

Minimos Cuadrados

$$t_n = \varphi(x_n)^T w + b_n$$

$$w = \arg \min J(w)$$

↳ Función de Costo

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

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

Nota: los factores $\frac{1}{2}$ son para simplificar la derivada

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

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

$$\nabla_w J(w) = 0 - t^T \phi + \phi^T \phi w$$

- Igualamos a Cero

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

$$\phi^T \phi w = t^T \phi$$

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

$$\hat{w} = (\phi^T \phi)^{-1} t^T \phi //$$

Minimos Cuadrados Regularizados

$$t_n = \varphi(X_n)^T w + \eta_n \quad \eta_n \sim N(0, \sigma^2)$$

$$\hat{w}_{ridge} = \arg \min J(w)$$

↳ Función de Costo

- Agregamos la norma L2 para penalizar pesos grandes

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

Nota: los factores $\frac{1}{2}$ son para simplificar la derivada

$$\hat{w}_{ridge} = \arg \min \left(\frac{1}{2} \|t - \phi w\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \right)$$

$$\|t - \phi w\|_2^2 = (t - \phi w)(t - \phi w)^T$$

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

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

$$\nabla_w J(w) = 0 - t^T \phi + \phi^T \phi w + \lambda w$$

$$-t^T \phi + \phi^T \phi \omega + \lambda \omega = 0$$

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

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

$$\omega = (\phi^T \phi + \lambda I)^{-1} t^T \phi$$

$$\hat{\omega}_{ridge} = (\phi^T \phi + \lambda I)^{-1} t^T \phi$$

Maxima Verosimilitud

$$t_n = \varphi(x_n)^T w + \eta_n \quad \eta_n \sim N(0, \sigma_\eta^2) \quad | \quad t_n$$

- Sigue una distribución normal.

$$p(t_n | x_n, w, \sigma_\eta^2) = N(t_n | \varphi(x_n)^T w, \sigma_\eta^2)$$

- Asumiendo i.i.d. la verosimilitud conjunta, es el producto de todas las probabilidades individuales.

$$p(t | X, w, \sigma_\eta^2) = \prod_{n=1}^N N(t_n | \varphi(x_n)^T w, \sigma_\eta^2)$$

- Forma de una densidad gaussiana unidimensional.

$$N(t_n | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(t_n - \mu^2)}$$

$$p(t | X, w, \sigma_\eta^2) = \prod_{n=1}^N p(t_n | x_n, w, \sigma_\eta^2)$$

$$p(t | X, w, \sigma_\eta^2) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_\eta^2}} e^{-\frac{1}{2\sigma_\eta^2}(t_n - \varphi(x_n)^T w)^2}$$

- El producto de las constantes $(2\pi\sigma_\eta^2)^{-1/2}$ se convierte en:

$$\prod_{n=1}^N (2\pi\sigma_\eta^2)^{-1/2} = (2\pi\sigma_\eta^2)^{-N/2}$$

$$p(t | X, w, \sigma_\eta^2) = (2\pi\sigma_\eta^2)^{-N/2} \prod_{n=1}^N e^{-\frac{1}{2\sigma_\eta^2}(t_n - \varphi(x_n)^T w)^2}$$

$$\prod_n e^{a_n} = e^{\sum_n a_n}$$

$$p(+|X, w, \sigma_v^2) = (2\pi\sigma_v^2)^{-N/2} e^{\left(-\frac{1}{2\sigma_v^2} \sum_{n=1}^N (t_n - \varphi(X_n)^T w)^2\right)}$$

$$\sum_{n=1}^N (t_n - \varphi(X_n)^T w)^2 = (t - \phi w)^T (t - \phi w)$$

$$p(+|X, w, \sigma_v^2) = (2\pi\sigma_v^2)^{-N/2} e^{\left(-\frac{1}{2\sigma_v^2} \|t - \phi w\|^2\right)}$$

