

# Introducción

Queremos predecir un valor real  $t_n \in \mathbb{R}$  asociado a una entrada  $X_n \in \mathbb{R}^p$  usando un modelo

lineal definido como

$$t_n = \Phi(X_n) W^T + \eta_n$$

*Nota:*  $\Phi(X_n) \in \mathbb{R}^{1 \times Q}$  (vector fila)  
 $W \in \mathbb{R}^{1 \times Q} \rightarrow W^T \in \mathbb{R}^{Q \times 1}$  (vector columna)

Luego  
 $\Phi(X_n) \cdot W^T \in \mathbb{R}^{1 \times Q} \cdot \mathbb{R}^{Q \times 1} = \mathbb{R}^{1 \times 1}$

donde

- $\Phi: \mathbb{R}^p \rightarrow \mathbb{R}^Q$  es una función base
- $W \in \mathbb{R}^Q$  son los parámetros del modelo
- $\eta_n \sim \mathcal{N}(0, \sigma_n^2)$  representa ruido gaussiano.

Asumimos que los datos  $\{(X_n, t_n)\}_{n=1}^N$  son independientes e idénticamente distribuidos.

## Mínimos cuadrados

Buscamos estimar el vector de parámetros  $w$  que minimiza el error cuadrático entre las predicciones del modelo y los valores observados  $t_n$ .

## Función objetivo.

La función a minimizar es el error cuadrático medio total

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

Definimos

$$\Phi \in \mathbb{R}^{N \times Q} = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_n) \end{bmatrix}$$

$$t \in \mathbb{R}^N = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$$

Re escribiendo

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

Luego queremos encontrar

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

Derivando e igualando a cero obtenemos

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

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

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

igualando a cero

$$\frac{\partial J}{\partial w} = -\Phi^T (t - \Phi w^T) = 0$$

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

Despejando  $w$  (si  $\Phi^T \Phi \in \mathbb{R}^{Q \times Q}$  es invertible)

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T t$$

Esta solución proporciona el estimador lineal de mínimos cuadrados, que minimiza el error cuadrático sobre los datos observados. Es la base de modelos más avanzados como la máxima verosimilitud, la regresión regularizada o el enfoque bayesiano. Su simplicidad permite obtener una solución cerrada, pero es sensible al sobreajuste y a problemas de acondicionamiento si  $\Phi^T \Phi$  no es invertible o está mal condicionada.

## Mínimos cuadrados regularizados (Regresión Ridge)

### Función objetivo

Se busca minimizar el error cuadrático más una penalización sobre la norma de los parámetros.

$$J(w) = \frac{1}{N} \sum_{i=1}^N (y_i - \phi_i^T w)^2 + \lambda \|w\|^2$$

$$L(w) = \frac{1}{2} \sum_{n=1}^N (t_n - \Phi(x_n)^T w)^2 + \frac{\lambda}{2} \|w\|_2^2$$

En forma matricial:

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

donde:

- $\Phi \in \mathbb{R}^{N \times Q}$  es la matriz de diseño con  $\Phi(x_n)^T$  como filas.
- $t \in \mathbb{R}^N$  es el vector de salidas.
- $\lambda \geq 0$  es el hiperparámetro de regularización

Luego el estimado de los parámetros  $w$  es

$$w^* = \underset{w}{\operatorname{argmin}} \left( \frac{1}{2} \|t - \Phi w\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \right)$$

derivando e igualando a cero

$$\frac{\partial L}{\partial w} = -\Phi^T (t - \Phi w) + \lambda w = 0$$

$$\rightarrow (\Phi^T \Phi + \lambda I) = \Phi^T \epsilon$$

de donde

$$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \epsilon$$

La regresión Ridge agrega una penalización cuadrática sobre los parámetros al error de mínimos cuadrados. Esto mejora su estabilidad numérica cuando  $\Phi^T \Phi$  está mal condicionada o no invertible, y reduce el riesgo de sobreajuste. Cuando  $\lambda=0$  se recupera la solución de mínimos cuadrados.

## Máxima verosimilitud

Bajo el supuesto de ruido gaussiano, la salida  $t_n$  condicionada a la entrada  $x_n$  tiene distribución:

$$p(t_n | x_n, w) = \mathcal{N}(t_n | \phi(x_n)^T w, \sigma^2)$$

La verosimilitud conjunta de los datos

es

$$p(t | \Phi, w) = \frac{1}{N} \prod_{n=1}^N \mathcal{N}(t_n | \Phi(x_n)^T w, \sigma^2)$$

luego, el estimador de máxima verosimilitud  
es

$$w^* = \underset{w}{\operatorname{argmax}} p(t | \Phi, w)$$

Maximizar la verosimilitud es equivalente a minimizar la log-verosimilitud negativa. Tomando logaritmo y descartando términos constantes:

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

Como  $\sigma^2$  es constante, minimizar  $L$  equivale a minimizar el error cuadrático ordinario:

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{2} \|t - \Phi w\|_2^2$$

Derivando e igualando a cero

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

→

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T t$$

La estimación por máxima verosimilitud produce el mismo estimador que mínimos cuadrados ordinarios, pero con una motivación probabilística: se busca el valor  $w$  que **maximiza la probabilidad de observar los datos** bajo el supuesto de ruido gaussiano.

