

Parcial #1

Teoría de Aprendizaje de Máquina

Yonnier Alexander Muñoz Salazar

2.1. Sea el modelo de regresión $x_n = \phi(x_n)W^T + \eta_n$
Con: $\{x_n \in \mathbb{R}, \mathcal{X}_n \in \mathbb{R}^P\}_n^N$, $W \in \mathbb{R}^Q$, $\phi: \mathbb{R}^P \rightarrow \mathbb{R}^Q$
 $Q \geq P$, $\eta_n \sim N(\eta_n | \vec{0}, \sigma_n^2)$

Presentar el problema de optimización y solución:

◆ Mínimos Cuadrados

Se tiene el problema de optimización:

$$\hat{W} = \underset{W}{\operatorname{argmin}} \|x_n - f(\phi | W)\|_2^2$$

Se resuelve para hallar W^T . De la norma se tiene:

$$\|x_n - f(\phi | W^T)\|_2^2 = \langle x_n - f(\phi | W^T), x_n - f(\phi | W^T) \rangle$$

Se supone ϕ como datos transformados con regresión no lineal.

Se distribuye el producto interno

$$\langle x_n - f(\phi | W^T), x_n - f(\phi | W^T) \rangle$$

$$= \langle x_n - \phi W^T, x_n - \phi W^T \rangle$$

$$= x_n^T x_n - x_n^T \phi W^T - (\phi W^T)^T x_n + (\phi W^T)^T \phi W^T$$

$$= x_n^T x_n - 2x_n^T \phi W^T + W \phi^T \phi W^T$$

Se deriva respecto a W^T y se iguala a cero

$$\frac{d}{dW} (x_n^T x_n - 2x_n^T \phi W^T + W \phi^T \phi W^T)$$

$$= 0 - 2x_n^T \phi + 2\phi^T \phi W^T = 0$$

Se despeja W^T

$$\phi \phi^T \phi W^T = \phi x_n^T \phi$$

$$\boxed{W^T = x_n^T (\phi^T \phi)^{-1} \phi}$$

Finalmente, Se tiene el predict para MC como:

$$\hat{x}_{new} = \phi(x_{new}) \hat{W}^T$$

◆ Mínimos Cuadrados Regularizados

Se tiene el problema de optimización

$$\hat{W} = \underset{W^T}{\operatorname{argmin}} \frac{1}{N} \|x_n - F(\phi(x_n) | W)\|_2^2 + \lambda \|W\|_2^2$$

Donde $\lambda \in \mathbb{R}$, λ es un hiperparametro de regularización

Se distribuye el producto interno del primer termino

$$\|x_n - F(\phi_n | W^T)\|_2^2 = [x_n^T x - 2x_n^T \phi W + W^T \phi^T \phi W] \frac{1}{N}$$

Del segundo termino se tiene:

$$\lambda \|W\|_2^2 = \lambda \langle W, W \rangle = \lambda W^T W$$

Se deriva ambos terminos respecto a W y se iguala a 0

$$\frac{d}{dW} \|x_n - \phi W\|_2^2 = [-2 \phi^T x_n + 2 \phi^T \phi W] \frac{1}{N}$$

$$\frac{d}{dW} \lambda \|W\|_2^2 = 2 \lambda W$$

$$\Rightarrow -\frac{2}{N} \phi^T x_n + \frac{2}{N} \phi^T \phi W + 2 \lambda W = 0$$

Se multiplica a ambos lados por $N/2$ y $(\phi^T \phi + \lambda N I)^{-1}$

$$\frac{N}{2} \left[-\frac{2}{N} \phi^T x_n + \frac{2}{N} \phi^T \phi W + 2 \lambda W \right] = 0 \left[\frac{N}{2} \right]$$

$$= (\cancel{\phi^T \phi} + \lambda N I)^{-1} (\cancel{\phi^T \phi} + \lambda N I) W = (\phi^T \phi + \lambda N I)^{-1} \phi^T x_n$$

Finalmente Se tiene:

$$\begin{aligned} W &= (\phi^T \phi + \lambda N I)^{-1} \phi^T x_n \\ \hat{W}^T &= x_n^T (\phi^T \phi + \lambda N I)^{-1} \phi \end{aligned}$$

