

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  $\Sigma$ ;

**Solution.** (a) Since  $X_i \in \mathbb{R}^p \stackrel{\text{iid}}{\sim} N(\mu, \Sigma)$ , ( $i = 1, \dots, n$ )

$$\text{the density function is } f(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\},$$

$$\text{the likelihood function is } P(X_1, \dots, X_n | \mu, \Sigma) = \prod_{i=1}^n \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) \right\},$$

$$\begin{aligned} \text{so the log-likelihood is } \ln(\mu, \Sigma) &= \log P(X_1, \dots, X_n | \mu, \Sigma) = \sum_{i=1}^n \left( \log ((2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}}) - \frac{1}{2} (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) \right) \\ &= -\frac{1}{2} \sum_{i=1}^n (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) - \frac{n}{2} \log |\Sigma| + C, \end{aligned}$$

where  $C = -\frac{pn}{2} \log(2\pi)$  is a constant.

$$\begin{aligned} \text{obviously, } l_n(\mu, \Sigma) &= \text{trace}(\ln(\mu, \Sigma)) \\ &= -\frac{n}{2} \text{tr} \left( \frac{1}{n} \sum_{i=1}^n (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) \right) - \frac{n}{2} \log |\Sigma| + C, \\ \left( S_n = \frac{1}{n} \sum_{i=1}^n (x_i-\mu)(x_i-\mu)^T \right) &= -\frac{n}{2} \text{tr} (\Sigma^{-1} S_n) - \frac{n}{2} \log \det(\Sigma) + C. \end{aligned}$$

(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  $(I + X)^{-1} \approx I - X$ );

(b) Suppose  $\Delta = hV$ , where  $V$  is a matrix.  $A, X$  are P.S.D.

$$\begin{aligned} \text{so } f(X + \Delta) - f(X) &= \text{tr}(A(X + hV)^{-1}) - \text{tr}(AX^{-1}) \\ &= \text{tr}(A(I + hX^{-1}V)^{-1}X^{-1}) - \text{tr}(AX^{-1}) \\ &= \text{tr} \left\{ [ (I + hX^{-1}V)^{-1} - I ] X^{-1} A \right\} \\ (\text{h is small}) \quad &\approx \text{tr} \left\{ (I - hX^{-1}V - I) X^{-1} A \right\} \\ &= -\text{tr}(X^{-1} hV X^{-1} A) = -\text{tr}(X^{-1} \Delta X^{-1} A) = -\text{tr}(X^{-1} A X^{-1} \Delta) \end{aligned}$$

$\Rightarrow f(X)$  has first-order approximation:  $f(X + \Delta) \approx f(X) - \text{tr}(X^{-1} A X^{-1} \Delta)$

$$\text{Hence } df(X)/dX = -X^{-1} A X^{-1}$$

(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  $X^{-1/2}\Delta X^{-1/2}$ );

(c) Suppose  $\Delta = hV$ ,  $h$  is small.  $X$  and  $A$  are S.P.D.

$$\begin{aligned} g(X + \Delta) - g(X) &= \log \det(X + hV) - \log \det(X) \\ &= \log \frac{\det(X) \det(I + hX^{-1}V)}{\det(X)} \\ &= \log \det(I + hX^{-1}V) \end{aligned}$$

Since  $X^{-1}, V$  are S.P.D matrix, suppose its eigenvalue decomposition is  $X^{-1}V = U\Lambda V^T$ ,  $V^TV = U^TU = I$ ,  $\Lambda$  is diagonal matrix,  $\Lambda_{ii} = t_i$ .

$$\begin{aligned} \text{then } \log \det(I + hX^{-1}V) &= \log \det(I + h\Lambda) \\ &= \log \prod_{i=1}^n (1 + ht_i) \\ &= \log (1 + h + \text{tr}(\Lambda) + o(h)) \\ &\approx h \text{tr}(\Lambda) + o(h) \\ &= h \text{tr}(X^{-1}V) + o(h) \\ &= \text{tr}(X^{-1}\Delta) + o(h) \end{aligned}$$

So  $g(X)$  has first-order approximation

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

$$\text{hence } dg(X)/dX = X^{-1}$$

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

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

(d) Obviously,  $S_n, \Sigma$  are S.P.D.

