

PARCIAL #1

Preguntas

1. Sea el modelo de regresión $t_n = \phi(x_n)w^T + \eta_n$, con $\{t_n \in \mathbb{R}, x_n \in \mathbb{R}^P\}_{n=1}^N$, $w \in \mathbb{R}^Q$, $\phi: \mathbb{R}^P \rightarrow \mathbb{R}^Q$, $Q \geq P$, y $\eta_n \sim \mathcal{N}(\eta_n|0, \sigma_\eta^2)$. Presente el problema de optimización y la solución del mismo para los modelos mínimos cuadrados, mínimos cuadrados regularizados, máxima verosimilitud, máximo a-posteriori, Bayesiano con modelo lineal Gaussiano, regresión rígida kernel y mediante procesos Gaussianos. Asuma datos i.i.d. Discuta las diferencias y similitudes entre los modelos estudiados.

1) Sea el modelo de regresión

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

Con $\{t_n \in \mathbb{R}, x_n \in \mathbb{R}^P\}_{n=1}^N$

$$Q > P \\ \eta_n \sim \mathcal{N}(0, \sigma_n^2)$$

Donde $\{t_n \in \mathbb{R}, x_n \in \mathbb{R}^P\}_{n=1}^N$, $w \in \mathbb{R}^Q$ y $\phi: \mathbb{R}^P \rightarrow \mathbb{R}^Q$ es un vector de función base ($Q \geq P$)

En donde definimos:

$$t = \begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}, \quad \Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{pmatrix} \in \mathbb{R}^{N \times Q}$$

De manera que el modelo completo se puede escribir como:

$$t = \Phi w + \eta, \quad \therefore \eta \sim \mathcal{N}(0, \sigma_n^2 I_N)$$

Ahora para los siguientes:

- * MÍNIMOS CUADRADOS
- * MÍNIMOS CUADRADOS REGULARIZADOS
- * MÁXIMA VEROSIMILITUD
- * MÁXIMA A-POSTERIORI
- * BAYESIANO CON MODELO LINEAL GAUSSIANO
- * REGRESIÓN RÍGIDA KERNEL
- * PROCESOS GAUSSIANOS

→ MÍNIMOS CUADRADOS:

Asumimos que las observaciones t_n están relacionadas con x_n , con el modelo:

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

En donde:

- * $\phi(x_n)^T w \rightarrow$ es la predicción para la entrada x_n
- * w son los parámetros
- * η_n es el ruido blanco gaussiano

SUPONEMOS LAS OBSERVACIONES t_n ESTÁN DISTRIBUIDAS COMO:

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

→ Para t_n la probabilidad conjunta es:

$$P(t | \phi, w) = \prod_{n=1}^N P(t_n | x_n, w)$$

EN DONDE ϕ ES LA MATRIZ DE CARACTERÍSTICAS

→ El error e_n (CUANTO X DESVIA LA PREDICCIÓN)

$$e_n = t_n - \hat{t}_n \quad \therefore t_n \in \mathbb{R}^n$$

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

AHORA PARA PLANTEAR EL PROBLEMA, LA DISTRIBUCIÓN GAUSSIANA DE t_n Y e_n

$$P(t | \phi, w) = \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp \left[-\frac{1}{2\sigma_n^2} \|t_n - \phi w\|^2 \right]$$

DOnde:

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

NOS INTERESA MINIMIZAR LA PERDIDA PROMEDIO DEL ERROR ENTONCES:

$$\text{ASUMIENDO i.i.d} = E\{x_n\} \approx \frac{1}{N} \sum_{n=1}^N x_n$$

DEBEMOS MINIMIZAR w O $\|t - \phi w\|^2$:

$$w^* = \operatorname{argmin} E\{e_n\}$$

$$w^* = \operatorname{argmin} E\{\|t - \phi w\|^2\}$$

$$E\{\|t - \phi w\|^2\} = \frac{1}{N} [\|t - \phi w\|^2]$$

$$\begin{aligned} &= \frac{1}{N} (t - \phi w)^T (t - \phi w) = \frac{1}{N} [t^T t - t^T \phi w - \phi w^T + (\phi w)^T \phi w] \\ &= \frac{1}{N} [t^T t - 2t^T \phi w + w^T \phi^T \phi w] \end{aligned}$$

