

PARCIAL 1

Modelos de Regresión:

$w \in \mathbb{R}^Q$; donde $\phi: \mathbb{R}^P \rightarrow \mathbb{R}^Q$

$t_n = \phi(x_n)w^T + \eta_n$; $\{t_n \in \mathbb{R}, x_n \in \mathbb{R}^P\}_{n=1}^N$; $\eta_n \sim N(\eta_n | 0, \sigma_n^2)$

target
Datos
Vector de funciones

Mínimos Cuadrados:

$$\hat{E} = \arg \min_w \|t - \phi w^T\|^2$$

Se minimiza

$$\frac{\partial E}{\partial w} = \langle t - \phi w^T, t - \phi w^T \rangle \rightarrow (t - \phi w^T)^T (t - \phi w^T)$$

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

$$\frac{\partial (||t||^2 - 2t^T \phi w^T + \phi^T w^T \phi w^T)}{\partial w} = 0$$

$$0 - 2t^T \phi^T + 2\phi^T \phi w = 0$$

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

Pesos

Mínimos Cuadrados Regularizados: (Ridge Regression)

Don los min. Cuadrados con regularización λ
Penaliza pesos grandes.

$$E = \arg \min_w \|t - \phi w^T\|_2^2 + (\lambda \|w\|_2^2) \rightarrow \text{Regularización}$$

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

$$\frac{\partial (||t||^2 - 2t^T \phi^T w + \phi^T w \phi w^T + \lambda w w^T)}{\partial w} = 0$$

$$0 - 2t^T \phi^T + 2\phi^T \phi w + 2\lambda w^T = 0$$

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

Pesos grandes

Máxima Verosimilitud (MLE):

Es un modelo probabilístico, describe t_n según la entrada x_n
 w se describe según el ruido η_n

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

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

$$P(t_n | \phi(x_n)w^T, \sigma_n^2) = N(t_n | \phi(x_n)w^T, \sigma_n^2)$$

distribución Normal para t_n
 μ_x
 σ_x^2

$$P(t_n | \phi(x_n)w^T, \sigma_n^2) = \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(t_n - \phi(x_n)w^T)^2}{2\sigma_n^2}}$$

func. de densidad para cada t_n

$$w_{ML} = \arg \max_{w, \sigma_n^2} \log \left(\prod_{n=1}^N N(t_n | \phi(x_n)w^T, \sigma_n^2) \right)$$

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

Observaciones independientes
 \rightarrow iid, por eso \prod

$$L(w, \sigma_n^2) = \log \left(\prod_{n=1}^N \left(\frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(t_n - \phi(x_n)w^T)^2}{2\sigma_n^2}} \right) \right)$$

con $\log \prod = \sum$

...

$$L(w, \sigma_n^2) = \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_n^2) - \frac{(t_n - \phi(x_n)w^T)^2}{2\sigma_n^2} \right)$$

$$L(w, \sigma_n^2) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_n^2) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N (t_n - \phi(x_n)w^T)^2$$

para σ_n^2 grandes
Error cuadrático respecto w_n, σ_n^2

Maximización respecto w_{ML} (Pesos óptimos):

$$\frac{\partial L(w, \sigma_n^2)}{\partial w} = 0 - 0 - \frac{\partial}{\partial w} \left(\frac{1}{2\sigma_n^2} \|t_n - \phi(x_n)w^T\|_2^2 \right)$$

= Mínimos Cuadrados

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

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

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

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

$$\phi(x_n)^T \phi(x_n)w = \phi(x_n)^T t$$

$$w_{ML} = (\phi(x_n)^T \phi(x_n))^{-1} \phi(x_n)^T t$$

Pesos óptimos

Maximizamos respecto σ_{ML}^2 (Varianza óptima)

$$\frac{\partial L(w, \sigma_n^2)}{\partial \sigma_n^2} = 0 - \frac{N}{2\sigma_n^2} + \frac{1}{2(\sigma_n^2)^2} \|t_n - \phi(x_n)w^T\|_2^2 = 0$$

$$\frac{\|t_n - \phi(x_n)w^T\|_2^2}{2(\sigma_n^2)^2} = \frac{N}{2\sigma_n^2}$$

$$\sigma_n^2 = \frac{1}{N} \|t_n - \phi(x_n)w^T\|_2^2 \rightarrow \text{Promedio Er residual}$$

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

Varianza óptima
Distancia o diferencia entre predicciones y target

Máximo a-posteriori: (MAP)

Con prior $P(w)$

Maximiza distribución a-posteriori $P(w|t)$

$$P(t_n | \phi(x_n)w^T, \sigma_n^2) = N(t_n | \phi(x_n)w^T, \sigma_n^2)$$

Gaussianas

$$P(t_n | \phi(x_n)w^T, \sigma_n^2) = \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{-\frac{(t_n - \phi(x_n)w^T)^2}{2\sigma_n^2}}$$