$$\text{Let } \frac{d \ln(\mu, \Sigma)}{d \Sigma} = -\frac{n}{2} \frac{d \text{tr}(S_n \Sigma^{-1})}{d \Sigma} - \frac{1}{2} \frac{dg(\Sigma)}{d \Sigma} = 0$$

$$\text{then } \frac{d \text{tr}(S_n \Sigma^{-1})}{d \Sigma} + \frac{dg(\Sigma)}{d \Sigma} = -\Sigma^{-1} S_n \Sigma^{-1} + \Sigma^{-1} = 0.$$

$$\Rightarrow -\Sigma^{-1} S_n + I = 0$$

$$\Sigma^{-1} S_n = I$$

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

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

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

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

Solution (a)

$$I(\mu) = \frac{1}{2} \sum_{i=1}^p (y_i - \mu_i)^2 + \frac{\lambda}{2} \|\mu\|^2, \quad \text{Ridge regression is to } \min_{\mu} I(\mu)$$

$$\frac{\partial I(\mu)}{\partial \mu_i} = (\mu_i - y_i) + \lambda \mu_i = 0$$

$$\Rightarrow \hat{\mu}_i^{ridge} = \frac{1}{1 + \lambda} y_i \quad \text{so } \hat{\mu}_n = \frac{1}{1 + \lambda} y, \quad \mathbb{E} \hat{\mu}_n = \frac{1}{1 + \lambda} \mathbb{E} y = \frac{1}{1 + \lambda} \mu$$

$$\text{The risk is } \mathbb{E} \|\hat{\mu}_n - \mu\|^2 = \mathbb{E} \|\hat{\mu}_n - \mathbb{E} \hat{\mu}_n\|^2 + \|\mathbb{E} \hat{\mu}_n - \mu\|^2$$

$$= \mathbb{E} \left\| \frac{1}{1 + \lambda} y - \frac{1}{1 + \lambda} \mu \right\|^2 + \left\| \frac{1}{1 + \lambda} \mu - \mu \right\|^2$$

$$= \frac{1}{(1 + \lambda)^2} \mathbb{E} [(y - \mu)^T (y - \mu)] + \left( \frac{1}{1 + \lambda} - \frac{1}{1 + \lambda} \right)^2 \mu^T \mu$$

$$= \frac{1}{(1 + \lambda)^2} \text{tr}(Z) + \left( \frac{\lambda}{1 + \lambda} \right)^2 \mu^T \mu$$

$$= \frac{p}{(1 + \lambda)^2} + \frac{\lambda^2}{(1 + \lambda)^2} \mu^T \mu$$

(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

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

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

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

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.

(b) To solve  $\min_{\mu} J(\mu) = \frac{1}{2} \|y - \mu\|_2^2 + \lambda \|\mu\|_1$ ,

we compute the subgradient of  $J(\mu)$ :

$$\partial_{\mu} J(\mu) = \mu - y + \lambda \text{sign}(\mu), \quad \text{where } \text{sign}(\mu) \begin{cases} = 1, & \mu > 0 \\ \in [-1, 1], & \mu = 0 \\ = -1, & \mu < 0 \end{cases}, \quad \text{and } 0 \in \partial_{\mu} J(\mu)$$

if  $\hat{\mu}_i > 0$ , since  $0 = \hat{\mu}_i - y_i + \lambda$ , then  $\hat{\mu}_i = y_i - \lambda > 0$ ,  $y_i > \lambda > 0$   
so  $\hat{\mu}_i = \text{sign}(y_i) (|y_i| - \lambda)_+$

if  $\hat{\mu}_i < 0$ , since  $0 = \hat{\mu}_i - y_i - \lambda$ , then  $\hat{\mu}_i = y_i + \lambda < 0$ ,  $y_i < -\lambda < 0$   
so  $\hat{\mu}_i = -(-y_i - \lambda) = \text{sign}(y_i) (|y_i| - \lambda)_+$

if  $\hat{\mu}_i = 0$ , then  $\hat{\mu}_i = 0 \in y_i + \lambda [-1, 1]$ , so  $y_i \in [-\lambda, \lambda]$ ,  $(|y_i| - \lambda)_+ = 0$

In conclusion,  $\hat{\mu}_i^{\text{soft}} = \text{Soft}(y_i; \lambda) = \text{sign}(y_i) (|y_i| - \lambda)_+$

