

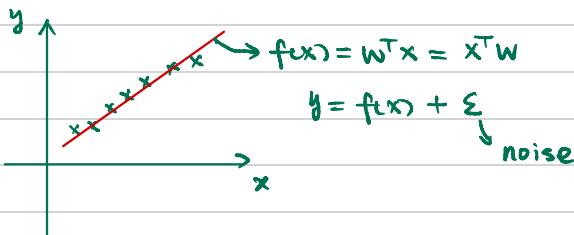


线性回归

Data = $\{(x_i, y_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^P$, $y_i \in \mathbb{R}$

$$X = (x_1 \ x_2 \ \dots \ x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & & & \\ x_{N1} & x_{N2} & \dots & x_{NP} \end{pmatrix}_{N \times P}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

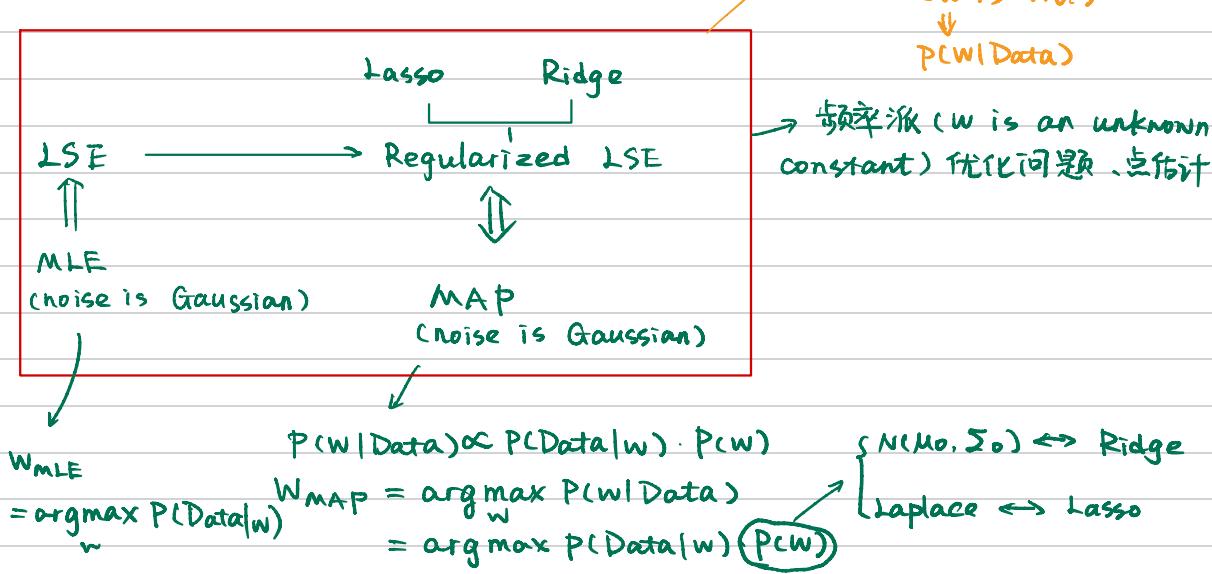


x, y, ε are r.v.

$x \in \mathbb{R}^P$, $y \in \mathbb{R}^1$, $\varepsilon \sim N(0, \sigma^2)$

Bayesian method → 贝叶斯法

(w is r.v.)
 \downarrow
 $P(w|Data)$



Bayesian Linear regression

Data: $\{f(x_i, y_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^P$, $y_i \in \mathbb{R}$

$$X = (x_1, x_2, \dots, x_N)^T$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

Model:

$$\begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \varepsilon \end{cases} \quad x, y, \varepsilon \text{ are r.v.}$$