$P(t_n | w)$

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

Pesos indep. con var. σ_w^2
Controla regularización

$$P(w) = \frac{1}{\sqrt{2\pi\sigma_w^2}} e^{-\frac{\|w\|_2^2}{2\sigma_w^2}}$$

El posteriori $P(w|t_n)$ se maximiza con verosimilitud y prior del teorema de Bayes.

$$P(w|t_n) = \frac{P(t_n|w)P(w)}{P(t_n)}$$

Maximizamos sólo numerador (Maximo a-posteriori)

Para pasar a Σ

$$W_{MAP} = \arg \max_w \log \left(\prod_{n=1}^N N(t_n | \phi(x_n)W^T, \sigma_n^2) \prod_{q=1}^Q N(w_q | 0, \sigma_w^2) \right)$$

$$\log P(w|t_n) = \underbrace{\log P(t_n|w)}_{\text{Verosimilitud}} + \underbrace{\log P(w)}_{\text{Prior}} \rightarrow \text{Procede verosimilitud de MLE}$$

$$\frac{d \log P(w|t_n)}{dw} = \frac{d}{dw} \left(-\frac{1}{2\sigma_n^2} \|t_n - \phi(x_n)w\|_2^2 + \text{cte} \right) + \frac{d \log P(w)}{dw}$$

$$\frac{d \log P(w)}{dw} = \frac{d}{dw} \left(\sum_{q=1}^Q \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma_w^2 - \frac{\|w\|_2^2}{2\sigma_w^2} \right) \right)$$

Asumimos iid =

$$W_{MAP} = \arg \max_w \left(-\frac{1}{2\sigma_n^2} \|t_n - \phi(x_n)w\|_2^2 - \frac{1}{2\sigma_w^2} \|w\|_2^2 \right)$$

Maximizar $\log P(w|t) = \min$.

$$W_{MAP} = \arg \max_w \left(\|t_n - \phi(x_n)w\|_2^2 + \frac{\sigma_n^2}{\sigma_w^2} \|w\|_2^2 \right)$$

$$W_{MAP} = \arg \max_w \left(\|t_n - \phi(x_n)w\|_2^2 + \lambda \|w\|_2^2 \right)$$

Mismo result. de min. Cuadrados regularizados

$$W_{MAP}^* = (\phi(x_n)^T \phi(x_n) + \lambda I)^{-1} \phi(x_n)^T t_n$$

Bayesiano con modelo lineal Gaussiano:

No trata de encontrar un único valor para W_{ML} como MAP, sino que se calcula toda la distribución.

$$P(w|t_n, x_n) = N(w | W_{MAP}, \Sigma_{MAP})$$

Distribución Posteriori de W en Bayesiano Gaussiano

Punto de MAP Covarianza a-posteriori

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

Captura la Incertidumbre

$$W_{MAP} = m_N = \Sigma_N \left(\Sigma_0^{-1} m_0 + \frac{1}{\sigma_n^2} \phi^T t \right)$$

Se asume $m_0 = 0 \rightarrow$ Punto de inicio

Matriz de Covarianza $\left\{ \begin{array}{l} \Sigma_0 \rightarrow \text{Incertidumbre inicial sobre } W \\ \text{(Correlaciones entre pesos)} \end{array} \right.$

$$P(w) = N(w | m_0, \Sigma_0) \rightarrow \text{Prior sobre pesos } (w)$$

$$P(t_n|w) = N(t_n | \phi(x_n)w^T, \sigma_n^2 I) \rightarrow \text{verosimilitud}$$

$$P(w|t_n, x_n) = P(t_n|w) \cdot P(w)$$

$$= \frac{1}{\sqrt{2\pi\sigma_n^2}} e^{\left(-\frac{1}{2\sigma_n^2} \|t_n - \phi w\|_2^2 - \frac{1}{2} \|w - w_0\|_2^2 \right)}$$

Min. Cuadrados
Min. Cuadrados Regularizados

$$P(w|t_n, x_n) = -\frac{1}{2\sigma_n^2} \left(\|t_n\|_2^2 - 2t_n \phi(x_n)^T w + \phi(x_n)^T w \phi(x_n) w^T \right) - \frac{1}{2} \Sigma_0^{-1} (w - m_0)^T (w - m_0)$$

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

$$= -\frac{1}{2} \left(\underbrace{\Sigma_0^{-1} w^T w}_{\text{Cuadrático}} - \underbrace{2 \Sigma_0^{-1} m_0^T w}_{\text{Lineal}} + \underbrace{\Sigma_0^{-1} m_0^T m_0}_{\text{cte}} \right)$$

Optimizamos $P(w|t_n, x_n)$ respecto W :

$$\frac{d P(w|t_n, x_n)}{dw} = -\frac{1}{2\sigma_n^2} \left(0 - 2t_n \phi(x_n)^T + 2\phi(x_n)^T \phi(x_n) w \right) - \frac{1}{2} \left(2\Sigma_0^{-1} w - 2\Sigma_0^{-1} m_0^T + 0 \right)$$

