

1. 极大似然估计

首先要明确极大似然估计要解决的问题。假设某个随机变量的概率分布形式已知，但是要确定这个分布的具体形式还需要知道分布模型的参数。极大似然估计就是一种从观测样本来估计出分布参数的技术；其核心思想就是，分布的参数在取什么具体值的时候会使得观测以最大可能的概率发生。

例：设样本 (X_1, X_2, \dots, X_n) 来自正态总体 $X \sim N(\mu, \sigma^2)$ ， μ, σ^2 未知，求 μ, σ^2 的极大似然估计。

解：

设 (x_1, x_2, \dots, x_n) 为对应的样本观察值，则关于 μ, σ^2 的似然函数为

$$L(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}$$

因此，

$$\ln L(\mu, \sigma^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{令 } \frac{\partial \ln L}{\partial \mu} = 0, \frac{\partial \ln L}{\partial \sigma^2} = 0 \text{ 得到}$$

$$\begin{cases} \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \\ -\frac{n}{2} \frac{1}{\sigma^2} - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \cdot \left(-\frac{1}{\sigma^4}\right) = 0 \end{cases}, \text{ 得到 } \begin{cases} \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \end{cases}$$

例：设样本 $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n), \mathbf{x}_i \in \mathbb{R}^{d \times 1}$ 来自正态总体 $X \sim N(\boldsymbol{\mu}, \Sigma)$ ， $\boldsymbol{\mu}, \Sigma$ 未知，求 $\boldsymbol{\mu}, \Sigma$ 的极大似然估计。

解： d 维随机变量的高斯分布概率密度函数为，

$$f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \text{ 其中 } \boldsymbol{\mu} \in \mathbb{R}^{d \times 1}, \Sigma \in \mathbb{R}^{d \times d} \text{ 为协方差}$$

矩阵，是一个半正定（当然也是实对称）矩阵。

则 $\boldsymbol{\mu}, \Sigma$ 的似然函数为，

$$\begin{aligned} L(\boldsymbol{\mu}, \Sigma) &= \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_i - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-\frac{nd}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_i - \boldsymbol{\mu})\right) \end{aligned}$$

其对数似然函数为，

$$l(\boldsymbol{\mu}, \Sigma) = -\frac{nd}{2} \ln(2\pi) - \frac{n}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \quad (3-1)$$

现在要计算最优的 $\boldsymbol{\mu}^*, \Sigma^*$ 来使得 $l(\boldsymbol{\mu}, \Sigma)$ 达到最大，因此要计算 $l(\boldsymbol{\mu}, \Sigma)$ 的驻点。

$$\frac{dl}{d\boldsymbol{\mu}} = -\frac{1}{2} \frac{d \left\{ \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\}}{d\boldsymbol{\mu}} = -\frac{1}{2} \sum_{i=1}^n \frac{d \{ (\boldsymbol{\mu} - \mathbf{x}_i)^T \Sigma^{-1} (\boldsymbol{\mu} - \mathbf{x}_i) \}}{d\boldsymbol{\mu}}$$

根据第二章（矩阵论）结论 7.5 第 5 条，

$$\frac{dl}{d\boldsymbol{\mu}} = -\frac{1}{2} \sum_{i=1}^n (\Sigma^{-1} + \Sigma^{-T}) (\boldsymbol{\mu} - \mathbf{x}_i) = -\frac{1}{2} \sum_{i=1}^n (\Sigma^{-1} + \Sigma^{-1}) (\boldsymbol{\mu} - \mathbf{x}_i) = \Sigma^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})$$

令上式等于 $\mathbf{0}$ ，则我们有

$$\Sigma^{-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{0}, \text{ 由于 } \Sigma^{-1} \text{ 可逆, 则 } \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{0}, \text{ 因此,}$$

$$\boldsymbol{\mu}^* = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

接下来要计算最优的 Σ^* 。如果要直接计算式（3-1）对 Σ 的导数不太容易，我们把（3-

1）代换成 Σ^{-1} 的等价函数，求出可使得 $l(\boldsymbol{\mu}, \Sigma^{-1})$ 取得最大值时的 $(\Sigma^{-1})^*$ ，当然就可以

求得 Σ^* 。

令 $A = \Sigma^{-1}$ ，则 $A^{-1} = \Sigma$ ，且容易知道， $A = A^T$ ，式（3-1）可等价转换为，

$$l(\boldsymbol{\mu}, A) = -\frac{nd}{2} \ln(2\pi) - \frac{n}{2} \ln(|A^{-1}|) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T A (\mathbf{x}_i - \boldsymbol{\mu})$$

接下去计算 $l(\boldsymbol{\mu}, A)$ 关于 A 这部分的梯度为零的点。

$$\begin{aligned} \frac{d \{ l(\boldsymbol{\mu}, A) \}}{dA} &= -\frac{n}{2} \frac{d \{ \ln(|A^{-1}|) \}}{dA} - \frac{1}{2} \frac{d \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T A (\mathbf{x}_i - \boldsymbol{\mu})}{dA} \\ &= -\frac{n}{2} \frac{d \{ -\ln|A| \}}{dA} - \frac{1}{2} \sum_{i=1}^n \frac{d \{ (\mathbf{x}_i - \boldsymbol{\mu})^T A (\mathbf{x}_i - \boldsymbol{\mu}) \}}{dA} \\ &= \frac{n}{2} \frac{1}{|A|} |A| A^{-T} - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \\ &= \frac{n}{2} A^{-1} - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \end{aligned}$$

$$\text{令 } \frac{d\{l(\boldsymbol{\mu}, A)\}}{dA} = \mathbf{0}, \text{ 则有}$$

$$\frac{n}{2}A^{-1} - \frac{1}{2}\sum_{i=1}^n(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \mathbf{0}, \text{ 则 } A^{-1} = \frac{1}{n}\sum_{i=1}^n(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T, \text{ 则}$$

$$\Sigma^* = \frac{1}{n}\sum_{i=1}^n(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T。$$