LUEGO TENEMOS QUE:

$$\frac{\partial}{\partial w} E\{\|t - \phi w\|^2\} = 0$$

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

$$2\phi^T \phi w = 2(t^T \phi)^T$$

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

$$\hookrightarrow \text{AHORA CON LA PSEUDODINVERSA: } w^* = (\phi^T \phi)^{-1} \phi^T t$$

→ MÍNIMOS CUADRADOS REGULARIZADOS

Teniendo el modelo:

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

y el error cuadrático

$$e_n^2 = \|t - w\phi\|_2^2$$

Debido a que mínimos cuadrados es un modelo sensible en datos atípicos, es necesario un término de regularización.

La regularización L^2 se define como:

$$\|w\|_2^2 = \sum_{j=1}^D w_j^2$$

Luego el término de regularización es:

$$\lambda \|w\|_2^2$$

Por lo que tenemos:

$$\|t - \phi^T w\|_2^2 + \lambda \|w\|_2^2$$

Antes el problema de optimización es:

$$w = \underset{w}{\operatorname{Argmin}} \|t - \phi w\|_2^2 + \|w\|_2^2 \lambda$$

De mínimos cuadrados tenemos

$$\begin{aligned} \|t - \phi w\|_2^2 &= (t - \phi w)^T (t - \phi w) \\ &= t t^T - 2 t^T \phi w + w^T \phi^T \phi w \end{aligned}$$

$$w^* \operatorname{argmin} [t t^T - 2 t^T \phi w + w^T \phi^T \phi w + w^T w \lambda] = 0$$

$$\frac{d}{dw} [t t^T - 2 t^T \phi w + (w \phi)^T \phi w + w^T w \lambda] = 0$$

$$\frac{d}{dw} [t t^T - 2 t^T \phi w + w^T (\phi^T \phi + \lambda I) w] = 0$$

$$[-2 \phi^T t] + 2 (\phi^T \phi + \lambda I) w = 0$$

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

$$2 (\phi^T \phi + \lambda I) w = \phi^T t$$

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

→ MÁXIMA VEROSIMILITUD

SE TIENE EL MODELO DE REGRESIÓN

$$t_n = \phi(x_n) w^T + \eta_n \quad \therefore \eta_n \sim N(\eta_n | 0, \sigma_n) \rightarrow \text{RUIDO BLANCO GAUSSIANO}$$

ENTONCES

Y DATOS i.i.d
(SUPONEMOS)

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

AHORA PARA UNA SOLA OBSERVACIÓN (x_n, t_n) TENEMOS LA VEROSIMILITUD

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

TENIENDO EL SUPUESTO DE i.i.d LA VEROSIMILITUD CONJUNTA

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

AHORA COMO ASUMIMOS UNA NORMAL TENEMOS $\log(N(t_n | \phi(x_n) w^T, \sigma_n^2))$

$$N(t | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[\frac{-(t - \phi(x_n) w^T)^2}{2\sigma_n^2} \right]$$

ENTONCES LA LOG-VEROSIMILITUD ES LA SUMA DE LAS LOG-VEROSIMILITUDES INDIVIDUALES

$$\begin{aligned} \log(P(x)) &= \log \left(\prod_{n=1}^N N(x_n | \mu, \sigma^2) \right) \\ &= \log \left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(\frac{-\|x_n - \mu\|^2}{2\sigma^2} \right) \right) \\ &= \log \left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left(\exp \left(\frac{-\|x_n - \mu\|^2}{2\sigma^2} \right) \right) \end{aligned}$$

LUEGO VEMOS QUE:

$$\log \left(\frac{1}{(2\pi\sigma^2)^{N/2}} \right) + \log \left(\exp \left(- \left[\frac{\|x_1 - \mu\|^2}{2\sigma^2} + \frac{\|x_2 - \mu\|^2}{2\sigma^2} + \dots + \right] \right) \right)$$

