phi^T \Phi$  es invertible (lo que significa que las columnas de  $\Phi$  son linealmente independientes y  $N \geq M$ ), podemos multiplicar ambos lados por  $(\Phi^T \Phi)^{-1}$

$$\mathbf{W} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

## 2. Mínimos cuadrados regularizados

Para prevenir el sobre ajuste, se añade un término de penalización a la función de costo. La regularización L2 penaliza la suma de los cuadrados de los cuadrados de los coeficientes  $\mathbf{W}$ .

$$J_R(\mathbf{W}) = \sum_{n=1}^N (t_n - \phi(\mathbf{x}_n)^T \mathbf{W})^2 + \lambda \|\mathbf{W}\|_2^2$$

$\lambda$  es el parámetro de regularización, y  $\|\mathbf{W}\|_2^2 = \mathbf{W}^T \mathbf{W}$  en forma matricial

$$\begin{aligned} J_R(\mathbf{W}) &= (\mathbf{t} - \Phi \mathbf{W})^T (\mathbf{t} - \Phi \mathbf{W}) + \lambda \mathbf{W}^T \mathbf{W} \\ &= \mathbf{t}^T \mathbf{t} - 2 \mathbf{W}^T \Phi^T \mathbf{t} + \mathbf{W}^T \Phi^T \Phi \mathbf{W} + \lambda \mathbf{W}^T \mathbf{W} \end{aligned}$$

Para los primeros tres términos la derivada es la misma que en OLS. Solo necesitamos la derivada del término de regularización:

$$\frac{\partial (\lambda \mathbf{W}^T \mathbf{W})}{\partial (\mathbf{W})} = \lambda \frac{\partial (\mathbf{W}^T \mathbf{I} \mathbf{W})}{\partial (\mathbf{W})} = \lambda (2 \mathbf{I} \mathbf{W}) = 2 \lambda \mathbf{W}$$

Entonces la derivada completa es:

$$\frac{\partial J_R(W)}{\partial W} = -Z\Phi^T t + Z\Phi^T \Phi W + 2\lambda W$$

Iguálamos a cero y despejamos  $W$

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

$$\Phi^T \Phi W + \lambda W = \Phi^T t$$

$$(\Phi^T \Phi + \lambda I) W = \Phi^T t$$

La matriz  $(\Phi^T \Phi + \lambda I)$  es invertible para  $\lambda > 0$

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

### 3. Máxima verosimilitud (MLE)

Problema de optimización:

El modelo de regresión es  $t_n = \phi(X_n)^T W + \eta$  con  $\eta_n \sim N(0, \sigma_n^2)$ .  
esto implica que  $t_n$  dado  $X_n, W, \sigma_n^2$  sigue una distribución gaussiana:

$$p(t_n | X_n, W, \sigma_n^2) = N(t_n | \phi(X_n)^T W, \sigma_n^2)$$

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

Dado que los términos son iid, la función de verosimilitud para todos los datos  $t$  (dado  $X, W$  y  $\sigma_n^2$ ) es el producto de las probabilidades individuales:

$$L(W, \sigma_n^2) = p(t | X, W, \sigma_n^2) = \prod_{n=1}^N p(t_n | X_n, W, \sigma_n^2)$$

$$L(W, \sigma_n^2) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(t_n - \phi(X_n)^T W)^2}{2\sigma_n^2}\right)$$

$$L(W, \sigma_n^2) = \left(\frac{1}{2\pi\sigma_n^2}\right)^{N/2} \exp\left(-\frac{1}{2\sigma_n^2} \sum_{n=1}^N (t_n - \phi(X_n)^T W)^2\right)$$

Es computacional más conveniente maximizar la log-verosimilitud (log-likelihood):

$$\mathcal{L}(W, \sigma_n^2) = \ln \mathcal{L}(W, \sigma_n^2)$$

$$\mathcal{L}(W, \sigma_n^2) = \frac{N}{2} \ln(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N (t_n - \phi(X_n)^T W)^2$$

En forma matricial, la suma de cuadrados es  $(t - \Phi W)^T (t - \Phi W)$ .

$$\mathcal{L}(W, \sigma_n^2) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\sigma_n^2) - \frac{1}{2\sigma_n^2} (t - \Phi W)^T (t - \Phi W)$$

El problema de optimización es encontrar  $(W$  y  $\sigma_n^2)$  que maximicen  $\mathcal{L}(W, \sigma_n^2)$ .

