

Parcial #1

Sea el modelo de regresión $t_n = \phi(x_n)w^T + n_n$, 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 $n_n \sim N(n_n | 0, \sigma_n^2)$

- mínimos cuadrados

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

$w \rightarrow$ parámetros a estimar

$\phi(x_n) \rightarrow$ transformación de las variables de entrada

$n_n \rightarrow$ ruido gaussiano media 0 varianza constante

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

vector de predicción

$$\Phi w$$

$$t = [t_1, t_2, \dots, t_N]^T \in \mathbb{R}^N$$

vector de errores

$$t - \Phi w$$

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix} \in \mathbb{R}^{N \times q}$$

función de costo

$$J(w) = \|t - \Phi w\|^2$$

$$w \in \mathbb{R}^q$$

Derivamos función de costo

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

se deriva con respecto a w

$$\frac{\partial J}{\partial w} = -2\Phi^T t + 2\Phi^T \Phi w$$

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

$$\Phi^T \Phi w = \Phi^T t$$

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

$\Phi^T \Phi$ matriz de correlaciones entre las variables de entrada

$\Phi^T t$: correlacion entre variable de entrada y objetivo

$(\Phi^T \Phi)^{-1} \Phi^T t$: vector de peso que minimiza la suma de los errores al cuadrado

- Minimos cuadrados regularizados

se agrega un termino mas que penaliza la magnitud de w para evitar sobreajuste

Funcion de costo con Regularizacion

$$J(w) = \sum_{n=1}^N (t_n - \phi(x_n)^T w)^2 + \lambda \|w\|^2$$

pasamos a matriz

t = vector columna de todas las salidas t_n

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$$

Φ = matriz de disenio con todas las transformaciones $\phi(x_n)^T$

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \phi(x_2)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}$$

w = vector de peso

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_q \end{bmatrix}$$

$\hat{t} = \phi w$ valor predicho para todos los datos

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

$$\|a\|^2 = \sum_{i=1}^N a_i^2$$

entonces

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

$$\|w\|^2 = w^T w \quad \lambda \rightarrow \lambda \|w\|^2 = \lambda w^T w$$

$$J(w) = \|t - \phi w\|^2 + \lambda \|w\|^2$$

$\phi \in \mathbb{R}^{N \times q}$ matriz de diseño

$t \in \mathbb{R}^N$ vector objetivo

$\lambda > 0$ parametro regularizacion controla el trade-off entre complejidad del modelo y ajuste

expandimos la funcion de costo

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

primer termino $= t^T t - 2t^T \phi w + w^T \phi^T \phi w$

segundo termino $= \lambda w^T w$

funcion completa

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

Derivamos respecto a w

derivada de $t^T t$ respecto a w es 0

derivada de $-2t^T \phi w$ es $-2\phi^T t$

derivada de $w^T \phi^T \phi w$ es $2\phi^T \phi w$

derivada $\lambda w^T w$ es $2\lambda w$

$$\frac{\partial J}{\partial w} = -2\Phi^T t + 2\Phi^T \Phi w + 2\lambda w$$

nos daría

$$\frac{\partial J}{\partial w} = -2\Phi^T t + 2\Phi^T \Phi w + 2\lambda w$$

despejamos

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

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

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

$\Phi^T \Phi$ matriz de correlación entre las variables de entrada

λI término de regularización, controla la magnitud del coeficiente w

$\lambda >$ mayor la penalización \rightarrow Se reduce las magnitudes de w

- máxima verosimilitud

función de verosimilitud

$$p(t|w, \sigma^2) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (t_n - \Phi(x_n)^T w)^2\right)$$

pasamos a log

$$\log p(t|w, \sigma^2) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \Phi(x_n)^T w)^2$$

$$J(w) = \sum_{n=1}^N (t_n - \Phi(x_n)^T w)^2$$

$$\sum_{n=1}^N (t_n - \Phi(x_n)^T w)^2 = \|t - \Phi w\|^2$$

① matriz de características $N \times Q$

t vector $N \times 1$

w vector $Q \times 1$

Expandimos el error cuadrático medio

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

Derivamos respecto a w

$$\frac{\partial}{\partial w} (t^T t) = 0$$

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

$$\frac{\partial}{\partial w} (w^T \phi^T \phi w) = 2\phi^T \phi w$$

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

despejamos

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

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

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

- máximo a-posteriori

$$P(t|w, \sigma^2) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (t_n - \phi(x_n)^T w)^2\right)$$

$$P(w) = \mathcal{N}(0, \alpha^{-1} I)$$

$$P(w) = \frac{1}{(2\pi\alpha^{-1})^{Q/2}} \exp\left(-\frac{\alpha}{2} w^T w\right)$$

α = parámetro de precisión

Posterior Sin Constante

$$p(w/t) \propto p(t/w) \cdot p(w)$$

$$\log p(w/t) \propto \log p(t/w) + \log p(w)$$

Verosimilitud

$$\log p(t/w) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|t - \Phi w\|^2$$

prior

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

$$\log p(w/t) \propto -\frac{1}{2\sigma^2} \|t - \Phi w\|^2 - \frac{\alpha}{2} w^T w$$