PODEMOS ESCRIBIR EL ARGUMENTO DEL EXPONENTE COMO LA SUMATORIA DE LAS MUESTRAS

$$\log \left(\frac{1}{(2\pi\sigma^2)^{N/2}} \right) + \log \left(\exp \left(- \frac{1}{2\sigma^2} \sum_{n=1}^N \|x_n - \mu\|^2 \right) \right)$$

LUEGO TENEMOS QUE:

$$\log(p(x)) = \frac{-N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N \|x_n - \mu\|^2$$

EL OBJETIVO ES ENCONTRAR LOS PARÁMETROS QUE MAXIMIZAN LA VEROSIMILITUD Y ENCONTRAR LOS PESOS Y VARIANZAS BAJO LAS SUPOSICIONES ANTERIORES, POR LO TANTO EL PROBLEMA DE OPTIMIZACIÓN ES EL SIGUIENTE

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

Resolviendo y usando log-likelihood

$$w_{ml} | \sigma_n^2 = \arg \max_{w, \sigma_n^2} \left[-\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_n^2) \right] - \frac{1}{2\sigma^2} \|t - \phi w^T\|$$

Ahora derivando respecto a σ^2

$$\frac{d}{d\sigma^2} \log p(t|x, w, \sigma_n^2) = \left(\frac{d}{d\sigma^2} \left[-\frac{N}{2} \log(2\pi) \right] - \frac{d}{d\sigma^2} \left[\frac{N}{2} \log(\sigma_n^2) \right] \right)$$

$$\frac{d}{d\sigma^2} \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N \|t - \phi w^T\|^2 \right]$$

Resolviendo

$$\frac{d}{d\sigma^2} p(t|w, \sigma^2) = \frac{-N}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{n=1}^N \|t_n - \phi(x_n) w^T\|^2 = 0$$

$$-N\sigma_n^2 + \sum \|t_n - \phi(x_n) w^T\|^2 = 0$$

$$\sigma_{ml}^2 = \frac{1}{N} \sum \|t_n - \phi(x_n) w^T\|^2$$

Derivando con respecto a w :

$$\frac{d}{dw} [\log p(t|x, w, \sigma_n^2)] = -\frac{1}{2\sigma^2} \sum_{n=1}^N 2(t_n - \phi(x_n) w^T)(-\phi(x_n))$$

$$\frac{d}{dw} [\log p(t|x, w, \sigma_n^2)] = 0$$

$$\frac{-1}{2\sigma^2} \sum_{n=1}^N 2(t_n - \phi(x_n) w^T)(-\phi(x_n)) = 0$$

$$\boxed{\frac{1}{\sigma^2} (t_n - \phi(x_n) w^T) \phi(x_n) = 0}$$

De forma matricial

$$\frac{1}{\sigma^2} (t - \phi w^T) \phi^T = 0$$

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

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

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

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

→ Máximo A-Posteriori (MAP)

Este modelo busca estimar los parámetros probabilísticos con la interpretación de información previa antes de observar datos en estadística bayesiana, de aquí tenemos el teorema de Bayes

$$P(x, y) = \frac{\overbrace{f(y|x)}^{\text{Verosimilitud}} \overbrace{f(x)}^{\text{Prior}}}{\underbrace{f(y)}_{\text{Evidencia}}}$$

POSTERIOR

Entonces

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

* Con evidencia

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

* Con prior

$$P(w) = \prod_{n=1}^Q N(w|0, \sigma_w^2)$$

* Verosimilitud

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

$$P(t | \phi, w, \sigma_n) = \prod_{n=1}^N N(t_n | \phi(x_n)w, \sigma_n^2) \rightarrow \text{MAP}$$

