

# Teoría de Aprendizaje de máquina / Juan Jerónimo Castaño Rivera

1) Sea el modelo

$$t_n = \phi(x_n) w^T + \eta_n, \quad t_n \in \mathbb{R}, x_n \in \mathbb{R}^p$$
$$w \in \mathbb{R}^p; \quad \phi \in \mathbb{R}^p \rightarrow \mathbb{R}^p$$
$$\eta_n \sim N(\eta_n | 0, \sigma_n^2)$$

De forma matricial  $t = \Phi w + \eta$

Donde

$\Phi \in \mathbb{R}^{N \times p}$  = Datos

$t \in \mathbb{R}^N$  = Vector de salida

$w \in \mathbb{R}^p$  = Modelo

- Mínimos Cuadrados

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

El valor mínimo se puede hallar derivando e igualando a cero

$$\|t - \Phi w\|_2^2 = \langle t - \Phi w, t - \Phi w \rangle = (t - \Phi w)^T (t - \Phi w)$$

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

$$\frac{\partial}{\partial w} \left\{ \frac{1}{2} (t^T t - 2t^T \Phi w + (\Phi w)^T (\Phi w)) \right\} = 0$$

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

$$-2\Phi^T t + 2\Phi^T \Phi w = 0 \quad \Rightarrow \quad 2\Phi^T t = 2\Phi^T \Phi w$$

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

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



- Mínimos Cuadrados regularizados

$$W = \underset{W}{\operatorname{argmin}} \|t - \Phi W\|_2^2 + \lambda \|W\|_2^2$$

$$W^* = \frac{\partial}{\partial W} \{ \|t - \Phi W\|_2^2 + \lambda \|W\|_2^2 \} = 0$$

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

$$(2\Phi^T \Phi + 2\lambda I)W = 2\Phi^T t$$

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

- Máxima verosimilitud

$$t_n = \phi(x_n)^T W + \eta_n \quad ; \quad \eta \sim N(0, \sigma_n^2)$$

Se puede ver que  $p(t_n | x_n, W) = N(t_n | \phi(x_n)^T W, \sigma_n^2)$

Dado que son iid  $\prod_{n=1}^N p(t_n | \phi(x_n)W, \sigma_n^2)$

$$\log \left( \prod_{n=1}^N p(t_n | \phi_n W, \sigma_n^2) \right) = \log \left( \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp \left( -\frac{|t_n - \phi_n W|^2}{2\sigma_n^2} \right) \right)$$

$$= N \log \left( \frac{1}{\sqrt{2\pi\sigma_n^2}} \right) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N |t_n - \phi_n W|^2$$

$$= -\frac{N}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N |t_n - \phi_n W|^2$$

Se deriva e iguala a cero, el problema corresponde a mínimos cuadrados

$$W^* = \underset{W}{\operatorname{argmin}} -\frac{1}{2\sigma_n^2} \|t - \Phi W\|_2^2 - \frac{N}{2} \log(2\pi\sigma_n^2)$$

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



- Máximo a-posteriori

En lugar de encontrar el valor de  $w$  que maximiza la verosimilitud, se asume un prior sobre  $w$  y buscamos el valor que maximiza la posterior

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

$$p(t, w) = p(w, t) \Rightarrow p(t|w)p(w) = p(w|t)p(t)$$

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

Para simplificar  $p(t) = 1$

Al igual que en máxima verosimilitud

$$p(t|w) = \prod_{n=1}^N \mathcal{N}(t_n | \phi(x_n)w, \sigma_n^2) \quad p(w) = \mathcal{N}(w | 0, \sigma_w^2)$$

$$= \log \left( \prod_{n=1}^N p(t_n | \phi_n w, \sigma_n^2) \prod_{q=1}^Q p(w_q | 0, \sigma_w^2) \right) = \prod_{n=1}^Q \mathcal{N}(w_q | 0, \sigma_w^2)$$

$$= \log \left( \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp \left( -\frac{|t_n - \phi_n w|^2}{2\sigma_n^2} \right) \right) + \log \left( \prod_{q=1}^Q \frac{1}{\sqrt{2\pi\sigma_w^2}} \exp \left( -\frac{|w_q|^2}{2\sigma_w^2} \right) \right)$$