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

$$\frac{\phi(x_n)^T \phi(x_n) w}{\sigma_n^2} + \Sigma_0^{-1} w = \frac{t_n \phi(x_n)^T}{\sigma_n^2} + \Sigma_0^{-1} m_0^T$$

Cuadrático Lineal

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

$$W = \left(\underbrace{\Sigma_0^{-1} I + \frac{\phi(x_n)^T \phi(x_n)}{\sigma_n^2}}_{\Sigma_N} \right)^{-1} \left(\underbrace{\Sigma_0^{-1} m_0^T + \frac{t_n \phi(x_n)^T}{\sigma_n^2}}_{m_N} \right)$$

Pesos óptimos

$$P(w) = N(w | 0, \sigma_w^2) \rightarrow \frac{P(w|t_n, x_n)}{P(w|t_n)} = N(w | \bar{m}_N, \bar{\Sigma}_N)$$

$$W = \left(\underbrace{\frac{1}{\sigma_w^2} I + \frac{\phi(x_n)^T \phi(x_n)}{\sigma_n^2}}_{\bar{\Sigma}_N} \right)^{-1} \frac{\phi(x_n)^T t_n}{\sigma_n^2} = \bar{m}_N$$

$$\bar{\Sigma}_N = \left(\frac{1}{\sigma_w^2} \right)^{-1} \left(\frac{\sigma_n^2}{\sigma_w^2} I + \frac{\phi(x_n)^T \phi(x_n)}{\sigma_n^2} \right)^{-1} \phi(x_n)^T t_n$$

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

Min. Cuadrados Reg. $\lambda = \frac{\sigma_n^2}{\sigma_w^2}$

Kernel Ridge: (KRR)

Para No linealidades
Mapea x_n a un espacio de dimensión mayor con $\phi(x_n)$

Mejora:

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

Partimos de Bayesiano completo:

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

$$w^* = \arg \min_w \|t_n - \phi(x_n)w\|_2^2 + \lambda \|w\|_2^2$$

$$w^* = (\phi^T(x_n) \phi(x_n) + \lambda I)^{-1} \phi^T(x_n) t_n$$

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

α_i

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

Necesitamos predicción para un nuevo punto: Proyectando sobre solución W^*

$$t_* = \phi(X_*)^T W^*$$
$$t_* = \underbrace{\phi(X_*)^T \phi^T(X_n)}_{K(*)^T} \underbrace{(\phi^T(X_n) \phi(X_n) + \lambda I)^{-1} t_n}_{\alpha_i}$$
$$t_* = (K + \lambda I)^{-1} t_n K(*)^T \rightarrow [K(X_*, X_1), \dots, K(X_*, X_N)]$$

Procesos Gaussianos: (GPs / GPR)

Generaliza la regresión kernel modelando distribuciones sobre func.

Define por \rightarrow func. media $m(x) = 0$
 \rightarrow func. covarianza $K(x, x')$

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

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

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

$K(x, x') = Cov(f(x), f(x'))$

$K(x, x') = E\{f(x), f(x')\} = E\{\phi(x)^T w \phi^T(x') w^T\}$
 $= \phi(x)^T E\{w w^T\} \phi(x') = \phi(x)^T \Sigma_w \phi(x') \rightarrow$ Cov

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

$t_n = f(x) + \eta_n \rightarrow$ Ruido
 $\rightarrow N(0, \sigma_n^2)$

$P(f(x)) = N(f(x)|0, K) \rightarrow$ Prior sobre $f(x)$

$p(t_n|f(x)) = N(t_n|f(x), \sigma_n^2 I_N) \rightarrow$ Verosimilitud

Hallamos distribución marginal de t_n :

$P(t_n) = \int p(t|f(x)) p(f(x)) df(x) = N(t|0, K + \sigma_n^2 I_N)$

$P(t_n) = \int N(t_n | f(x), \underbrace{\sigma_n^2 I_N}_{S_0}) \underbrace{(\sigma_n^2 I_N)^{-1}}_{S_0} (t_n - f(x))$
 $- \frac{1}{2} f(x)^T K^{-1} f(x)) df(x)$

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

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

$P(t_n) = \int e^{(-\frac{1}{2} a f(x)^T f(x) + b f(x) + c)} df(x) \rightarrow M_f = a^{-1} b$

$P(t_n) = \int e^{(-\frac{1}{2} (f(x) - M_f)^T a (f(x) - M_f) + \frac{1}{2} b^T a^{-1} b + c)} df(x)$

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

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

$P(t_n) = e^{(-\frac{1}{2} t_n^T (-(\sigma_n^2 I_N)^{-1} (K + \sigma_n^2 I_N)^{-1} K + (\sigma_n^2 I_N)^{-1}) t_n)}$

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

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