Por lo que el modelo simplifica la aplicación de Bayes mediante la proporcionalidad

$$P(w|t) \propto P(t|w) P(w)$$

Ahora para determinar la función de costo

$$\mathcal{J}(t, t(x)) = P(w|t, \phi, \sigma_n^2)$$

$$= \prod_{n=1}^N N(t_n | \phi(x_n)w, \sigma_n^2) \prod_{j=1}^Q N(w_j | 0, \sigma_w^2)$$

$$\log(\mathcal{J}(t, t(x))) = \log\left(\prod_{n=1}^N N(t_n | \phi(x_n)w, \sigma_n^2)\right) \prod_{j=1}^Q N(w_j | 0, \sigma_w^2)$$

$$= \log\left(\prod_{n=1}^N N(t_n | \phi(x_n)w, \sigma_n^2)\right) + \log\left(\prod_{j=1}^Q N(w_j | 0, \sigma_w^2)\right)$$

$$= \log\left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{\|t_n - \phi(x_n)w\|^2}{2\sigma_n^2}\right)\right) + \log\left(\prod_{j=1}^Q \frac{1}{\sqrt{2\pi\sigma_w^2}} \exp\left(-\frac{\|w_j - 0\|^2}{2\sigma_w^2}\right)\right)$$

$$= \frac{N}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} \sum_{n=1}^N \|t_n - \phi(x)w\|^2 + \log\left(\frac{1}{\sqrt{2\pi\sigma_w^2}}\right) + \log\left(\exp\left(-\frac{1}{2\sigma_w^2} \sum_{j=1}^Q \|w_j\|^2\right)\right)$$

$$= \frac{N}{2} \log(2\pi\sigma_n^2) - \frac{Q}{2} \log(2\pi\sigma_w^2) - \frac{1}{2\sigma_n^2} \|t - \phi w\|^2 - \frac{1}{2\sigma_w^2} \sum_{j=1}^Q \|w_j\|^2$$

Ahora para fines de optimización se ignoran los términos que no dependen de w , el problema queda

$$w_{\text{MAP}} = \underset{w}{\operatorname{argmax}} \quad -\frac{1}{2\sigma_n^2} \|t - \phi w^T\|_2^2 - \frac{1}{2\sigma_w^2} \|w\|_2^2$$

Como los factores de escala no modifican el punto máximo o mínimo, el problema se transforma

$$w_{\text{MAP}} = \underset{w}{\operatorname{argmin}} \quad \|t - \phi w^T\|_2^2 + \frac{\sigma_n^2}{\sigma_w^2} \|w\|_2^2$$

En otras palabras, el problema se reduce a una optimización por mínimos cuadrados con $\lambda = \frac{\sigma_n^2}{\sigma_w^2}$, el cual se resolvió anteriormente y queda:

$$w_{\text{MAP}}^* = \left(\phi^T \phi + \frac{\sigma_n^2}{\sigma_w^2} I \right)^{-1} \phi^T t$$

→ Bayesiano con modelo lineal Gaussiano

Como los demás modelos se requiere incorporar conocimiento previo en el prior tenemos la verosimilitud

$$p(t|w) = \mathcal{N}(t | \phi w, \sigma_n^2)$$

y el prior

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

La distribución del prior

$$p(w|t) \propto p(t|w) p(w)$$

Asumimos el prior

$$p(w) = \mathcal{N}(w_0 | S_0)$$

y la posterior como

$$p(w|t) = \mathcal{N}(w | m_n, S_n)$$

$$S, \quad p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$

$$p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$$

$$\text{con } y = f(x|A, b) = Ax + b$$

$$\text{Luego } p(y|x) = \mathcal{N}(x | \mu_x | y, \Sigma_x)$$

$$\mu_x | y = (\Lambda + A^T L A)^{-1} (A^T L (y - b) + \Lambda \mu)$$

$$\Sigma_x | y = (\Lambda + A^T L A)^{-1}$$

Es el caso $t_n = w^T \phi(x_n) t_n = N(t_n | w^T \phi(x_n), \beta^{-1})$

$$\hookrightarrow p(t|w) = N(t | \phi w, \beta^{-1} I)$$

$$\text{Luego } w \sim p(w) = N(w | m_0, S_0)$$

Es posterior $p(w|t) = N(w | m_n, S_n)$

$$M_n = M_w | t = (S_0^{-1} + \phi^T B I \phi)^{-1} (\phi^T B I (t - 0) + S_0^{-1} m_0)$$