$$= \log \left( \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \right) + \log \left( \prod_{n=1}^N \exp \left( -\frac{|t_n - \phi_n w|^2}{2\sigma_n^2} \right) \right) + \log \left( \prod_{q=1}^Q \frac{1}{\sqrt{2\pi\sigma_w^2}} \right)$$

$$+ \log \left( \prod_{q=1}^Q \exp \left( -\frac{|w_q|^2}{2\sigma_w^2} \right) \right)$$

$$= \underbrace{-\frac{N}{2} \log(2\pi\sigma_n^2)}_{\text{No aporta}} - \underbrace{\frac{Q}{2} \log(2\pi\sigma_w^2)}_{\text{No aporta}} - \frac{1}{2\sigma_n^2} \sum_{n=1}^N |t_n - \phi_n w|^2 - \frac{1}{2\sigma_w^2} \sum_{q=1}^Q |w_q|^2$$

$$w^* = \min_w \left[ \|t - \Phi w\|^2 + \frac{2\sigma_n^2}{2\sigma_w^2} \|w\|^2 \right] \quad \text{Min. cuad. regularizados}$$

$$w^* = \left( \Phi^T \Phi + \frac{\sigma_n^2}{\sigma_w^2} I \right)^{-1} \Phi^T t$$



## - Bayesiano con modelo lineal Gaussiano

Este modelo no solo busca un único vector  $w$ , sino la distribución completa posterior sobre  $w$  dados los datos observados

$$p(w|t) = N(w|m_N, \Sigma_N)$$

$\Sigma_N$ : Matriz covarianza

$$\Sigma_N = \left( \frac{1}{\sigma_w^2} I + \frac{1}{\sigma_n^2} \Phi^T \Phi \right)$$

$$m_N = \frac{1}{\sigma_n^2} \Sigma_N \Phi^T t$$

Optimización: Maximizar la evidencia

$$p(t) = \int p(t|w) p(w) dw$$

Demostración:

Sea un vector  $x \in \mathbb{R}^q$  con prior Gaussiano

$$p(x) = N(x|\mu, \Delta^{-1})$$

Sea el modelo  $y = Ax + b$ ;  $p(y|x) = N(y|Ax+b, L^{-1})$

$$p(x, y) = p(x) p(y|x) \Rightarrow p(x|y) = N(x|\mu_{x|y}, \Sigma_{x|y})$$

$$\Sigma = (\Delta + A^T L A)^{-1}; \quad \mu_{x|y} = \Sigma_{x|y} (A^T L (y-b) + \Delta \mu)$$

Para nuestro caso

$$\epsilon_n = \phi(x_n) w^T + \eta_n \quad t = \phi w + \eta; \quad p(t|w) = N(t|\phi w, \sigma_n^2)$$

Se asume  $p(w) = N(w|m_0, S_0)$

$$\text{Sabemos que } \log p(t|w) = -\frac{1}{2\sigma_n^2} \|t - \phi w\|_2^2 + \text{cte.}$$

$$\log(N(x|\mu, \Sigma)) = -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) - \frac{1}{2} \log |\Sigma| - \frac{n}{2} \log 2\pi$$

$$\text{Por tanto } \log(p(w)) = -\frac{1}{2} (w-m_0)^T S_0^{-1} (w-m_0) + \text{cte}$$

$$\log(p(t|w) p(w)) = \log(p(t|w)) + \log(p(w))$$



$$= -\frac{1}{2} \left[ \frac{1}{\sigma_n^2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) + (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right] + cte$$

Agrupando términos

$$\mathbf{w}^T \left( \frac{1}{\sigma_n^2} \Phi^T \Phi + \mathbf{S}_0^{-1} \right) \mathbf{w} - \mathbf{w}^T \left( \frac{1}{\sigma_n^2} \Phi^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{m}_0 \right)$$