Derivamos

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

$$t^T t - 2t^T \Phi w + w^T \Phi^T \Phi w$$

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

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

$$-\frac{1}{2\sigma^2} w^T \Phi^T \Phi w = -\frac{1}{2} \Phi^T \Phi w$$

$$-\frac{\alpha}{2} w^T w = -\alpha w$$

Despejamos

$$\frac{1}{\sigma^2} \Phi^T t \Rightarrow \frac{1}{\sigma^2} \Phi^T \Phi w + \alpha w = 0$$

$$\left(\frac{1}{\sigma^2} \Phi^T \Phi + \alpha I \right) w = \frac{1}{\sigma^2} \Phi^T t$$

$$(\Phi^T \Phi + \sigma^2 \alpha I) w = \Phi^T t$$

$$x = \sigma^2 \alpha$$

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

w se agrega un termino de penalizacion cuadratica

$\lambda > 0$, mas se penaliza los valores grandes de w

- Bayesiano con modelo lineal gaussiano

funcion de verosimilitud

$$p(t|w) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(t_n - \Phi(x_n)^T w)^2}{2\sigma^2}\right)$$

$$p(t|w) = N(\Phi w, \sigma^2 I)$$

Φ matriz $N \times Q$ con filas $\Phi(x_n)^T$

t = vector $N \times 1$

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

$$m_0 = 0$$

$$S_0 = \alpha^{-1} I$$

$$p(w|t) = p(t|w) \cdot p(w)$$

$$p(w|t) = N(m_N, S_N)$$

$$S_N^{-1} = S_0^{-1} + \frac{1}{\sigma^2} \Phi^T \Phi$$

$$m_N = S_N \left(S_0^{-1} m_0 + \frac{1}{\sigma^2} \Phi^T t \right)$$

$$m_0 = 0$$

$$S_0 = \alpha^{-1} I$$

$$S_0^{-1} = \alpha I$$

$$S_N^{-1} = \alpha I + \frac{1}{\sigma^2} \Phi^T \Phi$$

$$m_N = S_N \left(\frac{1}{\sigma^2} \Phi^T t \right)$$

para un nuevo punto

$$p(t_x | t) = N(\Phi(x_x)^T m_N, \sigma_x^2)$$

Varianza

$$\sigma_x^2 = \sigma^2 + \Phi(x_x)^T S_N \Phi(x_x)$$

- Regression ridge kernel

Regression lineal ordinaria

$$\hat{w} = \arg \min \|t - \Phi w\|^2$$

Φ matriz característica

t vector

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

Penalización sobre el tamaño de los pesos

$$\hat{w} = \arg \min_w [\|t - \Phi w\|^2 + \lambda \|w\|^2]$$

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

$$\hat{w} = \Phi^T a$$

$$\hat{w} = \Phi^T a$$

$$\hat{t} = \Phi \hat{w} = \Phi \Phi^T a$$

Resolvemos para a

$$\|t - \Phi \Phi^T a\|^2 + \lambda a^T \Phi \Phi^T a$$

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

$$a = (C \Phi \Phi^T + \lambda I)^{-1} e$$

introducimos Kernel

$$K(x_i, x_j) = \Phi(x_i) \Phi(x_j)$$

$$K = \Phi \Phi^T$$

$$a = (K + \lambda I)^{-1} e$$

para nuevos puntos

$$\hat{f}_x = \sum_{n=1}^N a_n K(x_n, x)$$

en vectorial

$$\hat{f}_x = K_x a$$

- procesos Gaussianos

$$f(x) \sim GP(m(x), k(x, x'))$$

$m(x)$ = media

$K(x, x')$ = función de covarianza

matriz de covarianza puntos de entrenamiento

$$K = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_N) \\ \vdots & \ddots & \vdots \\ K(x_N, x_1) & \dots & K(x_N, x_N) \end{bmatrix}$$

vector de covarianza

$$K_x = \begin{bmatrix} K(x_x, x_1) \\ \vdots \\ K(x_x, x_N) \end{bmatrix}$$

varianza nuevo punto

$$K_{xx} = K(x_x, x_x)$$

$$C = K + \sigma_n^2 I$$

media

$$\mu_* = K_*^T C^{-1} t$$

t = vector de observacion

varianza

$$\sigma_*^2 = K_{xx} - K_*^T C^{-1} K_*$$

Prediccion x_*

$$p(t_* | x_*, x, t) = N(\mu_*, \sigma_*^2)$$

- la regresion lineal por maximo verosimilitud predice valores puntuales sin regularizacion ni incertidumbre
- la regresion ridge añade regularizacion λ
- la regresion lineal bayesiana incorpora un prior sobre los pesos y permite obtener media y varianza en las predicciones con kernel. la regresion ridge kernel mantiene la regularizacion pero sigue siendo determinista
- la regresion bayesiana kernel añade prior sobre funciones y predice media y varianza
- procesos gaussianos modela toda la funcion como una distribucion gaussiana sobre funciones, predice media y varianza y maneja incertidumbre natural