$$M_n = M_w | t = (S_0^{-1} + B \phi^T \phi)^{-1} (B \phi^T t + S_0^{-1} m_0)$$

$$M_n = M_n | t = S_n (S_0^{-1} m_0 + B \phi^T t)$$

Luego se tiene

$$S_n = (S_0^{-1} + B \phi^T \phi)^{-1}$$

$$S_n^{-1} = S_0^{-1} + B \phi^T \phi$$

$$\text{Ahora } B = \frac{1}{\sigma_n^2}$$

La media posterior

$$m_n = S_n \left(S_0^{-1} m_0 + \frac{1}{\sigma_n^2} \phi \phi^T t \right)$$

La covarianza posterior

$$S_n = \left(S_0^{-1} + \frac{1}{\sigma_n^2} \phi \phi^T \phi \right)^{-1}$$

Como suponemos $p(w|t) = N(w | \tilde{m}_n, \tilde{S}_n)$

Suponemos $m_0 = 0$

$$m_n = S_n \left(\frac{1}{\sigma_n^2} \phi \phi^T t \right) (*)$$

$$S_n = \sigma_w^2 I_Q$$

$$S_n = \left((\sigma_w^2 I_Q)^{-1} + \frac{1}{\sigma_n^2} \phi \phi^T \phi \right)^{-1}$$

$$\tilde{S}_n = \left(\frac{1}{\sigma_w^2} I_Q + \frac{1}{\sigma_n^2} \phi \phi^T \phi \right)^{-1}$$

$$\tilde{S}_n = \left[\frac{1}{\sigma_n^2} \left(\frac{\sigma_n^2}{\sigma_w^2} I_Q + \phi \phi^T \phi \right) \right]^{-1} = \sigma_n^2 \left(\frac{\sigma_n^2}{\sigma_w^2} I_Q + \phi \phi^T \phi \right)^{-1} (*)$$

Antes (*) es (*) (*)

$$\bar{w}_n = \frac{1}{\sigma_z^2} \sigma_n^2 \left(\frac{\sigma_n^2}{\sigma_w^2} I + \Phi \Phi^T \right)^{-1} \Phi \Phi^T \epsilon$$

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

Por lo que la solución del modelo es equivalente de \bar{w}_n a la solución de mínimos cuadrados regularizados, teniendo

$$\lambda = \frac{\sigma_n^2}{\sigma_w^2}$$

→ REGRESION RÍGIDA KERNEL

A diferencia de la regresión lineal clásica, RK incluye un término de regularización que penaliza la magnitud de los coeficientes para evitar el sobreajuste

Tenemos

$$y \in \mathbb{R}^N; \quad X \in \mathbb{R}^{N \times P} \quad H \subseteq \mathbb{R}^Q \quad Q \rightarrow \infty$$

$$\hat{y} = f(u(x)|w) = \phi(x)w$$

$$f: \mathbb{R}^Q \rightarrow \mathbb{R} \quad w \in H \subseteq \mathbb{R}^Q$$

El problema se resume en:

$$w^* = \arg \min_w \frac{1}{N} (y, f(\phi(x)|w)) + R(t, \lambda)$$

Consiste en minimizar el error más un parámetro de regularización, similar al problema de mínimos cuadrados regularizados

$$e = \|y - \phi(x)w\|^2 \quad r(t, \lambda) = \lambda \|w\|^2$$

Entonces

$$w^* = \arg \min_w \frac{1}{N} \|y - \phi(x)w\|^2 + \lambda \|w\|^2 \quad (i)$$

donde $\lambda \in \mathbb{R}^+$ (fijo)

$$\|y - \phi(x)w\|^2 = (y - \phi(x)w)(y - \phi(x)w)^T$$

$$yy^T = -2y^T \phi(x)w + (\phi(x)w)^T (\phi(x)w)$$

$$yy^T - 2y^T \phi(x)w + w^T \phi^T \phi w \quad (ii)$$

Ahora haciendo (i) en (ii), derivando la función de costo

$$\frac{d}{d\omega} \{ L(y, f) + R(f | \lambda) \} = \frac{1}{N} \frac{d}{d\omega} [y^T y - 2\phi(x)y + \omega^T \phi^T(x) \phi(x) \omega] + \frac{d}{d\omega} [\lambda \omega^T \omega]$$