Calculamos la derivada de  $\mathcal{L}(W, \sigma_n^2)$  con respecto a  $W$ .

Los términos  $-\frac{N}{2} \ln(2\pi)$  y  $-\frac{N}{2} \ln(\sigma_n^2)$  no dependen de  $W$ , por lo que su derivada es cero.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W} &= \frac{\partial}{\partial W} \left( -\frac{1}{2\sigma_n^2} (t - \Phi W)^T (t - \Phi W) \right) \\ &= -\frac{1}{2\sigma_n^2} \frac{\partial}{\partial W} (t^T t - 2W^T \Phi^T t + W^T \Phi^T \Phi W) \\ &= -\frac{1}{2\sigma_n^2} (-2\Phi^T t + 2\Phi^T \Phi W) \\ &= \frac{1}{\sigma_n^2} (\Phi^T t - \Phi^T \Phi W) \end{aligned}$$

Igualemos la derivada a cero y resolvamos para  $W$

$$\frac{1}{\sigma_n^2} (\Phi^T t - \Phi^T \Phi W) = 0$$

Assumiendo  $\sigma_n^2 \neq 0$ :

$$\begin{aligned} \Phi^T t - \Phi^T \Phi W &= 0 \\ \Phi^T \Phi W &= \Phi^T t \end{aligned}$$

$$W_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T t$$

Notamos que  $W_{MLE} = W_{OLS}$



Solución para  $\sigma_u^2$ , MLE

Calculamos la derivada de  $\mathcal{L}(W, \sigma_u^2)$  con respecto a  $\sigma_u^2$

$$\mathcal{L}(W, \sigma_u^2) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\sigma_u^2) - \frac{1}{2\sigma_u^2} \text{SSE}(W)$$

donde  $\text{SSE}(W) = (t - \Phi W)^T (t - \Phi W)$ .

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \sigma_u^2} &= -\frac{N}{2} \frac{1}{\sigma_u^2} - \left( -\frac{1}{2(\sigma_u^2)^2} \right) \text{SSE}(W) \\ &= -\frac{N}{2\sigma_u^2} + \frac{\text{SSE}(W)}{2(\sigma_u^2)^2} \end{aligned}$$

Iguamos a cero y resolvemos para  $\sigma_u^2$

$$-\frac{N}{2\sigma_u^2} + \frac{\text{SSE}(W)}{2(\sigma_u^2)^2} = 0$$

Multiplcamos por  $2(\sigma_u^2)^2$

$$-N\sigma_u^2 + \text{SSE}(W) = 0$$

$$N\sigma_u^2 = \text{SSE}(W)$$

Solucionamos para  $\sigma_u^2$ , MLE

$$\sigma_{u, \text{MLE}}^2 = \frac{\text{SSE}(W_{\text{MLE}})}{N}$$

$$= \frac{1}{N} (t - \Phi W_{\text{MLE}})^T (t - \Phi W_{\text{MLE}})$$

$$= \frac{1}{N} \sum_{n=1}^N (t_n - \phi(x_n)^T W_{\text{MLE}})^2$$

#### 4. Máximo a-posteriori (Maximum a-posteriori - MAP)

Problema de optimización:

En la estimación MAP, incorporamos conocimiento previo sobre los parámetros  $W$  a través de una distribución a priori  $p(W)$ .

El objetivo es maximizar la probabilidad a posteriori  $p(W | t, X)$ , que por el teorema de Bayes, es proporcional al producto de la verosimilitud y la prior:

$$p(W | t, X) \propto p(t | X, W) p(W)$$

↳ "es proporcional a"

Maximizar la posterior es equivalente a maximizar su logaritmo:

$$\ln p(W | t, X) = \ln p(t | X, W) + \ln p(W) + \text{constante}$$

Assumimos una prior gaussiana para  $W$  con media cero y covarianza isotrópica  $\alpha^{-1} I_Q$ :

$$\begin{aligned} p(W | \alpha) &= \mathcal{N}(W | 0, \alpha^{-1} I_Q) \\ &= \left(\frac{\alpha}{2\pi}\right)^{Q/2} \exp\left(-\frac{\alpha}{2} W^T W\right) \end{aligned}$$

La log-prior es:

$$\ln p(W | \alpha) = \frac{Q}{2} \ln\left(\frac{\alpha}{2\pi}\right) - \frac{\alpha}{2} W^T W$$

