

Maximum Likelihood Estimate (极大似然估计)

Fisher

Statistical Model: $f(x|\theta)$, $x \in \mathbb{R}^d$ etc. $\theta \in \mathbb{R}^p$
probability model

Data: i.i.d. $x_1, \dots, x_n \sim f(x|\theta)$ $\theta_0 \in \mathbb{R}^p$

目标: $\hat{\theta} = G(x_1, \dots, x_n) \xrightarrow{n \rightarrow \infty} \theta_0$?

$$\begin{aligned} \text{MLE } \hat{\theta}^{\text{MLE}} &= \arg \max_{\theta \in \Theta} \prod_{i=1}^n f(x_i|\theta) \quad (\text{M-estimate}) \\ &= \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ln f(x_i|\theta) \end{aligned}$$

例子: $f(x|\theta) = \frac{1}{\sqrt{2\pi}|\Sigma|} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right)$, $\theta = (\mu, \Sigma)$
 $x \in \mathbb{R}^p$

x_1, \dots, x_n i.i.d.

MLE \rightarrow ?

Log Likelihood:

$$\log f(x|\theta) = -\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu) - \frac{1}{2} \ln |\Sigma| + \text{Const}$$

$$I_n = \frac{1}{n} \sum_{i=1}^n \log f(x_i|\theta) = -\frac{1}{2n} \sum_{i=1}^n (x_i-\mu)^\top \Sigma^{-1}(x_i-\mu) - \frac{1}{2n} \ln |\Sigma| + C$$

1st order condition

$$0 = \frac{\partial I_n}{\partial \mu} = -\frac{1}{n} \sum_{i=1}^n \Sigma^{-1}(x_i-\mu) \Rightarrow \hat{\mu}^{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i$$

sample mean !

$$\text{trace} I_n(\Sigma) = -\frac{1}{2n} \sum_{i=1}^n \text{tr} \left[(X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right] - \frac{1}{2} \log |\Sigma| + C$$

linear, cyclic property $\text{tr}(AB) = \text{tr}(BA)$
 $\text{tr}(ABC) = \text{tr}(BCA) = \dots$

$$\frac{1}{2n} \sum_{i=1}^n \text{trace} \left[(X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right] = \frac{1}{2n} \sum_{i=1}^n \text{tr} \left[\Sigma^{-1} (X_i - \mu) (X_i - \mu)^T \right]$$

$$= \frac{1}{2} \text{tr} \left[\Sigma^{-1} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu) (X_i - \mu)^T \right) \right]$$

$$\hat{S}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}) (X_i - \hat{\mu})^T, \quad S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu) (X_i - \mu)^T$$

$$= \frac{1}{2} \text{tr} (\Sigma^{-1} S_n) \quad S_n^{\frac{1}{2}} S_n^{\frac{1}{2}}$$

$$= \frac{1}{2} \text{tr} (S_n^{\frac{1}{2}} \Sigma^{-1} S_n^{\frac{1}{2}})$$

对称. Symmetric, p.s.d.

$$S = S_n^{\frac{1}{2}} \Sigma^{-1} S_n^{\frac{1}{2}} = U \Lambda U^T \quad \Lambda = \text{diag}(\lambda_i)_{i=1, p} \quad \lambda_i \geq 0$$

$$\Sigma = S_n^{-\frac{1}{2}} S^{-1} S_n^{\frac{1}{2}} \quad \det(AB) = |AB| = |A| \cdot |B|$$

$$\log |\Sigma| = + \log |S_n| - \log |S| \quad S(\Sigma) \text{ 变}$$

S_n sample cov. Σ 未知

$$\arg \max I_n(\Sigma) = -\frac{1}{2} \text{trace}(S) + \frac{1}{2} \log |S| + C(S_n, \mu)$$

$$= -\frac{1}{2} \sum_{i=1}^p \lambda_i + \frac{1}{2} \sum_{i=1}^p \log \lambda_i + C$$

$$\frac{\partial I_n}{\partial \lambda_i} = -\frac{1}{2} + \frac{1}{2\lambda_i} \Rightarrow \lambda_i = 1$$

$$S = I_p = S_n^{-\frac{1}{2}} \Sigma^{-1} S_n^{\frac{1}{2}} \Rightarrow \Sigma^{\text{MLE}} = S_n^{-1}$$

$$= \frac{1}{n} \sum_{i=1}^n (X_i - \mu) (X_i - \mu)^T$$

Note $\hat{\Sigma}_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})(X_i - \hat{\mu})^T$

总结: $\hat{\mu}^{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$ sample mean

$\hat{\Sigma}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}^{MLE})(X_i - \hat{\mu}^{MLE})^T$ sample covariance

为什么 MLE.