## Maximo a posteriori

Assumimos el mismo modelo lineal con ruido gaussiano:

$$t_n = \phi(x_n)^T w + \eta_n, \eta_n \sim \mathcal{N}(0, \sigma^2)$$

Además, incorporamos una **distribución a priori** sobre los parámetros

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

**Función objetivo**



Se busca el valor de  $w$  que maximiza la distribución posterior

$$w^* = \underset{w}{\operatorname{argmax}} p(w | t, \Phi)$$

donde

$$p(w | t, \Phi) = \frac{p(t | \Phi, w) \cdot p(w)}{p(t | \Phi)}$$

Como el denominador no depende de  $w$ , maximizar la posterior es equivalente a maximizar el producto de la verosimilitud y el prior:

$$w^* = \underset{w}{\operatorname{argmax}} [p(t | \Phi, w) \cdot p(w)]$$

Tomando logaritmos y descartando constantes, obtenemos:

$$L(w) = \frac{1}{2\sigma^2} \|t - \Phi w\|_2^2 + \frac{1}{2\tau^2} \|w\|_2^2$$

Luego

$$w^* = \arg \min_w \left( \frac{1}{2} \|t - \Phi w\|_2^2 + \frac{\sigma^2}{2\tau^2} \|w\|_2^2 \right)$$

Entonces  $\lambda = \frac{\sigma^2}{\tau^2} :$

$$w^* = \left( \Phi^T \Phi + \lambda I \right)^{-1} \Phi^T t$$

La estimación MAP incorpora conocimiento previo sobre los parámetros, asumiendo que  $w$  sigue una distribución gaussiana centrada en cero. Esto conduce al mismo estimador que regresión Ridge, pero con una interpretación Bayesiana: la regularización aparece como consecuencia natural del prior. En el caso límite  $\tau^2 \rightarrow \infty$  se recupera la estimación de máxima verosimilitud.

## Modelo Bayesiano Completo (Regresión lineal Bayesiana con Prior Gaussiana)

### Función objetivo

A diferencia de MAP o ML, aquí no buscamos un único valor  $w$ , sino que

inferimos la distribución posterior completa:

$$p(w | t, \Phi) = \frac{p(t | \Phi, w) \cdot p(w)}{p(t, \Phi)}$$

Tanto el prior como la verosimilitud son distribuciones normales, por lo que la posterior también es gaussiana.

Sea forma es

$$p(w | t, \Phi) = \mathcal{N}(w | \mu_N, \Sigma_N)$$

con

$$\Sigma_N = \left( \frac{1}{\sigma^2} \Phi^T \Phi + \frac{1}{\tau^2} I \right)^{-1} \quad y$$

$$\mu_N = \frac{1}{\sigma^2} \Sigma_N \Phi^T t$$

Predicción.

Dada una nueva entrada  $x_*$ , se transforma como  $\Phi_* = \Phi(x_*)$ . La salida  $t_*$  también es una variable aleatoria, cuya distribución predictiva es gaussiana:

$$p(t_* | x_*, t, \Phi) = \mathcal{N}(t_* | \Phi_*^T \mu_N, \Phi_*^T \Sigma_N \Phi_* + \sigma^2)$$

Esto nos da

- Una media predictiva  $E[t_*] = \Phi_*^T \mu_N$
- Una varianza predictiva  $V[t_*] = \Phi_*^T \Sigma_N \Phi_* + \sigma^2$

Este enfoque no entrega un único vector  $w$ , sino una distribución posterior completa que refleja la incertidumbre sobre los parámetros dado el conjunto de datos observado. La predicción sobre nuevos puntos también es una **distribución**, no un valor fijo, lo que permite cuantificar la incertidumbre del modelo.

El modelo Bayesiano completo puede verse como una generalización natural de MAP, donde en lugar de maximizar la posterior se la integra completamente. Es también el antecedente directo de los procesos gaussianos, en los cuales se coloca una distribución directamente sobre funciones en vez de sobre parámetros  $w$ .

## Regresión Ridge con Kernel

Partimos del mismo supuesto funcional que en regresión lineal regularizada:

$$\epsilon_n = \Phi(X_n)^T w + \eta_n$$

Pero en este caso:

- El espacio de características  $\Phi(\cdot)$ , puede ser de muy alta (incluso infinita) dimensión.
- En lugar de trabajar explícitamente con  $\Phi$ , se utiliza un  $\text{kernel}(X, x')$  definido como:

$$K(X, x') = \Phi(X)^T \Phi(x')$$

Función objetivo.