El log-posterior (ignorando constante que no dependen de  $W$ ), es:

$$\ln p(W | t, X, \sigma_n^2, \alpha) \propto -\frac{1}{2\sigma_n^2} (t - \Phi W)^T (t - \Phi W) - \frac{\alpha}{2} W^T W$$

Maximizar esta cantidad es equivalente a minimizar su negación:

$$J_{\text{MAP}}(W) = \frac{1}{2\sigma_n^2} (t - \Phi W)^T (t - \Phi W) + \frac{\alpha}{2} W^T W$$

Si multiplicamos por  $2\sigma_n^2$  (asumiendo  $\sigma_n^2$  es una constante conocida), minimizar esto:

$$(t - \Phi W)^T (t - \Phi W) + \alpha \sigma_n^2 W^T W$$

esto es idéntico a la función de costo de la Regresión Ridge si definimos  $\lambda = \alpha \sigma_u^2$ .

Tomamos la derivada de la función objetivo (negativo de log-posterior sin constantes) con respecto a  $W$ .

$$\text{Sea } f(W) = \frac{1}{2\sigma_u^2} (\mathbf{t} - \Phi W)^T (\mathbf{t} - \Phi W) + \frac{\alpha}{2} W^T W$$

$$\begin{aligned} \frac{\partial f(W)}{\partial W} &= \frac{1}{2\sigma_u^2} (-2 \Phi^T \mathbf{t} + 2 \Phi^T \Phi W) + \frac{\alpha}{2} (2W) \\ &= -\frac{1}{\sigma_u^2} \Phi^T \mathbf{t} + \frac{1}{\sigma_u^2} \Phi^T \Phi W + \alpha W \end{aligned}$$

Iguémoslo a cero y resolvámoslo para  $W$

$$-\frac{1}{\sigma_u^2} \Phi^T \mathbf{t} + \frac{1}{\sigma_u^2} \Phi^T \Phi W + \alpha W = 0$$

Multiplicamos por  $\sigma_u^2$

$$-\Phi^T \mathbf{t} + \Phi^T \Phi W + \alpha \sigma_u^2 W = 0$$

$$\Phi^T \Phi W + \alpha \sigma_u^2 W = \Phi^T \mathbf{t}$$

$$(\Phi^T \Phi + \alpha \sigma_u^2 \mathbf{I}) W = \Phi^T \mathbf{t}$$

Obtenemos la solución para  $W_{\text{MAP}}$

$$W_{\text{MAP}} = (\Phi^T \Phi + \alpha \sigma_u^2 \mathbf{I})^{-1} \Phi^T \mathbf{t}$$

Notese que es idéntico a  $W_{\text{Ridge}}$  con  $\lambda = \alpha \sigma_u^2$

### 5. Modelo Bayesiano con Modelo Lineal Gaussiano

En un enfoque Bayesiano completo, no se busca una estimación puntual de  $W$ , si no que calcular la posterior completa

$$p(W | \mathbf{t}, X, \sigma_y^2, \alpha).$$

La verosimilitud es  $p(\mathbf{t} | X, W, \sigma_y^2) = \mathcal{N}(\mathbf{t} | \Phi W, \sigma_y^2 \mathbf{I}_N)$ .

La prior es  $p(W | \alpha) = \mathcal{N}(W | 0, \alpha^{-1} \mathbf{I}_Q)$ .

El logaritmo de la posterior es proporcional a:

$$\ln p(W | \mathbf{t}, X, \sigma_y^2, \alpha) \propto -\frac{1}{2\sigma_y^2} (\mathbf{t} - \Phi W)^T (\mathbf{t} - \Phi W) - \frac{\alpha}{2} W^T W$$

Expandiendo los términos cuadráticos en  $W$

$$(\mathbf{t} - \Phi W)^T (\mathbf{t} - \Phi W) = \mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \Phi W + W^T \Phi^T \Phi W$$

Así, el exponente de la posterior (Multiplicado por  $-\frac{1}{2}$ ) es:

$$\frac{1}{\sigma_y^2} (\mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \Phi W + W^T \Phi^T \Phi W) + \alpha W^T W$$

$$\frac{1}{\sigma_y^2} W^T \Phi^T \Phi W + \alpha W^T W - \frac{2}{\sigma_y^2} \mathbf{t}^T \Phi W + \frac{1}{\sigma_y^2} \mathbf{t}^T \mathbf{t}$$

$$W^T \left( \frac{1}{\sigma_y^2} \Phi^T \Phi + \alpha \mathbf{I}_Q \right) W - 2 W^T \left( \frac{1}{\sigma_y^2} \Phi^T \mathbf{t} \right) + \text{constantes.}$$