- Es mas conveniente trabajar con el logaritmo de la verosimilitud

$$\ln p(+|X, w, \sigma_v^2) = -\frac{N}{2} \ln(2\pi\sigma_v^2)$$

$$- \frac{1}{2\sigma_v^2} \|t - \phi w\|^2$$

$$\hat{w}_{ML} = \operatorname{argmax}_w \ln p(+|X, w, \sigma_v^2)$$

→ Maxima la lg-verosimilitud

- Como el primer termino no depende de w , es equivalente a minimizar el termino cuadrático

$$\hat{w}_{ML} = \operatorname{argmin}_w \|t - \phi w\|^2$$

$$\hat{w}_{ML} = (\phi^T \phi)^{-1} \phi^T t$$

$$\sigma_v^2 = \frac{1}{N} \|t - \phi \hat{w}_{ML}\|^2$$

Maximo a posteriori $\hat{w} = \arg \max_w p(t|X, w)$

$$t_n = \varphi(X_n)^T w + \eta_n, \quad \eta_n \sim N(0, \sigma_\eta^2)$$

- En base al MLE (Maxima Verosimilitud),

$$\hat{w}_{MLE} = \arg \max_w p(t|X, w)$$

$$p(t|X, w) = (2\pi\sigma_\eta^2)^{-N/2} e^{-\frac{1}{2\sigma_\eta^2} \|t - \Phi w\|^2}$$

$$\hat{w}_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T t = g_{MLE}$$

- Introducimos un prior sobre los pesos w ,
asumiendo que los pesos siguen una distribución
normal centrada en 0.

$$p(w) = N(w|0, \sigma_w^2 \Sigma)$$

$$p(w) = (2\pi\sigma_w^2)^{-N/2} e^{-\frac{1}{2\sigma_w^2} \|w\|^2}$$

- Se aplica el teorema de Bayes.

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

$$w_{map} = \arg \max_w \log p(t|X, w) + \log p(w)$$

$$\log p(t|X, w) = -\frac{N}{2} \log(2\pi\sigma_\eta^2) - \frac{1}{2\sigma_\eta^2} \|t - \Phi w\|^2$$

$$\log p(w) = -\frac{N}{2} \log(2\pi\sigma_w^2) - \frac{1}{2\sigma_w^2} \|w\|^2$$

$$\log p(w|t, X) = C - \frac{1}{2\sigma_\eta^2} \|t - \Phi w\|^2 - \frac{1}{2\sigma_w^2} \|w\|^2$$

$$J(\omega) = \frac{1}{2\sigma_y^2} \|t - \phi\omega\|^2 + \frac{1}{2\sigma_\omega^2} \|\omega\|^2$$

$$\lambda = \frac{\sigma_y^2}{\sigma_\omega^2}$$

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

$$\nabla_\omega J(\omega) = -\phi^T t + \phi^T \phi \omega + \lambda \omega = 0$$

$$\omega_{\text{map}} = (\phi^T \phi + \lambda I)^{-1} \phi^T t$$

↳ Mismo resultado que

con mínimos cuadrados
regulizados

Primer parcial TAM

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

Bayesiano con modelo lineal gaussiano:

En lugar de buscar un solo número w , el modelo devuelve una distribución completa sobre w , que refleja incertidumbre.

Expresión del posterior

$P(w|t, X) = \mathcal{N}(w|w_{MAP}, \Sigma_{MAP}) \rightarrow$ la distribución de w después de ver los datos es una normal multivariada con media y covarianza

Media:

$$w_{MAP} = \Sigma_0^{-1} m_0 + \frac{1}{\sigma_n^2} \Phi^T t$$

Covarianza:

$$\Sigma_{MAP} = \left(\Sigma_0^{-1} + \frac{1}{\sigma_n^2} \Phi^T \Phi \right)^{-1}$$

Ai tener más datos, más pequeño se vuelve Σ_n , o sea, el modelo se vuelve más seguro sobre los pesos w .

Si asumimos $m_0 = 0 \rightarrow w_{MAP} = \frac{1}{\sigma_n^2} \Sigma_n \Phi^T t$