$$-\frac{2}{N} \phi(x)^T y + \frac{2}{N} \omega^T \phi^T(x) \phi(x) + 2\lambda \omega = 0$$

$$\frac{2}{N} \omega^T \phi^T(x) \phi(x) + 2\lambda \omega = \frac{2}{N} \phi^T(x) y$$

$$\omega^T \phi^T(x) \phi(x) + N\lambda \omega = \phi^T(x) y$$

$$\omega [\phi^T(x) \phi(x) + N\lambda I] = \phi^T(x) y$$

$$\omega^* = [\phi^T(x) \phi(x) + N\lambda I]^{-1} \phi^T(x) y$$

Ahora la dimensión del espacio transformado es demasiado grande (∞), lo que hace que calcular la inversa cueste mucho computacionalmente o sea imposible, por lo que se debe reescribir $[\phi^T \phi(x) + N\lambda I]$ usando la identidad:

$$(I + AB)^{-1} A = A(I + BA)^{-1}$$

Con lo que simplificamos $(\phi(x)^T \phi(x) + N\lambda I)^{-1} \phi^T(x)$, ahora usando la identidad:

$$(\phi^T(x) \phi(x) + N\lambda I)^{-1} \phi^T(x) = \left[N\lambda \left(\frac{1}{N\lambda} (\phi^T \phi + I) \right) \right]^{-1} \phi^T$$

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

$$= \phi^T \left[N\lambda \left(I + \frac{\phi}{N\lambda} \phi^T \right) \right]^{-1} = \phi^T (N\lambda I + \phi \phi^T)$$

$$\omega^* = \phi^T(x) [N\lambda I + \phi(x) \phi^T(x)]^{-1} y$$

$$\phi(x_{new}) \in \mathbb{R}^Q$$

Ahora por kernel trick

$$K = \phi(x) \phi(x)^T \quad K_{ij} = K(x_i, x_j)$$

y a de la nueva predicción

$$K_{new}^T = [K(x_{max}, x_n)]_{n=1}^N$$

$$K_{new}^T = \phi(x) \phi(x_n)$$

$$\hat{y} = t(x_{new}) = K_{new}^T (N\lambda I + K)^{-1} y$$

→ Procesos Gaussianos

Este modelo busca incorporar a los parámetros del modelo como en las predicciones

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

Asumimos el prior de forma Gaussian

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

Consideramos $m_0 = 0$, $S_0 = \sigma_w^2 I_Q$ entonces $p(w)$

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

Ahora la distribución del posterior

$$p(w|t) \propto p(t|w) p(w)$$

La verosimilitud

$$p(t|w) = N(t | \phi w, \sigma_n^2)$$

y la log-verosimilitud

$$\log p(t|w) = -\frac{N}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (t - \phi w)^T (t - \phi w)$$

Por lo que el prior

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

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

Ahora sumando las expresiones

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

$$\log p(w|t) \propto -\frac{1}{2\sigma_n^2} (t - \phi w)^T (t - \phi w) - \frac{1}{2\sigma_w^2} w^T w$$

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

$$\log p(w|t) \propto -\frac{1}{2\sigma_n^2} (t^T t - 2t^T \phi w + w^T \phi^T w) - \frac{1}{2\sigma_w^2} w^T w$$

$$\log p(w|t) \propto \frac{1}{2\sigma_n^2} t^T t + \frac{1}{\sigma_n^2} t^T \phi w - \frac{1}{2\sigma_n^2} w^T \phi^T w - \frac{1}{2\sigma_w^2} w^T w$$

$$\log p(w|t) \propto -\frac{1}{2} w^T \left(\frac{1}{\sigma_n^2} \phi^T \phi + \frac{1}{\sigma_w^2} I_Q \right) w + \frac{1}{\sigma_n^2} w^T \phi^T t$$

$$\text{luego } \bar{S}_n^{-1} = \frac{1}{\sigma_n^2} \phi^T \phi - \frac{1}{\sigma_w^2} I \phi$$

$$\bar{S}_n^{-1} m_n = \frac{1}{\sigma_n^2} \phi^T t$$

Ahora, con el criterio MAP, el problema es:

$$w_{\text{MAP}} = \arg\max_w p(w|t) = \arg\max_w (\log p(t|w) + \log p(w))$$

$$w^* = \bar{m}_n$$

Quitamos calcular la distribución predictiva de t^*

$$p(t^*|x^*, D) = \int p(t^*|x^*, w) p(w|D) dw$$

Por tanto, de la integral obtenemos:

$$p(t^*|x^*, D) = N(t^* | \phi(x^*)^T m_n, \phi(x^*)^T S_n \phi(x^*) + \sigma_n^2)$$

$$\mu_{t^*} = \phi(x^*)^T m_n$$

$$\sigma_{t^*}^2 = \phi(x^*)^T S_n \phi(x^*) + \sigma_n^2$$

Ahora expresamos el ruido sobre la función $f(x)$

$$f(x) \sim N(0, K(x, x'))$$

Podemos considerar

$$p(y|D) \sim N(0, K + \sigma_n^2 I)$$

$$\begin{pmatrix} y \\ t^* \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K + \sigma_n^2 I & K_* \\ K_*^T & K(x^*, x^*) \end{pmatrix} \right)$$

Ahora usamos las propiedades gaussianas multivariadas

$$\begin{pmatrix} u \\ v \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_u \\ \mu_v \end{pmatrix}, \begin{pmatrix} \Sigma_{uu} & \Sigma_{uc} \\ \Sigma_{cu} & \Sigma_{cc} \end{pmatrix} \right)$$

$$p(u|v) \sim N(\mu_u + \Sigma_{uc} \Sigma_{cc}^{-1} (v - \mu_v), \Sigma_{uu} - \Sigma_{uc} \Sigma_{cc}^{-1} \Sigma_{cu})$$

Aplicando esto a nuestro caso

$$p(t^*|x^*, D) \sim N(K_*^T (K + \sigma_n^2 I)^{-1} y, K(x^*, x^*) - K_*^T (K + \sigma_n^2 I)^{-1} K_*)$$

Simplificando

$$p(t^*|y^*, x^*) \sim N(K_*^T (K + \sigma_n^2 I)^{-1} y, K(x^*, x^*) - K_*^T (K + \sigma_n^2 I)^{-1} K_*)$$

Donde las expresiones son:

$$\mu_* = K_*^T (K + \sigma_n^2 I)^{-1} y \quad \rightarrow \text{Predicción}$$

$$\sigma_*^2 = K(x^*, x^*) - K_*^T (K + \sigma_n^2 I)^{-1} K_* \quad \rightarrow \text{Incertidumbre}$$