Predict para MCB: $\hat{x}_{new} = \phi(x_{new}) \hat{W}^T$

◆ Máxima Verosimilitud

En el problema de optimización se agrega $WGN = \eta$ y se utiliza la función log para simplificar cálculos

$$\hat{W} = \operatorname{argmax}_W \log(P(x_n | F(X|W); \eta))$$

La información del ruido se toma como:

$$\hat{W} = \operatorname{argmax}_W \log\left(\prod_{n=1}^N N(x_n | F(X, W), \sigma_\eta^2)\right)$$

Con Φ como los datos transformados, se tiene la forma de la probabilidad como:

$$\log\left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_\eta^2}} \exp\left(-\frac{\|x_n - \Phi W\|_2^2}{2\sigma_\eta^2}\right)\right)$$

Se utilizan propiedades de log para simplificar la expresión

$$\begin{aligned} & \log\left(\prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma_\eta^2}}\right) + \log\left(\prod_{n=1}^N \exp\left(-\frac{\|x_n - \Phi W\|_2^2}{2\sigma_\eta^2}\right)\right) \\ &= \log_N\left(\frac{1}{(2\pi\sigma_\eta^2)^{N/2}}\right) + \log_N\left(\exp\left(-\sum_{n=1}^N \frac{\|x_n - \Phi W\|_2^2}{2\sigma_\eta^2}\right)\right) \\ &= \log_N\left((2\pi\sigma_\eta^2)^{N/2}\right) - \sum_{n=1}^N \frac{\|x_n - \Phi W\|_2^2}{2\sigma_\eta^2} \end{aligned}$$

Se obtiene

$$-\frac{N}{2} [\log(2\pi) + \log(\sigma_\eta^2)] - \frac{1}{2\sigma_\eta^2} \|x_n - \Phi W\|_2^2$$

Se deriva respecto a W y se iguala a cero, de MC se tiene la derivada del segundo término y la derivada del primer término es igual a cero. Finalmente, se tiene:

$$\begin{aligned} & \frac{d}{dW} \left(-\frac{N}{2} [\log(2\pi) + \log(\sigma_\eta^2)] - \frac{1}{2\sigma_\eta^2} \|x_n - \Phi W\|_2^2 \right) \\ &= 0 - \frac{1}{2\sigma_\eta^2} [-2\Phi^T x_n + 2\Phi^T \Phi W] \end{aligned}$$

Despejando W , se tiene:

$$\boxed{\hat{W} = (\Phi^T \Phi)^{-1} \Phi^T x_n}$$

Para encontrar σ_n^2 , se repite el procedimiento anterior, derivando y despejando para σ_n^2

$$\frac{d}{d\sigma_n^2} \left[-\frac{N}{2} (\log(2\pi) + \log(\sigma_n^2)) \right] - \frac{d}{d\sigma_n^2} \left[\frac{1}{2\sigma_n^2} \|x_n - \Phi w\|_2^2 \right]$$

$$= -\frac{N}{2\sigma_n^2} + \frac{1}{2(\sigma_n^2)^2} \|x_n - \Phi w\|_2^2 = 0$$

Se multiplica a ambos lados por $(\sigma_n^2)^2$

$$= -\frac{N}{2\sigma_n^2} (\sigma_n^2)^2 + \frac{1}{2(\sigma_n^2)^2} (\sigma_n^2)^2 \|x_n - \Phi w\|_2^2 = 0 (\sigma_n^2)^2$$

$$= -\frac{N}{2} \sigma_n^2 + \frac{1}{2} [-2\Phi^T x_n + 2\Phi^T \Phi w] = 0$$

Finalmente, se tiene:

$$\boxed{\sigma_n^2 = \frac{1}{N} [-2\Phi^T x_n + 2\Phi^T \Phi w]}$$

◆ Máximo A-posteriori

Se agrega un prior con probabilidad

$$P(w) = N(w | 0, \sigma_w^2) \rightarrow w \in \mathbb{R}^Q$$

Se tiene la optimización como:

$$\hat{w} = \underset{w}{\operatorname{Argmax}} \log(P(w) | x_n, \Phi, \sigma_n^2)$$

