

1. *Maximum Likelihood Method*: consider  $n$  random samples from a multivariate normal distribution,  $X_i \in \mathbb{R}^p \sim \mathcal{N}(\mu, \Sigma)$  with  $i = 1, \dots, n$ .

(a) Show the log-likelihood function

$$l_n(\mu, \Sigma) = -\frac{n}{2} \text{trace}(\Sigma^{-1} S_n) - \frac{n}{2} \log \det(\Sigma) + C,$$

where  $S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)(X_i - \mu)^T$ , and some constant  $C$  does not depend on  $\mu$  and  $\Sigma$ ;

proof:  $f(x) = \frac{1}{(2\pi)^p |\Sigma|} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\}$

$$\log f(x) = \sum_{i=1}^n \left[ -\frac{p}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right]$$

$$\begin{aligned} \Rightarrow \text{trace}(\log f(x)) &= \text{tr} \left( -\frac{np}{2} \log 2\pi \right) - \frac{n}{2} \log \det(\Sigma) - \frac{1}{2} \sum_{i=1}^n \text{tr} \left[ (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right] \\ &= -\frac{1}{2} \sum_{i=1}^n \text{tr} \left[ \Sigma^{-1} (x_i - \mu)(x_i - \mu)^T \right] \\ &= -\frac{n}{2} \text{tr}(\Sigma^{-1} S_n) - \frac{n}{2} \log [\det(\Sigma)] + C \end{aligned}$$

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

(b) Show that  $f(X) = \text{trace}(AX^{-1})$  with  $A, X \succeq 0$  has a first-order approximation,

$$f(X + \Delta) \approx f(X) - \text{trace}(X^{-1} A' X^{-1} \Delta)$$

hence formally  $df(X)/dX = -X^{-1} A X^{-1}$  (note  $(I + X)^{-1} \approx I - X$ );

proof:  $f(x+\Delta) = \text{tr} [A (X+\Delta)^{-1}] = \text{tr} \left\{ A [(I + \Delta X^{-1}) X]^{-1} \right\}$

$$\begin{aligned} &= \text{tr} \left\{ A X^{-1} (I + \Delta X^{-1})^{-1} \right\} \approx \text{tr} \left\{ A X^{-1} (I - \Delta X^{-1}) \right\} \\ &= \text{tr}(A X^{-1}) - \text{tr}(A X^{-1} \Delta X^{-1}) \\ &= \text{tr}(A X^{-1}) - \text{tr} \left\{ (X^{-1})^T \Delta^T (X^{-1})^T A^T \right\} \\ &= \text{tr}(A X^{-1}) - \text{tr}(X^{-1} \Delta X^{-1} A) \\ &= \text{tr}(A X^{-1}) - \text{tr}(X^{-1} A' X^{-1} \Delta) \\ &= f(x) - \text{tr}(X^{-1} A' X^{-1} \Delta) \end{aligned}$$

(c) Show that  $g(X) = \log \det(X)$  with  $A, X \succeq 0$  has a first-order approximation,

$$g(X + \Delta) \approx g(X) + \text{trace}(X^{-1}\Delta)$$

hence  $dg(X)/dX = X^{-1}$  (note: consider eigenvalues of  $X^{-1/2}\Delta X^{-1/2}$ );

proof:  $g(X+\Delta) = \log \det(X+\Delta)$

$$= \log \det [X^{\frac{1}{2}} (I + X^{-\frac{1}{2}} \Delta X^{-\frac{1}{2}}) X^{\frac{1}{2}}]$$

$$= \log \det X + \log \det (I + X^{-\frac{1}{2}} \Delta X^{-\frac{1}{2}})$$

$$= \log \det X + \sum_{i=1}^n \log (1 + \lambda_i)$$

$$= \log \det X + \text{tr} (X^{-\frac{1}{2}} \Delta X^{-\frac{1}{2}})$$

$$= \log \det X + \text{tr} (X^{-1} \Delta)$$

thus  $g(X+\Delta) = g(X) + \text{tr} (X^{-1} \Delta)$

(d) Use these formal derivatives with respect to positive semi-definite matrix variables to show that the maximum likelihood estimator of  $\Sigma$  is

$$\hat{\Sigma}_n^{MLE} = S_n.$$

proof:  $l(\Sigma) = -\frac{n}{2} \text{tr}(\Sigma^{-1} S_n) - \frac{n}{2} \log[\det(\Sigma)] + C$

$$\frac{\partial l(\Sigma)}{\partial \Sigma} = -\frac{n}{2} (-\Sigma^{-1} S_n \Sigma^{-1}) - \frac{n}{2} \Sigma^{-1} = 0$$

$$\Rightarrow \hat{\Sigma} = S_n.$$

2. *Shrinkage*: Suppose  $y \sim \mathcal{N}(\mu, I_p)$ .