Prior sobre pesos:

$$P(w) = \mathcal{N}(w|m_0, \Sigma_0)$$

Verosimilitud:

$$P(t|w) = \mathcal{N}(t|\Phi w, \sigma_n^2 I)$$

$$P(w|t_0, x_0) = P(t|w) \cdot P(w)$$

$$P(w|t_0, x_0) = -\frac{1}{2\sigma_n^2} \left(\|t_0\|_2^2 - 2t_0 \phi(x_0)^T w + \phi(x_0)^T w \phi(x_0) w^T \right)$$

$$= -\frac{1}{2} \Sigma_0^{-1} (w - m_0)^T (w - m_0)$$

$$= -\frac{1}{2} \left(w^T \Sigma_0^{-1} w - (w^T)^T \Sigma_0^{-1} (m_0)^T - m_0^T \Sigma_0^{-1} w + \Sigma_0^{-1} m_0^T m_0 \right)$$

$$= -\frac{1}{2} \left(\underbrace{\Sigma_0^{-1} w^T w}_{\text{cuadrático}} - \underbrace{2 \Sigma_0^{-1} m_0^T w}_{\text{lineal}} + \underbrace{\Sigma_0^{-1} m_0^T m_0}_{\text{constante}} \right)$$

Derivamos $P(w|t_0, x_0)$ respecto a w :

$$\frac{dP(w|t_0, x_0)}{dw} = -\frac{1}{2\sigma_n^2} \left(0 - 2t_0 \phi(x_0)^T + 2\phi(x_0)^T \phi(x_0) w \right) - \frac{1}{2} \left(2 \Sigma_0^{-1} w - 2 \Sigma_0^{-1} m_0^T + 0 \right)$$

$$-\frac{1}{2} \left(\frac{-2}{\sigma_n^2} t_0 \phi(x_0)^T + \frac{2\phi(x_0)^T \phi(x_0) w}{\sigma_n^2} - 2 \Sigma_0^{-1} w - 2 \Sigma_0^{-1} m_0^T \right) = 0$$

$$\underbrace{\frac{\phi(x_0)^T \phi(x_0) w}{\sigma_n^2}}_{\text{cuadrático}} + \underbrace{\Sigma_0^{-1} w}_{\text{lineal}} = \underbrace{\frac{t_0 \phi(x_0)^T}{\sigma_n^2}}_{\text{cuadrático}} + \underbrace{\Sigma_0^{-1} m_0^T}_{\text{lineal}}$$

$$w \left(\Sigma_0^{-1} I + \frac{\phi(x_n)^T \phi(x_n)}{\theta^2_n} \right) = \frac{t_n \phi(x_n)^T}{\theta^2_n} + \Sigma_0^{-1} m_0^T$$

Pesos optimos:

$$w = \left(\Sigma_0^{-1} I + \frac{\phi(x_n)^T \phi(x_n)}{\theta^2_n} \right)^{-1} \left(\Sigma_0^{-1} m_0^T + \frac{t_n \phi(x_n)^T}{\theta^2_n} \right)$$

$$p(w) = N(w | 0, \theta^2 w) \rightarrow \frac{p(w | t_n, x_n)}{p(w | t_n)} = N(w | \bar{m}_N, \bar{S}_N)$$

$$w = \left(\frac{1}{\theta^2_n} I + \frac{\phi(x_n)^T \phi(x_n)}{\theta^2_n} \right)^{-1} \frac{\phi(x_n)^T t_n}{\theta^2_n} = \bar{m}_N$$

\bar{S}_N

$$\bar{S}_N = \left(\frac{1}{\theta^2_n} \right)^{-1} \left(\frac{\theta^2_n}{\theta^2_n} I + \frac{\phi(x_n)^T \phi(x_n)}{\theta^2_n} \right)^{-1} \phi(x_n)^T t_n$$

$$w = \left(\frac{\phi(x_n)^T \phi(x_n)}{\theta^2_n} + \lambda I \right)^{-1} \phi(x_n)^T t_n$$

Regresion rigida kernel:

Este modelo mapea x_n a un espacio de dimension mayor con $\phi(x_n)$

$$k(x, x') = \phi(x)^T \phi(x') \quad \text{mejora}$$

Arrancamos de Bayesiano

$$p(w | t_n) = \| t_n - \phi(x_n) w^T \|_2^2 + \lambda \| w \|^2_2$$

$$w^* = \underset{w}{\operatorname{argmin}} \| t_n - \phi(x_n) w^T \|_2^2 + \lambda \| w \|^2_2 = (\phi^T(x_n) \phi(x_n) + \lambda I)^{-1} \phi^T(x_n) t_n$$