Generally: $X_i \sim f(X_i | \theta_0)$ ← unknown.

$$\hat{\theta}^{MLE} = \arg \max_{\theta} L(X_{1:n} | \theta) = \prod_{i=1}^n f(X_i | \theta)$$

$\theta \in \mathbb{R}^p$ p fixed. $n \rightarrow \infty$ limitly properties

1) Consistency $\hat{\theta}^{MLE} \xrightarrow{n \rightarrow \infty} \theta_0$ (prob./almost sure)

2) Asymptotic Normality $\sqrt{n}(\hat{\theta}^{MLE} - \theta_0) \xrightarrow{d} \mathcal{N}(\theta_0, I^{-1})$

I : Fisher Information matrix

$$I_{ij} = - \mathbb{E}_X \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f(X | \theta_0) \right]_{p \times p} \geq 0$$

⇒ Asymptotic Efficiency: (second order)

$$\text{tr}(I^{-1}) = \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}^{MLE}) \leq \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}^1); \hat{\theta}^1 \text{ unbiased}$$

参数空间有限维 $\lim_{n \rightarrow \infty} \frac{p}{n} = 0$ 传统统计学

Big Data $n \rightarrow \infty$ $p_n \rightarrow \infty$ 高维统计学

$$\lim_{n \rightarrow \infty} \frac{p_n}{n} \rightarrow C \neq 0$$

MLE 还好估计? Stein's phenomenon.

有限 n , $p \geq 3$. \exists JS, MLE 好! www.ebanshu.com

Stein's Phenomenon

MLE $X_1, \dots, X_n \stackrel{iid.}{\sim} N(\mu, \Sigma)$

$$\hat{\mu}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu \quad \text{Consistency.}$$

$$E[\hat{I}^{-1}] = \lim_{n \rightarrow \infty} \text{var}(\hat{\mu}_n^{MLE}) \leq \lim_{n \rightarrow \infty} \text{var}(\hat{\mu}_n) \quad \left. \begin{array}{l} n \rightarrow \infty \\ \text{fixed } p \end{array} \right\}$$

Without loss of Generality, $\Sigma = U \Lambda U^T$ $\Lambda = \text{diag}(\lambda_i)$
 $Y_i = \Lambda^{-\frac{1}{2}} U^T X_i$ P.C.A.

$$Y_i \sim N(\mu, I_p)$$

Risk (Mean Square Error) MSE

Given $\hat{\mu}_n(Y_1, \dots, Y_n)$

$$\text{Risk } R(\hat{\mu}_n, \mu) = E_{Y_1, \dots, Y_n} L(\hat{\mu}_n(Y_1, \dots, Y_n), \mu)$$

$$\stackrel{\text{MSE}}{=} E \|\hat{\mu}_n - \mu\|^2 \quad \hat{\mu}_n, \mu \in \mathbb{R}^p$$

Bias-Variance

$$\begin{aligned} R(\hat{\mu}_n, \mu) &= E \|\hat{\mu}_n - E(\hat{\mu}_n) + E(\hat{\mu}_n) - \mu\|^2 \\ &= E \|\hat{\mu}_n - E(\hat{\mu}_n)\|^2 + \|E(\hat{\mu}_n) - \mu\|^2 + \cancel{2 E(\hat{\mu}_n - E(\hat{\mu}_n)) \cdot (E(\hat{\mu}_n) - \mu)} \\ &= \text{Var}(\hat{\mu}_n) + \text{Bias}(\hat{\mu}_n) \end{aligned}$$

Example $Y_i \sim N(\mu, \sigma^2 I_p)$ $\hat{\mu}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$ $E[Y_i] = \mu$

$$\text{Bias}(\hat{\mu}_n^{MLE}) = 0 \quad \text{unbiased!}$$

$$\text{Var}(\hat{\mu}_n^{MLE}) = \frac{\sigma^2}{n}$$

In particular, $n=1$, $\text{Var}(\hat{\mu}_1^{MLE}) = \sigma^2 p$

$$R(\hat{\mu}_1^{MLE}, \mu) = \sigma^2 p$$

Linear Estimator

$$\hat{\mu}_C(Y) = CY, \quad Y \sim N(\mu, \sigma^2 I_p) \quad n=1$$

$$C = I \rightarrow \text{MLE}$$

$$C = \text{diag}(c_i) \quad \min_{\theta} \frac{1}{2} \|Y - \theta\|^2 + \frac{\lambda}{2} \|\theta\|^2 \quad \text{Ridge Regression}$$

$$c_i = \frac{1}{1+\lambda} \quad \hat{\theta} = \frac{1}{1+\lambda} Y$$

$$\text{Bias}(\hat{\mu}_C) = \mathbb{E}[\hat{\mu}_C] - \mu = \|(I - C)\mu\|^2 \quad \mathbb{E}(CY) = C\mu$$