$\varepsilon \sim N(0, \sigma^2)$

Bayesian Method

$$\begin{cases} \text{Inference: posterior}(w) \\ \text{Prediction: } x^* \rightarrow y^* \end{cases}$$

$$P(w, Y|X) = P(Y|w, X) \cdot P(w|X)$$

Inference $P(w|Data)$

$$P(w|Data) = P(w|X, Y) = \frac{P(w, Y|X)}{P(Y|X)} = \frac{P(Y|w, X) \cdot P(w)}{\int P(Y|w, X) \cdot P(w) dw}$$

$$P(Y|w, X) = \prod_{i=1}^N P(y_i|w, x_i) = \prod_{i=1}^N N(y_i|w^T x_i, \sigma^2)$$

$$P(Y|X, w) = N(w^T x, \sigma^2)$$

likelihood \uparrow prior \nearrow

$$P(w) = N(0, \Sigma_p)$$

$$P(w|Data) \propto P(Y|w, X) \cdot P(w)$$

Gaussian Gaussian Gaussian

共轭: Gaussian 分布是自共轭的

$$\propto \prod_{i=1}^N N(y_i|w^T x_i, \sigma^2) \cdot N(0, \Sigma_p)$$

$$\rightarrow N(\mu_w, \Sigma_w) \quad \mu_w?, \Sigma_w?$$

$$\text{likelihood} = P(Y|X, W) = \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right\}$$

$$= \frac{1}{2\pi^{\frac{N}{2}} \sigma^N} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^T x_i)^2 \right\}$$

$$\sum_{i=1}^N (y_i - w^T x_i)^2 = (y_1 - w^T x_1, y_2 - w^T x_2, \dots, y_N - w^T x_N) \begin{pmatrix} y_1 - w^T x_1 \\ \vdots \\ y_N - w^T x_N \end{pmatrix}$$

$$(y_1, y_2, \dots, y_N) - w^T(x_1, x_2, \dots, x_N) \\ (Y^T - w^T x^T) (Y - XW)$$

$$\frac{1}{(2\pi)^{\frac{N}{2}} \sigma^N} \exp \left\{ -\frac{1}{2\sigma^2} (Y - XW)^T \Sigma^{-1} (Y - XW) \right\}$$

$$= N(XW, \sigma^2 I)$$

$$P(W|\text{Data}) \propto N(XW, \sigma^2 I) \cdot N(0, \Sigma_p)$$

$$\propto \exp \left\{ -\frac{1}{2} (Y - XW)^T \sigma^{-2} I (Y - XW) \right\} \exp \left\{ -\frac{1}{2} w^T \Sigma_p^{-1} w \right\}$$

$$\propto \exp \left\{ -\frac{1}{2\sigma^2} (Y^T - w^T x^T) (Y - XW) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\}$$

$$= \exp \left\{ -\frac{1}{2\sigma^2} (Y^T Y - 2Y^T XW + w^T x^T XW) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\}$$

$$P(X) = N(\mu, \Sigma)$$

$$\text{指數部分: } \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$-\frac{1}{2} ((x^T \Sigma^{-1} - \mu^T \Sigma^{-1})(x - \mu))$$

$$=-\frac{1}{2} (x^T \Sigma^{-1} x - 2\mu^T \Sigma^{-1} x + \Delta)$$

$$\Rightarrow \text{R 项: } -\frac{1}{2\sigma^2} w^T x^T XW - \frac{1}{2} w^T \Sigma_p^{-1} w = -\frac{1}{2} (w^T (\underbrace{\sigma^{-2} X^T X + \Sigma_p^{-1}}_{\Sigma_n^{-1} = A} W) w)$$

$$\Rightarrow \text{R 项: } + \frac{1}{2\sigma^2} (2) Y^T XW = \underbrace{\sigma^{-2} Y^T XW}_{\mu_W^T A}$$

$$\mu_W^T A = \mu_W^T \Sigma_W^{-1} = \sigma^{-2} Y^T X$$

$$A \mu_W = \sigma^{-2} X^T Y$$

$$\mu_W = \sigma^{-2} A^{-1} X^T Y$$

$$P(W | \text{Data}) = N(\mu_w, \Sigma_w)$$

$$\mu_w = \sigma^{-2} A^{-1} X^T Y$$

$$\Sigma_w = A^{-1}$$

$$A = \sigma^{-2} X^T X + \Sigma_p^{-1}$$

Prediction:

Given x^*, y^*

noise-free

$$\textcircled{1} \quad f(x^*) \Rightarrow f(x^*) = x^{*T} \underline{w} \rightarrow P(w | \text{Data}) = N(\mu_w, \Sigma_w)$$

$$P(f(x^*) | \text{Data}, X^*) = N(x^{*T} \mu_w, x^{*T} \Sigma_w x^*)$$

$$x^{*T} w \sim N(x^{*T} \mu_w, x^{*T} \Sigma_w x^*)$$

$$\textcircled{2} \quad \underbrace{y^*}_{\sim N(0, \sigma^2)} = f(x^*) + \varepsilon$$

$$P(y^* | X^*, \text{Data}) = N(x^{*T} \mu_w, x^{*T} \Sigma_w x^* + \sigma^2)$$

Gaussian Process

distribution 随机过程

Gaussian distribution \rightarrow multivariate Gaussian distribution

(一维)

$$p(x) = N(\mu, \sigma^2)$$

时间/空间

高斯过程是定义在连续域上的无限多个高斯随机变量所组成的随机过程

Gaussian Network
 $x \in \mathbb{R}^p$

$$p(x) = N(\mu, \Sigma), \Sigma_{p \times p} \Rightarrow p < +\infty$$

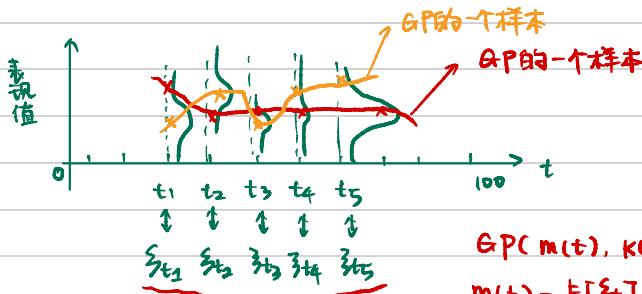
无限维 Gaussian Distribution

$\{\xi_t\}_{t \in T}$, $T \rightarrow$ 连续域, if $\forall n \in \mathbb{N}^+, t_1, t_2, \dots, t_n \in T$

Gaussian Process

s.t. $\{\xi_{t_1}, \xi_{t_2}, \dots, \xi_{t_n}\} \stackrel{\text{def}}{=} \xi_{t_1 \dots t_n} \sim N(\mu_{t_1 \dots t_n}, \Sigma_{t_1 \dots t_n})$. 那么

$\{\xi_t\}_{t \in T}$ 就是一个 Gaussian process



$$GP(m(t), k(s, t))$$

$$m(t) = E[\xi_t] \rightarrow \text{mean function}$$

$$k(s, t) = E[(\xi_s - m(s))(\xi_t - m(t))^T] \downarrow$$

kernel & covariance function

人的一生: $[0, 100]$ $N(\mu, \Sigma)$

$$t \in [0, 100], \xi_t \in N(\mu_t, \sigma_t^2)$$

这个人从 t 时刻他的表现 ($0 \sim 100$ 分)

$t=0$, 他的一生基本已定, $t>0$, 每一个时刻, μ_t, σ_t^2 是已经确定

Gaussian Process Regression (GPR)

Recall Bayesian Linear Regression

$$\textcircled{1} \quad P(w | \text{Data}) = N(w | \mu_w, \Sigma_w)$$

$$\left\{ \begin{array}{l} \mu_w = \sigma^{-2} A^{-1} X^T Y \\ \Sigma_w = A^{-1} \end{array} \right.$$

$$A = \sigma^{-2} X^T X + \Sigma_p^{-1}$$

\textcircled{2} Given x^*

$$P(f(x^*) | \text{Data}, x^*) = N(x^{*T} \mu_w, x^{*T} \Sigma_w x^*)$$

$$P(y^* | \text{Data}, x^*) = N(x^* \mu_w, x^{*T} \Sigma_w x^* + \sigma^2)$$

$$\phi: x \mapsto z \quad x \in \mathbb{R}^p, z \in \mathbb{R}^q, z = \phi(x)$$

Nonlinear $\rightarrow \left\{ \begin{array}{l} \textcircled{1} \text{ Non-linear Transformation} \\ \textcircled{2} \text{ Bayesian LR} \end{array} \right\} \xrightarrow{\text{kernel}} \text{kernel BLR}$

