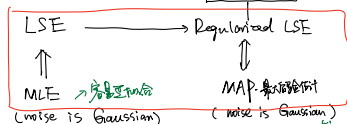
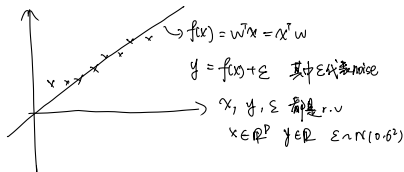


线性回归

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

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

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



→ 岭回归方法 (w is unknown constant)
↓
Hickin 问题
↓
点估计

→ Bayesian Method
↓
贝叶斯方法
↓
 w 是随机变量 (r.v.)
↓
不是点估计, 而是分布 ($w | \text{Data}$)

MLE = argmax_w P(Data|w)
MAP = argmax_w P(Data|w) P(w)
Wmap = argmax_w P(Data|w) P(w)
= argmax_w P(Data|w) P(w) → {N(μ, Σ) > Ridge
Laplace ← Lasso

Bayesian Linear Regression

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

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

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

$y = f(x) + \epsilon$ 其中 ϵ 是 noise
 x_i, y_i, ϵ 都是 r.v.
 $x \in \mathbb{R}^p, y \in \mathbb{R}, \epsilon \sim N(0, \sigma^2)$

Bayes Method 因为贝叶斯方法假设 w 是一个分布而不是一个常数

Inference: posterior (w)
Prediction: $x^* \rightarrow y^*$

Inference: $P(w | \text{Data}) = P(w | x, Y) = \frac{P(w, Y | x)}{P(Y | x)} = \frac{P(Y | w, x) P(w)}{\int P(Y | w, x) P(w) dw}$

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

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

其中: Gaussian 分布是高维的

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

假设 $N(\mu, \Sigma)$
- 可以通过计算得到
对 μ, Σ 求导

对 μ, Σ 求导

likelihood: $P(Y | x, w) = \frac{1}{(2\pi)^{\frac{N}{2}}} \frac{1}{\sigma^N} \exp\{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - w^T x_i)^2\}$

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

$\sum_{i=1}^N (y_i - w^T x_i)^2 = (y_1 - w^T x_1, y_2 - w^T x_2, \dots, y_N - w^T x_N)$
 $= (Y - Xw)^T$
 $= (Y^T - w^T X^T)$

$\begin{pmatrix} y_1 - w^T x_1 \\ y_2 - w^T x_2 \\ \vdots \\ y_N - w^T x_N \end{pmatrix}$

$(Y^T - w^T X^T)^T = (Y - Xw)$

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

$P(w | \text{Data}) \propto N(Xw, \sigma^2 I) N(0, \Sigma_p)$

$\propto \exp\{-\frac{1}{2} (Y - Xw)^T \sigma^{-2} I (Y - Xw)\} \cdot \exp\{-\frac{1}{2} w^T \Sigma_p^{-1} w\}$

与 w 相关的项全都在 exp 里面的

$$P(w|Data) \propto N(xw, \sigma^2 I) N(y, \Sigma_p)$$

$$\begin{aligned} &\downarrow \\ N(\mu_w, \Sigma_w) &\propto \exp \left\{ -\frac{1}{2} (y - xw)^T \sigma^2 I (y - xw) \right\} \cdot \exp \left\{ -\frac{1}{2} w^T \Sigma_p^{-1} w \right\} \\ \mu_w = ? &= \exp \left\{ -\frac{1}{2\sigma^2} (y^T - w^T x^T) (y - xw) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\} \\ \Sigma_w = ? &= \exp \left\{ -\frac{1}{2\sigma^2} (y^T y - y^T xw - \underbrace{w^T x^T y}_{\text{二次项}} + w^T x^T xw) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} (y^T y - 2y^T xw + w^T x^T xw) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\} \end{aligned}$$

用更方便的表示。

设 $P(x) = N(\mu, \Sigma)$, 其参数部分为:

$$\begin{aligned} &\exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right) \\ &= -\frac{1}{2} (x^T \Sigma^{-1} - \mu^T \Sigma^{-1}) (x - \mu) \\ &= -\frac{1}{2} (\underbrace{x^T \Sigma^{-1} x}_{\text{二次项}} - \underbrace{2\mu^T \Sigma^{-1} x}_{\text{一次项}} + \underbrace{\mu^T \Sigma^{-1} \mu}_{\text{常数}}) \end{aligned}$$

再回到原来形式:

$$\begin{aligned} &\exp \left\{ -\frac{1}{2\sigma^2} (y^T y - 2y^T xw + w^T x^T xw) - \frac{1}{2} w^T \Sigma_p^{-1} w \right\} \\ \text{二次项} & \left(\text{因 } w \text{ 与 } x \text{ 有关, 故 } w \text{ 与 } x \text{ 有关} \right): -\frac{1}{2\sigma^2} w^T x^T x w - \frac{1}{2} w^T \Sigma_p^{-1} w \\ &= -\frac{1}{2} (w^T (\underbrace{\sigma^2 x^T x + \Sigma_p^{-1}}_{\Sigma_w^{-1}}) w) \end{aligned}$$

$$\text{一次项: } -\frac{1}{2\sigma^2} \cdot (-2) y^T x w = \sigma^2 y^T x w$$

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

$$A \mu_w = \sigma^2 x^T y$$

$$\mu_w = \sigma^2 A^{-1} x^T y$$

$$\therefore P(w|Data) = N(\mu_w, \Sigma_w) \quad \mu_w = \sigma^2 A^{-1} x^T y$$

$$\Sigma_w = A^{-1} \quad (A = \sigma^2 x^T x + \Sigma_p^{-1})$$

Prediction

Inference: $P(w|Data) \quad w|Data \sim N(\mu_w, \Sigma_w)$

$$\begin{aligned} \mu_w &= \sigma^2 A^{-1} x^T y \\ \Sigma_w &= A^{-1} \\ A &= \sigma^2 x^T x + \Sigma_p^{-1} \end{aligned}$$

Prediction: Given $x^* \in \mathcal{X}^*$ Model $f(x) = w^T x = x^T w$
 $\{y = f(x) + \epsilon \mid \epsilon \sim N(0, \sigma^2)\}$

$$\textcircled{1} f(x^*) = x^{*T} w \xrightarrow{\text{noise-free}} P(w|Data) = N(\mu_w, \Sigma_w)$$

$$w \sim N(\mu_w, \Sigma_w) \xrightarrow{\text{噪声}} \mu_w, \Sigma_w$$

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

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

$$\textcircled{2} y^* \quad y^* = f(x^*) + \epsilon \quad \text{其中 } \epsilon \sim N(0, \sigma^2)$$

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

Summary

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

Model $\begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \epsilon, \quad \epsilon \sim N(0, \sigma^2) \end{cases}$

Bayesian Method:

参数 w 不是标量的参数, w 是一个概率分布.

① Inference: $P(w | \text{Data}) \rightarrow \text{posterior}$

$$P(w | \text{Data}) \propto \underbrace{\text{likelihood}}_{N(\Delta, \sigma)} \times \underbrace{\text{prior}}_{N(\Delta, b)}$$

$\mu_w = ?$
 $\Sigma_w = ?$

② prediction: given x^* , y^* ?

$$p(y^* | \text{Data}, x^*) = \int \underbrace{p(y^* | w, \text{Data}, x^*)}_{w \in \text{Data 集, 给定 } x^* \rightarrow p(y^* | w, x^*)} \underbrace{P(w | \text{Data}, x^*)}_{P(w | \text{Data}) \text{ 与 } x^* \text{ 无关, posterior}} dw$$