$$w^* = \underbrace{(\phi^T(x_n) \phi(x_n) + \lambda I)^{-1}}_{k(x, x')} t_n \underbrace{\phi^T(x_n)}_{\phi(x_i)^T}$$

$$w^* = \sum_i \alpha_i \phi(x_i)^T$$

Se necesita predicción para un nuevo punto proyectando sobre w^*

$$t_x = \phi(x_x)^T w^* = \phi(x_x)^T \phi^T(x_n) \underbrace{(\phi^T(x_n) \phi(x_n) + \lambda I)^{-1}}_{k(x, x')} t_n$$

$$t_x = (k + \lambda I)^{-1} t_n k(x)^T$$

$$k(x)^T = (k(x_x, x_1), \dots, k(x_x, x_n))$$

Procesos gaussianos

Este modelo generaliza la regresión de kernel modelando distribuciones sobre funciones y se define por la función media y función covarianza

$$GP: F(x) \sim GP(m(x), K(x, x'))$$

$$F(x) = \phi(x)^T w : w \sim N(0, \Sigma_{MAP}) \rightarrow \text{Sin ruido}$$

\downarrow
SN

media:

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

Covarianza:

$$K(x, x') = \text{Cov}(F(x), F(x'))$$

$$K(x, x') = E\{F(x), F(x')\} = E\{\phi(x)^T w \phi(x')^T w\}$$

$$K(x, x') = \phi(x)^T E\{ww^T\} \phi(x')$$

$$K(x, x') = \phi(x)^T \Sigma_w \phi(x')$$

$$F(x) \sim GP(F(x)|0, K) : K \in R^{n \times n} \sim K_{ij} = K(x_i, x_j)$$

$$t_n = F(x) + \eta_n \rightarrow \text{Ruido } N(0, \sigma_n^2)$$

Prior sobre $F(x)$:

$$P(F(x)) = N(F(x)|0, K)$$

Verosimilitud:

$$P(t_n|F(x)) = N(t_n|F(x), \sigma_n^2 I_N)$$

Distribución marginal de t_n :

$$P(t_n) = \int P(t_n|F(x)) P(F(x)) dF(x) = N(t_n|0, K + \sigma_n^2 I_N)$$

$$= \int N(t_n|F(x), \sigma_n^2 I_N) (\sigma_n^2 I_N)^{-1} (t_n - F(x)) - \frac{1}{2} F(x)^T K^{-1} F(x) dF(x)$$

$$= \int e^{-\frac{1}{2} ((\sigma_n^2 I_N)^{-1} (\|t_n\|^2 - 2 t_n^T F(x) + F(x)^T F(x)) + F(x)^T F(x) K^{-1})} dF(x)$$

$$= \int e^{-\frac{1}{2} \underbrace{((\sigma_n^2 I_N)^{-1} + K^{-1}) F(x)^T F(x)}_a + \underbrace{(\sigma_n^2 I_N)^{-1} t_n^T F(x)}_b - \underbrace{\frac{(\sigma_n^2 I_N)^{-1} t_n^T t_n}{2}}_c} dF(x)$$

$$= \int e^{-\frac{1}{2} a F(x)^T F(x) + b F(x) + c} dF(x) \rightarrow MF = a^{-1} b$$

$$= \int e^{-\frac{1}{2} (F(x) - MF)^T a (F(x) - MF) + \frac{1}{2} b^T a^{-1} b + c} dF(x)$$

$$= e^{\frac{1}{2} b^T a^{-1} b + c}$$

$$= e^{\frac{1}{2} (\sigma_n^4 I_N)^{-1} t_n^T a^{-1} t_n - \frac{(\sigma_n^2 I_N)}{2} t_n^T t_n}$$

Día _____ Mes _____ Año _____

Lun Mar Mie Jue Vie Sab

☐ ☐ ☐ ☐ ☐ ☐

$$p(t_n) = e^{-\frac{1}{2} t_n^T \left(-(\theta_n^2 I_N)^{-1} (k + \theta_n^2 I_N)^{-1} k + (\theta_n^2 I_N)^{-1} \right) t_n}$$

$$= e^{-\frac{1}{2} t_n^T \underbrace{(k + \theta_n^2 I_N)^{-1}}_{\text{Varianza}} t_n}$$

$$p(t_n) = N(t_n | 0, k + \theta_n^2 I_N)$$