Se sabe que  $\log(p(\mathbf{w})) = -\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)$

$$= -\frac{1}{2} \left[ \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{w} - 2 \mathbf{w}^T \mathbf{S}_0^{-1} \mathbf{m}_0 + \mathbf{m}_0^T \mathbf{S}_0^{-1} \mathbf{m}_0 \right] + cte$$

$$\mathbf{S}_N = \left( \frac{1}{\sigma_n^2} \Phi^T \Phi + \mathbf{S}_0^{-1} \right)^{-1} \quad \mathbf{S}_N^{-1} \mathbf{m}_N = \left( \frac{1}{\sigma_n^2} \Phi^T \mathbf{t} + \mathbf{S}_0^{-1} \mathbf{m}_0 \right)$$

$$\mathbf{m}_N = \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \frac{1}{\sigma_n^2} \Phi^T \mathbf{t} \right)$$

como  $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) = N(\mathbf{w} | \mathbf{0}, \sigma_w^2)$

$$\mathbf{S}_N = \left( \frac{1}{\sigma_w^2} \mathbf{I} + \frac{1}{\sigma_n^2} \Phi^T \Phi \right)^{-1} \quad \mathbf{m}_N = \frac{1}{\sigma_n^2} \mathbf{S}_N \Phi^T \mathbf{t}$$

$$\mathbf{m}_N = \left( \frac{1}{\sigma_n^2} \right) \left( \frac{1}{\sigma_w^2} \right)^{-1} \left( \frac{\sigma_n^2}{\sigma_w^2} \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t} = \left( \frac{\sigma_n^2}{\sigma_w^2} \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}$$

- Regresión rígida Kernel

Se extiende a espacios de funciones no lineales

$$\hat{y} = \Phi \mathbf{w} \quad ; \quad \Phi \in \mathbb{R}^{N \times Q} \quad ; \quad Q \rightarrow \infty \quad (\text{RKHS})$$

$$\hat{y} : \mathbb{R}^Q \rightarrow \mathbb{R}$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{N} \|\mathbf{t} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Por mínimos cuadrados  $\mathbf{w}^* = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t}$

Sin embargo  $\Phi^T \Phi \in \mathbb{R}^{Q \times Q} \quad ; \quad Q \rightarrow \infty$

$$(\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T = \left( \lambda \left( \frac{1}{\lambda} \Phi^T \Phi + \mathbf{I} \right) \right)^{-1} \Phi^T = \frac{1}{\lambda} \left( \mathbf{I} + \Phi^T \frac{\Phi}{\lambda} \right)^{-1} \Phi^T$$



$$\Phi^T \frac{1}{\lambda} (\mathbf{I} + \frac{1}{\lambda} \Phi \Phi^T)^{-1} = \Phi^T (\lambda \mathbf{I} + \Phi \Phi^T)^{-1}$$

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

Se hace predicción para nuevos puntos  $f(x_*) = \phi(x_*)^T \mathbf{W}$

$$f(x_*) = \phi(x_*)^T \Phi^T (\lambda \mathbf{I} + \Phi \Phi^T)^{-1} \mathbf{t}$$

$$\mathbf{K} = \Phi \Phi^T ; \quad \mathbf{K}_* = \phi(x_*)^T \Phi^T$$

$$f(x_*) = [\mathbf{K}_*(\cdot)^T] (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t}$$

### - Gaussian Process

Extienden el método paramétrico para definir la incertidumbre de los parámetros del regresor al imponer un prior sobre funciones directamente en  $\mathbb{R}^{K+1}S$

El GP se define por su media y covarianza

$$f(x) \in \mathbb{R} ; f(x) = \phi(x)^T \mathbf{w} ; p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_{\mathbf{w}}) ; \Sigma_{\mathbf{w}} \in \mathbb{R}^{Q \times Q}$$

$$m(x) = E\{f(x)\} = E\{\phi(x)^T \mathbf{w}\} = \phi(x)^T E\{\mathbf{w}\} = 0$$

$$\text{cov}(f(x), f(x')) = \kappa(x, x') = \phi(x)^T \Sigma_{\mathbf{w}} \phi(x')$$

$$f \sim \mathcal{GP}(f | \mathbf{0}, \mathbf{K}) ; \quad \mathbf{K} = [\kappa(x, x')] \in \mathbb{R}^{N \times N}$$

Este modelo busca maximizar la verosimilitud  $p(\mathbf{t} | \theta)$ , siendo  $\theta$  un hiperparámetro de  $\kappa(\cdot, \cdot | \theta)$

Se puede obtener que

$$m(x_*) = \mathbf{K}_*^T (\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} \mathbf{t} , \quad \mathbf{K}_* = [\kappa(x_*, x_1), \kappa(x_*, x_2), \dots, \kappa(x_*, x_N)]^T$$

$$\text{cov}(f(x), f(x')) = \kappa(x_*, x_*) + \sigma_e^2 - \mathbf{K}_*^T (\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} \mathbf{K}_*$$