$$\begin{aligned} \text{Var}(\hat{\mu}_C) &= \text{tr} \mathbb{E}(CY - C\mu)^T (CY - C\mu) \\ &= \mathbb{E} \text{tr}[(Y - \mu)^T C^T C (Y - \mu)] = \text{tr}[C^T C] \mathbb{E} \underbrace{(Y - \mu)(Y - \mu)^T}_{\sigma^2 I_p} \\ &= \sigma^2 \text{tr}(C^T C) \end{aligned}$$

$$C = \text{diag}(c_i)$$

$$R(\hat{\mu}_C, \mu) = \sum_{i=1}^p \sigma^2 c_i^2 + \sum_{i=1}^p (1 - c_i)^2 \mu_i^2 \leq \sum_{i=1}^p \sigma^2 c_i^2 + \sum_{i=1}^p (1 - c_i)^2 \tau_i^2$$

Statistical Decision theory: Minimax Risk

$|\mu_i| < \tau_i$ Rectangular class

$$\inf_{c_i} \sup_{|\mu_i| < \tau_i} R(\hat{\mu}_C, \mu) = \sum_{i=1}^p \frac{\sigma^2 \tau_i^2}{\tau_i^2 + \sigma^2} \leq \sum_{i=1}^p \sigma^2 \tau_i^2 \quad \text{MLE}$$

sparse family τ_i^2

Problem:

$\hat{\mu}_n$ better estimator?

Inadmissible: $\hat{\mu}_n$ is inadmissible

$$\exists \mu_n^* \text{ s.t. } \mathbb{E} \|\mu_n^* - \mu\|^2 \leq \mathbb{E} \|\hat{\mu}_n - \mu\|^2 \text{ for all } \mu \in \mathbb{R}^p$$

$$\exists \mu_0 \quad R(\mu_n^*, \mu_0) < R(\hat{\mu}_n, \mu_0)$$

$\hat{\mu}_n^{\text{MLE}}$

inadmissible \rightarrow
Yes

Stein '1966, James-Stein '1961

$$\hat{\mu}_n^{JS} = \left(1 - \frac{\sigma^2(p-2)}{\|\hat{\mu}_n^{MLE}\|^2} \right) \hat{\mu}_n^{MLE}, \quad Y \sim N(\mu, \sigma^2 I_p)$$

Thm

$$R(\hat{\mu}_n^{JS}, \mu) < R(\hat{\mu}_n^{MLE}, \mu), \quad \exists \mu \in \mathbb{R}^p, p \geq 3$$

几乎所有 μ

where $\hat{\mu}_n^{MLE} = Y, \quad \hat{\mu}_n^{JS} = \left(1 - \frac{\sigma^2(p-2)}{\|Y\|^2} \right) Y$

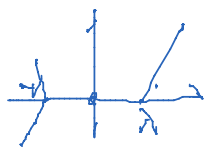
Stein's Unbiased Risk Estimates (SURE) : $Y \sim N(\mu, I_p)$

$$\hat{\mu}(Y) = Y + g(Y), \quad g \text{ nonlinear}$$

1) Linear Est. $\hat{\mu} = CY, \quad g(Y) = (C-I)Y$

2) Soft Thresholding $g_{ST}(Y) = \begin{cases} \lambda & Y_i > \lambda \\ -Y_i & |Y_i| \leq \lambda \\ \lambda & Y_i < -\lambda \end{cases}$

$$\min_{\hat{\mu}} \frac{1}{2} \|Y - \hat{\mu}\|^2 + \lambda \|\hat{\mu}\|_1 \Rightarrow \hat{\mu}_{ST} = Y - \hat{\mu}_{ST} + \lambda \partial \|\hat{\mu}\|_1 = 0$$



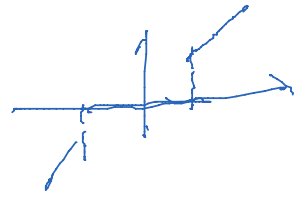
$$\hat{\mu}_i \neq 0: Y_i - \hat{\mu}_i + \lambda \partial g(\hat{\mu}_i) = 0$$

$$\hat{\mu}_j = 0: |Y_j - \hat{\mu}_j| < \lambda$$

3) Hard Thresholding $\frac{1}{2} \|Y - \hat{\mu}\|^2 + \lambda \|\hat{\mu}\|_0$

$$g_{HT}(Y) = \begin{cases} 0 & |Y_i| > \lambda \\ -Y_i & |Y_i| \leq \lambda \end{cases}$$