Noise-free:

$$f(x^*) | X, Y, x^* \sim N(x^{*T} (\sigma^{-2} A^{-1} X^T Y), x^{*T} A^{-1} x^*)$$

$$A = \sigma^{-2} X^T X + \Sigma_p^{-1}$$

$$\text{If } \phi: x \mapsto z, x \in \mathbb{R}^p, z = \phi(x) \in \mathbb{R}^q, q > p$$

$$\text{Define: } \Phi = \phi(X) = (\phi(x_1) \ \phi(x_2) \ \dots \ \phi(x_N))^T \quad N \times q$$

$$\Rightarrow f(x) = \phi^T w$$

$$f(x^*) | X, Y, x^* \sim N(\sigma^{-2} \phi(x^*)^T A^{-1} \Phi^T Y, \phi(x^*)^T A^{-1} \phi(x^*))$$

$$A = \sigma^{-2} \Phi^T \Phi + \Sigma_p^{-1}$$

compute $A^{-1} \rightarrow \text{woodbury formula}$

$$(A + UCV)^{-1} = A^{-1} - A^{-1} U (C^{-1} + V A^{-1} U)^{-1} V A^{-1}$$

$$A \Sigma_p \Phi^T = \sigma^{-2} \Phi^T \Phi \Sigma_p \Phi^T + \Phi^T$$

$$= \sigma^{-2} \Phi^T (k + \sigma^2 I) \quad \Rightarrow k = \Phi^T \Sigma_p \Phi$$

$$\Sigma_p \bar{\Phi}^T = \sigma^{-2} A^{-1} \bar{\Phi}^T (k + \sigma^2 I)$$

$$\sigma^{-2} A^{-1} \bar{\Phi}^T = \Sigma_p \bar{\Phi}^T (k + \sigma^2 I)^{-1}$$

$$\underbrace{\sigma^{-2} \phi(x^*)^T A^{-1} \bar{\Phi}^T}_\text{Expectation} Y = \phi(x^*) \Sigma_p \bar{\Phi}^T (k + \sigma^2 I)^{-1} Y$$

Expectation

$$\Rightarrow f(x^* | x, Y, X^*) \sim N(\phi(x^*)^T \Sigma_p \bar{\Phi}^T (k + \sigma^2 I)^{-1} Y, \phi(x^*)^T \Sigma_p \phi(x^*) - \phi(x^*)^T \Sigma_p \bar{\Phi}^T (k + \sigma^2 I)^{-1} \bar{\Phi} \Sigma_p \phi(x^*))$$

$$k = \bar{\Phi} \Sigma_p \bar{\Phi}^T$$

$$\phi(x^*)^T \Sigma_p \bar{\Phi}^T$$

$$\phi(x^*)^T \Sigma_p \bar{\Phi}^T$$

$$\bar{\Phi} \Sigma_p \phi(x^*)$$

$$\phi(x^*)^T \Sigma_p \phi(x^*)$$

$$\bar{\Phi} = \phi(X) = (\phi(x_1), \dots, \phi(x_N))^T$$

$$k(x, x') = \phi^T(x) \Sigma_p \phi(x') \leftrightarrow \text{kernel function}$$

$$\Sigma_p: \text{positive definite}, \Sigma_p = (\Sigma_p^{\frac{1}{2}})^2$$

$$k(x, x') = \phi^T(x) \Sigma_p^{\frac{1}{2}} \Sigma_p^{\frac{1}{2}} \phi(x') = (\Sigma_p^{\frac{1}{2}} \phi(x))^T \cdot \Sigma^{\frac{1}{2}} \phi(x') = \langle \psi(x), \psi(x') \rangle$$

$$\psi(x) = \Sigma_p^{\frac{1}{2}} \phi(x)$$

Bayesian Linear Regression + kernel trick (Non-linear Transformation
Inner product)

$$\text{GPR} \rightarrow \begin{cases} \textcircled{1} \text{ weight-space view } & \left\{ \begin{array}{l} f(x) = \phi^T(x) w \\ y = f(x) + \epsilon \end{array} \right. \\ \textcircled{2} \text{ function-space view } & f(x) \end{cases}$$

$f(x)$ is r.v.

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