(a) Consider the Ridge regression

$$\min_{\mu} \frac{1}{2} \|y - \mu\|_2^2 + \frac{\lambda}{2} \|\mu\|_2^2.$$

Show that the solution is given by

$$\hat{\mu}_i^{\text{ridge}} = \frac{1}{1 + \lambda} y_i.$$

Compute the risk (mean square error) of this estimator. The risk of MLE is given when  $C = I$ .

proof.  $RSS(\mu) = \frac{1}{2} (Y - \mu)^T (Y - \mu) + \frac{\lambda}{2} \mu^T \mu$

$$= \frac{1}{2} Y^T Y - Y^T \mu + \mu^T \mu + \frac{\lambda}{2} \mu^T \mu$$

$$\text{Let } \frac{\partial RSS(\mu)}{\partial \mu} = -\frac{1}{2} Y + \frac{1+\lambda}{2} \mu = 0$$

$$\Rightarrow \hat{\mu}^{\text{ridge}} = \frac{1}{1+\lambda} Y \Rightarrow \hat{\mu}_i^{\text{ridge}} = \frac{1}{1+\lambda} y_i$$

Def of Risk on  $\hat{\mu}$   $R(\mu, \hat{\mu}) = L(\mu, \hat{\mu})$

when  $L(\mu, \hat{\mu}) = \|\mu - \hat{\mu}\|^2$

$$\text{MSE Risk}(\mu, \hat{\mu}) = \text{Var}(\hat{\mu}) + \text{Bias}^2(\hat{\mu})$$

since  $\hat{\mu}^{\text{ridge}} = \hat{\mu}_c = CY$  with  $C = X(X^T X + \lambda I)^{-1} X^T$  where  $X = I$

so  $\text{Bias}(\hat{\mu}_c) = \|(CI - C)\mu\|^2 = 0$

$$\text{Var}(\hat{\mu}_c) = \sigma^2 \text{tr}(C^T C) = 1 \cdot p = p$$

$$\text{Risk}(\hat{\mu}^{\text{ridge}}) = p$$

(b) Consider the LASSO problem,

$$\min_{\mu} \frac{1}{2} \|y - \mu\|_2^2 + \lambda \|\mu\|_1.$$

Show that the solution is given by Soft-Thresholding

$$\hat{\mu}_i^{\text{soft}} = \mu_{\text{soft}}(y_i; \lambda) := \text{sign}(y_i)(|y_i| - \lambda)_+.$$

For the choice  $\lambda = \sqrt{2 \log p}$ , show that the risk is bounded by

$$\mathbb{E} \|\hat{\mu}^{\text{soft}}(y) - \mu\|^2 \leq 1 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, 1).$$

Under what conditions on  $\mu$ , such a risk is smaller than that of MLE? Note: see Gaussian Estimation by Iain Johnstone, Lemma 2.9 and the reasoning before it.

Proof:  $\min_{\mu} J(\mu) = \frac{1}{2} \|Y - \mu\|^2 + \lambda \|\mu\|_1$

$$2\|x\|_1 = \text{sign}(x) = \begin{cases} 1 & x > 0 \\ [-1, 1] & x = 0 \\ -1 & x < 0 \end{cases}$$

$$\Rightarrow 2J(u)|_{u^*} = u^* \gamma + \lambda \operatorname{sign}(u^*)$$

$$\text{if } u^* > 0 \quad \begin{aligned} u^* - \gamma + \lambda &= 0 \\ u^* &= \gamma - \lambda > 0 \end{aligned} \quad \gamma > \lambda$$

$$\text{if } u^* < 0 \quad \begin{aligned} u^* - \gamma - \lambda &= 0 \\ u^* &= \gamma + \lambda < 0 \end{aligned} \quad \gamma + \lambda < 0$$

$$\text{if } u^* = 0 \quad \begin{aligned} u^* &\in \gamma + \lambda \\ \gamma + \lambda &= 0 \end{aligned} \quad \lambda \in [-\lambda, \lambda]$$

$$\Rightarrow u^* = \begin{cases} \gamma - \lambda & \gamma > \lambda \\ 0 & -\lambda \leq \gamma \leq \lambda \\ \gamma + \lambda & \gamma < -\lambda \end{cases}$$

$$u^{\text{soft}} = \operatorname{sign}(\gamma) (|\gamma| - \lambda)^+$$

then we found the upper bound of the risk.