- not weakly differentiable



4) JS $g_{JS}(Y) = - \frac{\sigma^2(p-2)}{\|Y\|^2} Y, \quad \sigma^2 = 1$

1) weakly differentiable. $g(x_i, x_{-i}) \sim \forall i: x_{-i}$
absolutely continuous w.r.t. x_i

$$2) \sum_{i=1}^p \int |\partial_i g_i(x)| dx < \infty$$

$\hat{\mu}(Y + g(Y))$ suitable g above

Lemma (Stein '61)

$$R(\hat{\mu}, \mu) = \mathbb{E} [\rho + 2 \nabla^T g(Y) + \|g(Y)\|^2]$$

$$\nabla^T g(Y) := \sum_{i=1}^p \frac{\partial}{\partial Y_i} g_i(Y)$$

Proof (Integration by Parts)

$$\begin{aligned} \mathbb{E} \|\hat{\mu} - \mu\|^2 &= \mathbb{E} \|Y + g(Y) - \mu\|^2 = \mathbb{E} \|(Y - \mu) + g(Y)\|^2 \\ &= \mathbb{E} \|Y - \mu\|^2 + 2 \mathbb{E} (Y - \mu)^T g(Y) + \mathbb{E} \|g(Y)\|^2 \end{aligned}$$

$Y \sim \mathcal{N}(\mu, \sigma^2 I_p)$

$$\mathbb{E} [(Y - \mu)^T g(Y)] = \sum_{i=1}^p \int_{-\infty}^{+\infty} \frac{\partial g_i(Y)}{\partial Y_i} \phi(Y - \mu) dY = \mathbb{E} [\nabla^T g(Y)]$$

$y \sim \mathcal{N}(\mu, 1)$ $\phi(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2}}$ $\int_{-\infty}^{+\infty} (y-\mu) g(y) \phi(y-\mu) dy$

$\frac{\partial}{\partial y} \phi = -(y-\mu) \phi(y)$ $= - \int_{-\infty}^{+\infty} g(y) \frac{\partial}{\partial y} \phi(y-\mu) dy$

$= \underbrace{g(y) \phi(y-\mu)}_{\downarrow 0} \Big|_{-\infty}^{+\infty} + \int_{-\infty}^{+\infty} \phi(y-\mu) \frac{dg(y)}{dy} dy$

\therefore above = $\mathbb{E} [\rho + \nabla^T g(Y) + \|g(Y)\|^2]$

$U(Y) = \rho + \nabla^T g(Y) + \|g(Y)\|^2$

$$U(Y) = p + \frac{D^T(Y)}{\|Y\|^2} + \|Y\|^2$$

JS.

$$g(Y) = -\frac{p-2}{\|Y\|^2} Y$$

$$U(Y) = p + 2 \left[\sum_{i=1}^p \frac{\partial}{\partial Y_i} \left(\frac{p-2}{\|Y\|^2} Y_i \right) \right] + \frac{(p-2)^2}{\|Y\|^2}$$

$$\Rightarrow p - \frac{(p-2)^2}{\|Y\|^2} + \frac{(p-2)^2}{\|Y\|^2}$$

$$\underline{R(\hat{\mu}^{JS}, \mu) = \mathbb{E} U(Y) = p - \mathbb{E} \frac{(p-2)^2}{\|Y\|^2} < p = R(\hat{\mu}^{MLE}, \mu)}$$

$n=1, \sigma=1$

$\hat{\mu}^{MLE}$ inadmissible

$p \geq 3$

$\frac{2p-1}{n}$

Note

$$Y \sim N(\mu, I_p) \quad \forall \hat{\mu} = CY \quad (\text{Lin. est})$$

$\hat{\mu}$ is admissible iff

$$p \text{ C sym. } C=C^T$$

$$\geq^0 \quad 0 \leq \text{eigval}(C) \leq 1$$

$$\geq^1 \quad \text{eigval}_i(C) = 1 \text{ for at most two } i$$

Lemma 2.8. Johnstone (GE)

$$\text{例: } \hat{\mu}^{JS+} = \left(1 - \frac{p-2}{\|Y\|^2} \right) Y \quad \text{better than MLE JS.}$$

ST.

$$R(\hat{\mu}^{ST}, \mu) = 1 + (2 \log p + 1) \sum_{i=1}^p (\mu_i^2 / 1)$$

$$\text{注: } R(\hat{\mu}^{JS}, \mu) = 2 + c \left(\sum_{i=1}^p \mu_i^2 \wedge p \right) \quad c \in \left(\frac{1}{2}, 1 \right)$$

ST < JS.

sparse. $\mu = (*, 0, \dots, 0)$

dense. $\mu = (1, \dots, 1)$

$2 \log p + 1 < 0(c, p)$
www.ebanhu.com
 $R^{ST} < R^{JS}$

MLE. 有限 n . 有限 p . (≥ 3)

是 inadmissible

$$\text{MSE (JS)} < \text{MSE (MLE)}$$

$$\text{MSE (ST)} < \dots$$

(Lin)

"Shrinkage." better than MLE,