Como en Ridge, se busca minimizar

$$\mathcal{L}(w) = \frac{1}{2} \|t - \Phi\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

En lugar de resolver directamente para  $w$ , se utiliza la forma dual, que parte del hecho que la solución  $w$  puede expresarse como una combinación lineal de los vectores de entrenamiento:

$$w = \tilde{\Phi} \alpha$$

Este cambio permite reformular el problema en términos de un nuevo conjunto de parámetros  $\alpha \in \mathbb{R}^N$ , eliminando la necesidad de calcular  $\Phi(x)$  explícitamente.

Todo se puede expresar en función

de productos punto entre vectores transformados, que se sustituyen por evaluaciones del Kernel.

Reemplazando en la función objetivo y derivando respecto a  $\alpha$ , se obtiene el sistema dual:

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

donde  $K \in \mathbb{R}^{N \times N}$  es la matriz Kernel definida como:

$$K_{ij} = K(x_i, x_j)$$

## Predicción

Dado un nuevo punto  $x_*$ , su predicción se calcula como:

$$f_* = \sum_{n=1}^N \alpha_n K(x_n, x_*)$$

Es decir, la predicción se expresa como una combinación ponderada de las

Similitudes entre el nuevo punto y los puntos de entrenamiento.

Kernel Ridge Regression permite trabajar con espacios de características complejos o de dimensión infinita sin calcular explícitamente la transformación  $\Phi(y)$  gracias al truco del kernel. Este modelo sigue siendo lineal en un espacio de características (implícito), pero es capaz de capturar relaciones altamente no lineales en el espacio original. Su desventaja principal es que requiere almacenar y operar sobre la matriz kernel  $K \in \mathbb{R}^{N \times N}$  lo que limita su escalabilidad. Cuando se utiliza un kernel lineal  $K(x, x') = x^T x'$  se recupera el modelo Ridge Clásico.

## Procesos Gaussianos

Un proceso gaussiano (GP) es una distribución sobre funciones

$$f(\cdot) \sim GP(m(\cdot), K(\cdot, \cdot))$$

donde

- $m(x) = \mathbb{E}[f(x)]$  es la función media (usualmente se toma como cero)



$$\bullet K(x, x') = E[(f(x) - m(x))(f(x') - m(x')))]$$

Es una función de covarianza o kernel.

En regresión, se asume que los datos observados provienen de esta función con ruido gaussiano aditivo.

$$t_n = f(x_n) + \eta_n, \quad \eta_n \sim \mathcal{N}(0, \sigma^2)$$

## Supuestos

- El vector de salidas  $t$  está gobernado por una función  $f \sim GP(0, K)$
- Ruido gaussiano con varianza  $\sigma^2$ ,
- Kernel  $K(x, x')$ , simétrico y positivo.

## Inferencia y predicción

Dado el conjunto de entrenamiento

$$\mathcal{D} = \{(x_n, t_n)\}_{n=1}^N, \text{ se define la}$$

matriz kernel:

$$K \in \mathbb{R}^{N \times N}, \quad K_{ij} = K(x_i, x_j)$$

Para un nuevo punto  $x_*$ , se define

$$\bullet K_* = [K(x_1, x_*), \dots, K(x_N, x_*)]^T$$

$$\bullet K_{**} = K(x_*, x_*)$$

La distribución predictiva para  $t_*$  es gaussiana:

$$p(t_*, x_*, t) = \mathcal{N}(u_*, \sigma_*^2)$$

con

$$u_* = K_*^T (K + \sigma^2 I)^{-1} t$$

$$\sigma_*^2 = K_{**} - K_*^T (K + \sigma^2 I)^{-1} K_*$$

Los procesos gaussianos generalizan por completo el enfoque bayesiano. En lugar de asumir una forma paramétrica para la función (como una combinación de  $\phi(x)^T w$ ) colocan una distribución directamente sobre funciones, permitiendo capturar

relaciones altamente no lineales con control explícito de la incertidumbre.

Este modelo no calcula un vector de parámetros  $w$ , ni depende de una transformación explícita  $\Phi(w)$  todo se hace a través del kernel. De hecho, si se usa un kernel tipo  $k(x, x') = \Phi(x)^\top \Phi(x')$  se recuperan los mismos resultados predictivos que el modelo bayesiano lineal con prior gaussiano sobre  $w$ .

La principal limitación de los GPs es computacional: su cost es cúbico en el número de datos  $\mathcal{O}(N^3)$  lo cual restringe su uso a conjuntos relativamente pequeños.