Suppose  $\hat{u}_\lambda(\gamma) = \gamma + g_\lambda(\gamma)$   $g_\lambda(\gamma)$  is for soft thresholding

Consider  $Y = u + Z \sim N(u, 1)$

the risk function  $r_\lambda(u, u) = \int [\hat{u}_\lambda(u+Z) - u]^2 \phi(Z) dZ$

$$r_\lambda(u, 0) \leq 2\lambda^{-1} \phi(\lambda) \leq e^{-\lambda^2/2} \quad r_\lambda(u, \infty) = 1 + \lambda^2$$

$$\frac{2r_\lambda(u, u)}{2u} = 2u P(|u+Z| \leq \lambda) \leq 2u$$

$$\therefore r_\lambda(u, u) - r_\lambda(u, 0) \leq u^2$$

$$r_\lambda(u, u) \leq r_\lambda(u, 0) + \min(u^2, 1 + \lambda^2) \quad \text{Let } \lambda = \sqrt{2 \log P}$$

$$r_\lambda(u, u) \leq r_\lambda(u, 0) + (2 \log P + 1) \min(u^2, 1)$$

when  $Y \sim N(u, \epsilon^2 I)$

$$\text{Risk}(\hat{u}, u) \leq \epsilon^2 + (2 \log P + 1) \sum_{i=1}^P \min(u_i^2, \epsilon^2)$$

for  $\epsilon = 1$

$$\|\hat{\mu}^{\text{soft}}(y) - \mu\|^2 \leq (1 + (2 \log P + 1) \sum_{i=1}^P \min(\mu_i^2, 1))$$

the condition for  $R(\hat{\mu}^{\text{soft}}) < R(\hat{\mu}^{\text{MLE}})$

$$\text{is } 1 + (2 \log P + 1) \sum_{i=1}^P \min(\mu_i^2, 1) < P$$

(c) Consider the  $l_0$  regularization

$$\min_{\mu} \|y - \mu\|_2^2 + \lambda^2 \|\mu\|_0,$$

where  $\|\mu\|_0 := \sum_{i=1}^P I(\mu_i \neq 0)$ . Show that the solution is given by Hard-Thresholding

$$\hat{\mu}_i^{\text{hard}} = \mu_{\text{hard}}(y_i; \lambda) := y_i I(|y_i| > \lambda).$$

Proof:  $R_{\text{ss}}(\mu) = (y - \mu)^T (y - \mu) + \lambda^2 \left( \sum_{i=1}^P \mathbb{I}(\mu_i \neq 0) \right)$

$$= y^T y - y^T \mu - \mu^T y + \mu^T \mu + \lambda^2 \left( \sum_{i=1}^P \mathbb{I}(\mu_i \neq 0) \right)$$

$$= \begin{cases} (y_i - \mu_i)^2 + \lambda^2 & \mu_i \neq 0 \\ y_i^2 & \mu_i = 0 \end{cases}$$

$$y_i^2 < (y_i - \mu_i)^2 + \lambda^2 \quad \text{for all } \mu_i$$

$$\mu_i^2 - 2y_i \mu_i + \lambda^2 > 0 \quad \text{for all } \mu_i$$

$$\Leftrightarrow \Delta = 4y_i^2 - 4\lambda^2 < 0 \quad \Leftrightarrow y_i^2 < \lambda^2$$

so when  $y_i^2 < \lambda^2$  we adopt  $\mu_i = 0$ , otherwise  $\mu_i = y_i$

$$\hat{\mu}^{\text{hard}}(y_i) = (1 - g(y_i)) y_i = (1 - \mathbb{I}(|y_i| \leq \lambda)) y_i = y_i \mathbb{I}(|y_i| > \lambda)$$

$g(y) = \mathbb{I}(|y| \leq \lambda)$  is not weakly differentiable because  $g$  is not absolutely continuous.