Nombre	Problema Optimización	Solución analítica	Explicación
LS - least Squares	$W_{LS} = \underset{w}{\operatorname{argmin}} \operatorname{RSS}(w) = \underset{w}{\operatorname{argmin}} \sum_{n=1}^N (t_n - \phi(x_n) \cdot w^T)^2$	Invertible $(\Phi^T \Phi)^{-1} \Phi^T \cdot t = w$	Busca minimizar la distancia entre los datos observados y el dato predicho
RLS - Regularized Least Squares	$\underset{w}{\operatorname{argmin}} \sum_{n=1}^N (t_n - \phi(x_n) \cdot w^T)^2 + \lambda \ w\ _2^2$	Pseudoinversa $w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T t$	Se penaliza la norma de los pesos de w , con el objetivo de que mi modelo no sea tan sensible a cambios pequeños en x . λ es un hiperparámetro y controla la "fuerza" de la penalización y garantiza la invertibilidad. Esto evita overfitting \rightarrow Para evitar que mi modelo se 'aprenda' los datos y compense con grandes pesos las variaciones
MLE - Maximum Likelihood Estimation	$\underset{w, \sigma^2}{\operatorname{argmax}} \log p(t \Phi, w, \sigma^2)$ $= \underset{w, \sigma^2}{\operatorname{argmax}} \log \left(\prod_{n=1}^N N(t_n \phi(x_n) w^T, \sigma_n^2) \right)$	$(\Phi^T \Phi)^{-1} \Phi^T \cdot t = w$	La verosimilitud es la probabilidad de haber observado exactamente esa secuencia de datos, dado cierta hipótesis. El principio de máxima verosimilitud establece que el mejor modelo es el que hace que los datos observados sean lo más probable posible. La solución coincide con la de mínimos cuadrados, pero ya no es 'quiero minimizar el error' ahora quiero que mis datos sean lo más probables posibles.
MAP - Maximum a posteriori.	$= \underset{w}{\operatorname{argmin}} \left(\frac{1}{2\sigma^2} \ t - \Phi w\ ^2 + \frac{\sigma_n^2}{\sigma_w^2} \ w\ ^2 \right)$	$w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T t$ $\lambda = \sigma_n^2 / \sigma_w^2$	Aparece el concepto de prior, el concepto de incorporar conocimiento, creencias o experiencia que ya se tiene sobre el fenómeno antes de ver cualquier dato. Luego se lleva a cabo el experimento y se actualiza la creencia. Ese prior es una distribución de probabilidad sobre w .
Bayesiano Lineal Gaussiano	$= \underset{w}{\operatorname{argmax}} p(w t) = \underset{w}{\operatorname{argmax}} p(t w) p(w)$	$p(w t) = N(w m_N, S_N)$ Parámetros actualizados $S_N^{-1} = \lambda I + \frac{1}{\sigma^2} \Phi^T \cdot \Phi$ $m_N = \frac{1}{\sigma^2} \cdot S_N \cdot \Phi^T \cdot t$	Aquí aparece la "confianza de cada predicción" porque aparece la varianza (incertidumbre). Quiere decir a fin de una distribución. Mi posterior tiene: media, varianza. No busca minimizar el error, busca actualizar una distribución de probabilidad sobre los parámetros (media y varianza)
Regresión Rígida kernel	$\underset{a}{\operatorname{argmin}} (\ t - K a\ ^2 + \frac{\lambda}{2} a^T K a)$ $a \rightarrow$ Son los pesos (w) pero en otro espacio de representación $w = a_1 \phi(x_1) + a_2 \phi(x_2) \dots a_N \phi(x_N)$ $a =$ Qué tanto x_1 aporta al modelo $K_{ij} = \phi(x_i)^T \cdot \phi(x_j)$	$\hat{a} = (K + \lambda I_N)^{-1} \cdot t$ Predicción $t(x_*) = K(x_*)^T \cdot \hat{a}$	El kernel es una función que mide qué tan parecidas son dos cosas, pero en un espacio transformado no en un espacio original $K = \Phi \Phi^T$. $K_{ij} = K(x_i, x_j)$ mide la similitud entre dos puntos de entrenamiento \uparrow Puntos \uparrow Kernel Su principio parte de mínimos cuadrados regularizados. Es un modelo que extiende la regresión lineal a contextos no lineales. En su espacio original no se ajustan linealmente, pero se pueden volver lineales a un espacio de características mayor. Tiene un enfoque determinista.
Procesos Gaussianos	$= \underset{w}{\operatorname{argmax}} p(w t) = \underset{w}{\operatorname{argmax}} p(t w) p(w)$ $= \underset{w}{\operatorname{min}} \left(\frac{1}{2\sigma_n^2} \ t - \Phi w\ _2^2 + \frac{1}{2} w^T \Sigma_p^{-1} w \right)$	Solución Optima $\hat{a} = (K + \sigma_n^2 I)^{-1} \cdot t$ Media posterior. $\mathcal{M}_*(x_*) = K(x_*)^T \cdot (K + \sigma_n^2 I)^{-1} \cdot t$ Varianza asociada $\sigma_*^2(x_*) = K(x_*, x_*) - K(x_*)^T (K + \sigma_n^2 I)^{-1} K(x_*)$	Un proceso gaussiano es una distribución sobre funciones. Me da la varianza que se asocia a la incertidumbre del modelo. Lo que asume que los valores observados como los que quiero predecir siguen una distribución normal multivariada. $\begin{bmatrix} t \\ t_* \end{bmatrix} \sim N \left(0, \begin{bmatrix} K + \sigma_n^2 I & K(x_*) \\ K(x_*) & K(x_*, x_*) \end{bmatrix} \right)$