Se desarrolla la expresión

$$\log(P(w) | x_n, \Phi, \sigma_n^2) = \log\left(\prod_{n=1}^N (x_n | \Phi(x_n) | w, \sigma_n^2)\right) \prod_{q=1}^Q N(w_q | 0, \sigma_w^2)$$

$$= \left\{ -\frac{N}{2} [\log(2\pi) + \log(\sigma_n^2)] - \frac{1}{2\sigma_n^2} \|x_n - \Phi w\|_2^2 \right\} \Rightarrow A$$

$$+ \log\left(\prod_{q=1}^Q \left(\frac{1}{\sqrt{2\pi\sigma_w^2}}\right) \exp\left(-\frac{\|w_q - 0\|_2^2}{2\sigma_w^2}\right)\right)$$

$$= A + \log_2\left(\prod_{q=1}^Q \frac{1}{(2\pi\sigma_w^2)^{Q/2}}\right) + \log\left(\exp\left(-\sum_{q=1}^Q \frac{\|w_q\|_2^2}{2\sigma_w^2}\right)\right)$$

$$= A - \frac{Q}{2} [\log(2\pi) + \log(\sigma_w^2)] - \sum_{q=1}^Q \frac{\|w_q\|_2^2}{2\sigma_w^2}$$

$$= -\frac{N}{2} \left[\log(2\pi) + \log(\sigma_n^2) \right] - \frac{1}{2\sigma_n^2} \|x - \Phi w\|_2^2 - \frac{Q}{2} \left[\log(2\pi) + \log(\sigma_w^2) \right] - \frac{1}{2\sigma_w^2} \|w\|_2^2$$

Ahora, se tiene que:

$$= \text{Max} - \left[\frac{2\sigma_n^2}{2\sigma_n^2} \|x_n - \Phi w\|_2^2 + \frac{2\sigma_n^2}{2\sigma_w^2} \|w\|_2^2 \right] + c + c$$

$$= \text{min} \left[\|x_n - \Phi w\|_2^2 + \frac{\sigma_n^2}{\sigma_w^2} \|w\|_2^2 + c + c \right]$$

Finalmente \hat{w} esta dado por:

$$\hat{w} = (\Phi^T \Phi + \frac{\sigma_n^2}{\sigma_w^2} I)^{-1} \Phi^T x_n$$

Bayesiano lineal Gaussiano

Se define el prior con probabilidad de la forma:

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

Donde, $P(w)$ incluye una matriz de varianza anisotrópica

La probabilidad de WGN esta dada como:

$$P(\eta) = N(\eta | 0, B^{-1})$$

Donde, $B^{-1} = \sigma^2 I$, isotrópica $\rightarrow B \in \mathbb{R}^+$

El problema de optimización por Bayes es:

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

Donde: $P(w|x) \Rightarrow$ Posterior

$P(x|w) \Rightarrow$ Verosimilitud

$P(w) \Rightarrow$ Prior

$P(x) \Rightarrow$ Evidencia

Despejando el ruido de la verosimilitud $P(x|w, \eta)$

Se tiene: $\eta = t_n - w^T \phi(x_n) \sim N(t_n - w^T \phi(x_n) | 0, B^{-1})$

El termino $t_n - w^T \phi(x_n)$ traslada la media del termino asi:

$$P(x|w) = N(x | \phi w, B^{-1} I)$$

Para el posterior, se obtiene $M_{w|x}$ y $\Sigma_{w|x}$ a partir de la matriz de precision y completando cuadrados, se tiene:

$$\begin{aligned} M_{w|x} &= M_N = (S_0^{-1} + \phi^T B I \phi)^{-1} [\phi^T B I (t - 0) + S_0^{-1} M_0] \\ &= S_N (S_0^{-1} M_0 + \beta \phi^T x) \end{aligned}$$

y,

$$\Sigma_{w|x} = S_N$$

$$S_N = (S_0^{-1} + \beta \phi^T \phi)^{-1}; S_N^{-1} = S_0^{-1} + \beta \phi^T \phi$$

Finalmente se tiene la predictiva del modelo de la forma:
Como,

$$P(x_{\text{new}}|w) = N(x_{\text{new}} | w^T \phi(x_n), B^{-1})$$

y,

$$P(w|x) = N(w | m_N, S_N)$$

Con el operador convolucion, se tiene:

$$P(x_{\text{new}} | x, w) = \int P(x_{\text{new}} | w) P(w | x) dw$$

Queda:

$$P(x_{\text{new}} | x, w) = N(x_{\text{new}} | M_N^T \phi(x_*) \beta^{-1} + \phi^T(x_*) S_N \phi(x_*))$$

Sabiendo que:

$$M_N^T \phi(x_*) = \phi^T(x_*) M_N$$

La predictiva final es la predictiva de cada dato nuevo ponderada respecto a la esperanza del posterior

$$\hat{P}(x_{\text{new}} | w) = E_{w|x} \{ P(x_{\text{new}} | w) \}$$