Esta es la forma del exponente de una Gaussiana multivariada

$$\mathcal{N}(W | \mu_N, \Sigma_N),$$

cuyo logaritmo (ignorando constantes) es

$$-\frac{1}{2} (W - \mu_N)^T \Sigma_N^{-1} (W - \mu_N) = -\frac{1}{2} (W^T \Sigma_N^{-1} W - 2 W^T \Sigma_N^{-1} \mu_N + \mu_N^T \Sigma_N^{-1} \mu_N)$$

Comparando los términos cuadráticos en  $W$  (el término  $W^T \mathbf{A} W$ ), identificamos la precisión (inversa de la covarianza) de la posterior:

$$\Sigma_N^{-1} = \frac{1}{\sigma_y^2} \Phi^T \Phi + \alpha \mathbf{I}_Q$$



la covarianza de la posterior es:

$$\Sigma_N = (\alpha \mathbf{I}_Q + \frac{1}{\sigma_u^2} \Phi^T \Phi)^{-1}$$

Comparando los términos lineales en  $\mathbf{W}$  (el término  $-\mathbf{z} \mathbf{W}^T \mathbf{b}$ ), identificamos:

$$\Sigma_N^{-1} \mu_N = \frac{1}{\sigma_u^2} \Phi^T \mathbf{t}$$

la media de la posterior es:

$$\mu_N = \Sigma_N \left( \frac{1}{\sigma_u^2} \Phi^T \mathbf{t} \right)$$

$$\mu_N = (\alpha \mathbf{I}_Q + \frac{1}{\sigma_u^2} \Phi^T \Phi)^{-1} \frac{1}{\sigma_u^2} \Phi^T \mathbf{t}$$

La distribución posterior para  $\mathbf{W}$  es Gaussiana:

$$p(\mathbf{W} | \mathbf{t}, \mathbf{X}, \sigma_u^2, \alpha) = \mathcal{N}(\mathbf{W} | \mu_N, \Sigma_N)$$

con:

- Media:  $\mu_N = (\Phi^T \Phi + \alpha \sigma_u^2 \mathbf{I}_Q)^{-1} \Phi^T \mathbf{t}$  (que es igual a  $\mathbf{W}_{MAP}$ )
- Covarianza:  $\Sigma_N = (\alpha \mathbf{I}_Q + \frac{1}{\sigma_u^2} \Phi^T \Phi)^{-1}$

Predicción Bayesiana:

Para una nueva entrada  $\mathbf{x}_*$ , la distribución predictiva para  $t_* = \phi(\mathbf{x}_*)^T \mathbf{W} + \eta_*$  se obtiene integrando sobre la posterior de  $\mathbf{W}$ :

$$p(t_* | \mathbf{x}_*, \mathbf{t}, \mathbf{X}, \sigma_u^2, \alpha) = \int p(t_* | \mathbf{x}_*, \mathbf{W}, \sigma_u^2) p(\mathbf{W} | \mathbf{t}, \mathbf{X}, \sigma_u^2, \alpha) d\mathbf{W}$$

Esto resulta en una Gaussiana  $\mathcal{N}(t_* | \mu_*, \sigma_*^2)$  con:

- Media predictiva:  $\mu_* = \phi(\mathbf{x}_*)^T \mu_N$
- Varianza predictiva:  $\sigma_*^2 = \sigma_u^2 + \phi(\mathbf{x}_*)^T \Sigma_N \phi(\mathbf{x}_*)$

## 6. Regresión Rígida Kernel (Kernel Ridge Regression - KRR)

Partiendo de la regresión ridge  $W = (\Phi^T \Phi + \lambda I_Q)^{-1} \Phi^T t$

Se puede demostrar que  $W$  se puede expresar como una combinación lineal de las transformaciones de las características de entrada:  $W = \Phi^T a$  partir de algún vector  $a \in \mathbb{R}^N$ .  
Usando la identidad matricial:

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} y = P B^T (R + B P B^T)^{-1} y.$$

