

Stein's Phenomenon

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

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

$$\text{tr}(\hat{I}) = \lim_{n \rightarrow \infty} \text{var}(\hat{\mu}_n^{MLE}) \leq \lim_{n \rightarrow \infty} \text{var}(\hat{\mu}_n)$$

$n \rightarrow \infty$
fixed p

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

$$\text{Risk. } R(\hat{\mu}_n, \mu) = \mathbb{E}_{Y_1, \dots, Y_n} L(\hat{\mu}_n(Y_1, \dots, Y_n), \mu)$$

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

Bias-Variance

$$R(\hat{\mu}_n, \mu) = \mathbb{E} \|\hat{\mu}_n - \mathbb{E}(\hat{\mu}_n) + \mathbb{E}(\hat{\mu}_n) - \mu\|^2$$

$$= \mathbb{E} \|\hat{\mu}_n - \mathbb{E}(\hat{\mu}_n)\|^2 + \|\mathbb{E}(\hat{\mu}_n) - \mu\|^2 + \cancel{2 \mathbb{E} \langle \hat{\mu}_n - \mathbb{E}(\hat{\mu}_n), \mathbb{E}(\hat{\mu}_n) - \mu \rangle}$$

$$= \text{Var}(\hat{\mu}_n) + \text{Bias}(\hat{\mu}_n)$$

Examp

$$Y_i \sim N(\mu, \sigma^2 I_p)$$

$$\hat{\mu}_n^{MLE} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\mathbb{E}[Y_i] = \mu$$

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

unbiased!

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

In particular, $n=1$.

$$\text{Var}(\hat{\mu}_1^{MLE}) = \sigma^2 I_p$$

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

Linear Estimator

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

Copyright by 电子书

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

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

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

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

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

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

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

Statistical Decision theory : Minimax Risk

$|\mu_i| < \tau_i$ Rectangular class

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

sparse family $\tau_i \downarrow$

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

Yes

Stein '1956, James-Stein '1961

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

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

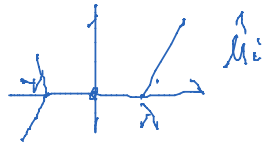
Stein's Unbiased Risk Estimates (SURE) : $Y \sim \mathcal{N}(\mu, I_p).$

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

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

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

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



$$\Leftrightarrow \begin{cases} \hat{\mu}_i \neq 0: & Y_i - \hat{\mu}_i + \lambda \partial g(\hat{\mu}_i) = 0 \\ \hat{\mu}_j = 0 & |Y_j - \hat{\mu}_j| \leq \lambda \end{cases}$$

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

not weakly differentiable

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



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

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

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

available g. above

Lemma (Stein '61)

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

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

Proof (Integration by Parts)

$$\mathbb{E} \|\hat{\mu} - \mu\|^2 = \mathbb{E} \|Y + g(Y) - \mu\|^2 = \mathbb{E} \|(Y - \mu) + g(Y)\|^2$$

$$= \mathbb{E} \|Y - \mu\|^2 + 2 \mathbb{E} (Y - \mu)^T g(Y) + \mathbb{E} \|g(Y)\|^2$$

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

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

$$y \sim N(\mu, 1) \quad \phi(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2}}$$

$$\frac{\partial}{\partial y} \phi = -(y - \mu) \phi(y)$$

$$\begin{aligned} & \int_{-\infty}^{+\infty} (y - \mu) g(y) \phi(y - \mu) dy \\ &= - \int_{-\infty}^{+\infty} g(y) \frac{\partial}{\partial y} \phi(y - \mu) dy \\ &= -g(y) \phi(y - \mu) \Big|_{-\infty}^{+\infty} + \int_{-\infty}^{+\infty} \phi(y - \mu) \frac{dg(y)}{dy} dy \end{aligned}$$

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

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

$$U(Y) = p + \nabla_{\mu}^T(Y) + \|Y\|^2$$

JS.

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

$$\begin{aligned} U(Y) &= p + 2 \left(-\sum_{i=1}^p \frac{\partial}{\partial Y_i} \left(\frac{p-2}{\|Y\|^2} Y_i \right) \right) + \frac{(p-2)^2}{\|Y\|^2} \\ &= p - 2 \frac{(p-2)}{\|Y\|^2} + \frac{(p-2)^2}{\|Y\|^2} \end{aligned}$$

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

$n=1, \sigma=1$

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

$p \geq 3$

$\frac{2p}{n}$

Note

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

$\hat{\mu}$ is admissible iff

$$1) \quad C \text{ sym.} \quad C = C^T$$

$$2) \quad 0 \leq \text{eigval}(C) \leq 1$$

$$3) \quad \text{eigval}_i(C) = 1 \text{ for at most two } i$$

Lemma 2.8. Johnstone (GE)

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

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

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

$ST < JS$.

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

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

$2 \log p + 1 < 0(p)$

$R^{ST} < R^{JS}$ www.ebanshu.com

MLE. 有限 n . 有限 p . (≥ 3)

是 inadmissible

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

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

"Shrinkage" better than MLE,