Aplicándola a  $W$

$$W = (\lambda I_Q + \Phi^T I_N^{-1} \Phi)^{-1} \Phi^T I_N^{-1} t \quad (\text{con } P = \lambda^{-1} I_Q, B = \Phi, R = I_N, y = t).$$

$$\begin{aligned} W &= \lambda^{-1} I_Q \Phi^T (I_N + \Phi (\lambda^{-1} I_Q) \Phi^T)^{-1} t \\ &= \Phi^T (\lambda I_N + \Phi \Phi^T)^{-1} t \end{aligned}$$

Definimos la matriz Kernel (o Gram)  $K = \Phi \Phi^T$ . Los elementos de  $K$  son  $K_{ij} = \phi(x_i)^T \phi(x_j) = k(x_i, x_j)$ , donde  $k$  es la función Kernel, entonces:

$$W = \Phi^T (K + \lambda I_N)^{-1} t$$

Si definimos:

$$a = (K + \lambda I_N)^{-1} t$$

Entonces,

$$W = \Phi^T a$$

Si definimos  $a = (K + \lambda I_N)^{-1} t$ , entonces  $W = \Phi^T a$

Problema de optimización (en términos de  $\alpha$ ):

La función de costo original es:

$$J_{\text{ridge}}(W) = (\mathbf{t} - \Phi W)^T (\mathbf{t} - \Phi W) + \lambda W^T W$$

Sustituyendo  $W = \Phi^T \alpha$ :

$$\begin{aligned} J_{\text{KRR}} &= (\mathbf{t} - \Phi \Phi^T \alpha)^T (\mathbf{t} - \Phi \Phi^T \alpha) + \lambda (\Phi^T \alpha)^T (\Phi^T \alpha) \\ &= (\mathbf{t} - K \alpha)^T (\mathbf{t} - K \alpha) + \lambda \alpha^T \Phi \Phi^T \alpha \\ &= (\mathbf{t} - K \alpha)^T (\mathbf{t} - K \alpha) + \lambda \alpha^T K \alpha \\ &= \mathbf{t}^T \mathbf{t} - 2 \mathbf{t}^T K \alpha + \alpha^T K^T K \alpha + \lambda \alpha^T K \alpha \end{aligned}$$

Dado que  $K$  es simétrica ( $K^T = K$ ):

$$J_{\text{KRR}}(\alpha) = \mathbf{t}^T \mathbf{t} - 2 \mathbf{t}^T K \alpha + \alpha^T K K \alpha + \lambda \alpha^T K \alpha$$

Calculamos la derivada de  $J_{\text{KRR}}(\alpha)$  con respecto a  $\alpha$

$$\frac{\partial J_{\text{KRR}}(\alpha)}{\partial \alpha} = -2 K^T \mathbf{t} + 2 K^T K \alpha + 2 \lambda K^T \alpha \quad (\text{usando } K^T \text{ por formalidad, aunque } K = K^T)$$

$$\frac{\partial J_{\text{KRR}}(\alpha)}{\partial \alpha} = -2 K \mathbf{t} + 2 K K \alpha + 2 \lambda K \alpha$$

Iguamos la derivada a cero y resolvemos para  $\alpha$

$$-2 K \mathbf{t} + 2 K K \alpha + 2 \lambda K \alpha = 0$$

$$K K \alpha + \lambda K \alpha = K \mathbf{t}$$

$$K (K + \lambda I_N) \alpha = K \mathbf{t}$$

Obtenemos la solución para  $\alpha$ , si  $K$  es invertible (no siempre es el caso, pero  $(K + \lambda I_N)$  sí lo es para  $\lambda > 0$ ):

$$(K + \lambda I_N) \alpha = \mathbf{t}$$

$$\alpha_{\text{KRR}} = (K + \lambda I_N)^{-1} \mathbf{t}$$

La predicción para una nueva entrada  $x_*$  es:

$$\mu_* = \phi(x_*)^T w = \phi(x_*)^T \Phi^T \alpha = K(x_*)^T \alpha$$

donde  $K(x_*)$  es un vector de  $N \times 1$  con elementos:

$$K_n(x_*) = K(x_n, x_*).$$

## 7. Procesos Gaussianos (Gaussian Processes - GP)

Un proceso gaussiano es una colección de variables aleatorias, cualquier subconjunto finito de las cuales tiene una distribución Gaussiana conjunta. Un GP define una distribución sobre funciones  $f(x)$ .

Assumimos que:

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

Un GP se especifica por una función de media  $m(x)$  (a menudo  $m(x)=0$ ) y una función de covarianza (kernel)  $K(x, x')$ .

$$f(x) \sim \text{GP}(m(x), K(x, x'))$$

Las salidas observadas  $\mathbf{t} = [t_1, \dots, t_N]^T$ , dado que  $f(x) + \eta$  sigue una distribución Gaussiana:

$$\mathbf{t} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma_\eta^2 \mathbf{I}_N)$$

donde  $\mathbf{K}$  es la matriz kernel  $K_{ij} = K(x_i, x_j)$ . Sea  $\mathbf{C} = \mathbf{K} + \sigma_\eta^2 \mathbf{I}_N$ .

## Problema de optimización (Hiperparámetros)

La "optimización" en GPs se refiere a encontrar los hiperparámetros  $\theta$  del kernel  $K$  y la varianza del ruido  $\sigma_\eta^2$ . Esto se hace maximizando la log-likelihood marginal (o evidencia):

$$\ln p(\mathbf{t} | \mathbf{X}, \theta) = -\frac{1}{2} \mathbf{t}^T \mathbf{C}(\theta)^{-1} \mathbf{t} - \frac{1}{2} \ln |\mathbf{C}(\theta)| - \frac{N}{2} \ln(2\pi)$$



Se usan métodos de optimización basados en gradiente para encontrar  $\theta$  que maximice esta expresión.

La derivada con respecto a un hiperparámetro  $\theta_j$  es:

$$\frac{\partial \ln p(\mathbf{t}|\theta)}{\partial \theta_j} = \frac{1}{2} \mathbf{t}^T \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_j} \mathbf{C}^{-1} \mathbf{t} - \frac{1}{2} \text{tr}(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_j})$$

Si definimos  $\alpha_{GP} = \mathbf{C}^{-1} \mathbf{t}$ :

$$\frac{\partial \ln p(\mathbf{t}|\theta)}{\partial \theta_j} = \frac{1}{2} \text{tr}((\alpha_{GP} \alpha_{GP}^T - \mathbf{C}^{-1}) \frac{\partial \mathbf{C}}{\partial \theta_j})$$

Predicción con GPs:

Para predecir el valor de la función latente  $f_* = F(\mathbf{x}_*)$  en un nuevo punto  $\mathbf{x}_*$ , consideramos la distribución conjunta  $\mathbf{t}$  y  $f_*$ :

$$\begin{pmatrix} \mathbf{t} \\ f_* \end{pmatrix} \sim \mathcal{N}\left(0, \begin{pmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I}_N & \mathbf{K}_* \\ \mathbf{K}_*^T & K_{**} \end{pmatrix}\right) = \mathcal{N}\left(0, \begin{pmatrix} \mathbf{C} & \mathbf{K}_* \\ \mathbf{K}_*^T & K_{**} \end{pmatrix}\right)$$

donde  $\mathbf{K}_*$  es un vector de  $N \times 1$  con elementos  $K(\mathbf{x}_n, \mathbf{x}_*)$ , y  $K_{**} = K(\mathbf{x}_*, \mathbf{x}_*)$ . La distribución predictiva condicionada  $p(f_* | \mathbf{t}, \mathbf{x}_*)$  es gaussiana con:

• Media:  $\bar{f}_* = E[f_* | \mathbf{t}] = \mathbf{K}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{t}$

• Varianza:  $\text{Var}(f_*) = K_{**} - \mathbf{K}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_* = K_{**} - \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{K}_*$   
 se se predice

$$t_* = f_* + \eta_*, \text{ la media es la misma, } \bar{t}_* = \bar{f}_*$$

y la varianza es

$$\text{Var}(t_*) = \text{Var}(f_*) + \sigma_n^2$$

## Discusiones

- OLS, MLE, Ridge, MAP: Producen una estimación puntual para los pesos  $W$ .
- El modelo Bayesiano Lineal, GP: Producen una distribución de probabilidad completa sobre los pesos o directamente sobre las funciones. Permitiendo cuantificar la incertidumbre.
- OLS, MLE: No incorporan regularización explícita.
- Ridge, MAP: Introducen regularización  $\Leftrightarrow$  lo una prior Gaussiana, lo que ayuda a prevenir el sobreajuste y estabilizar la solución.
- Modelo Bayesiano: en este modelo la prior actúa como forma de regularización.
- GP: la elección del kernel y sus hiperparámetros, junto con el término de ruido, inherentemente regularizan el modelo y controlan su suavidad y complejidad.