If we choose  $\lambda = \sqrt{2 \log p}$ ,

By Stein's unbiased risk estimate,  $\hat{\mu}^{\text{soft}} = y + g(y)$ ,

$$\text{so } g_i(y) = \begin{cases} -\lambda, & y_i > \lambda \\ -y_i, & |y_i| \leq \lambda \\ \lambda, & y_i < -\lambda \end{cases}, \quad \partial_i g_i(y) = -I(|y_i| \leq \lambda),$$

$$\text{the risk } \mathbb{E} \| \hat{\mu}^{\text{soft}}(y) - \mu \|^2 = \mathbb{E} (p - 2 \sum_{i=1}^p I(|y_i| \leq \lambda)) + \sum_{i=1}^p y_i^2 \wedge \lambda^2$$

$$(\frac{1}{2} a \wedge b \leq \frac{ab}{a+b} \leq a \wedge b) \leq 1 + (2 \log p + 1) \sum_{i=1}^p \mu_i^2 \wedge 1$$

the risk of MLE is  $p$ .

$$\text{If } \mathbb{E} \left( \sum_{i=1}^p y_i^2 \wedge \lambda^2 \right) \leq 2 \mathbb{E} \left( \sum_{i=1}^p I(|y_i| \leq \lambda) \right),$$

then the risk of LASSO is smaller than that of MLE

(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

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

Rewriting  $\hat{\mu}^{\text{hard}}(y) = (1 - g(y))y$ , is  $g(y)$  weakly differentiable? Why?

(c) To solve  $\min_{\mu} J_0(\mu) = \frac{1}{2} \|y - \mu\|_2^2 + \lambda \|\mu\|_0 = \begin{cases} \frac{1}{2} \|y - \mu\|_2^2 + \lambda, & \mu > 0 \\ \frac{1}{2} \|y - \mu\|_2^2 - \lambda, & \mu < 0 \end{cases}$   
we compute the gradient of  $J(\mu)$ :

$$1) \mu^* > 0, \text{ if } y > \lambda, \text{ then } \mu^* = y$$

$$2) \mu^* < 0, \text{ if } y < -\lambda, \text{ then } \mu^* = y$$

$$\text{otherwise, } \mu^* = 0$$

$$\hat{\mu}_i^{\text{hard}} = \mu_i^* = \begin{cases} y_i, & |y| > \lambda \\ 0, & |y| < \lambda \end{cases} = y_i I(|y_i| > \lambda).$$

$$\hat{\mu}^{\text{hard}} = (1 - g(y))y, \text{ then } g(y) = 1 - I(|y| > \lambda) = I(|y| < \lambda)$$

$g(y)$  is not weakly differentiable,

because  $g$  is not absolutely continuous

(d) Consider the James-Stein Estimator

$$\hat{\mu}^{JS}(y) = \left(1 - \frac{\alpha}{\|y\|^2}\right)y.$$

Show that the risk is

$$\mathbb{E}\|\hat{\mu}^{JS}(y) - \mu\|^2 = \mathbb{E}U_\alpha(y)$$

where  $U_\alpha(y) = p - (2\alpha(p-2) - \alpha^2)/\|y\|^2$ . Find the optimal  $\alpha^* = \arg \min_\alpha U_\alpha(y)$ . Show that for  $p > 2$ , the risk of James-Stein Estimator is smaller than that of MLE for all  $\mu \in \mathbb{R}^p$ .

$$(d) \quad \hat{\mu}^{JS}(y) = (1 - \frac{\alpha}{\|y\|^2})y = y + g(y) \Rightarrow g(y) = -\frac{\alpha}{\|y\|^2}y, \quad g \text{ is weakly differentiable}$$

$$\nabla^T g(y) = -\sum_i \frac{\partial}{\partial y_i} \left( \frac{\alpha}{\|y\|^2} y \right) = -\alpha \sum_i \left( \frac{1}{\|y\|^2} - \frac{2y_i^2}{\|y\|^4} \right)$$

$$= -\alpha \left( \frac{p}{\|y\|^2} - \frac{2\|y\|^2}{\|y\|^4} \right) = \frac{-2\alpha(p-2)}{\|y\|^2}.$$

$$\|g(y)\|^2 = \frac{\alpha^2}{\|y\|^2}$$

$$\text{so } U_\alpha(y) = p - \frac{2\alpha(p-2)}{\|y\|^2} + \frac{\alpha^2}{\|y\|^2} = p + 2\nabla^T g(y) + \|g(y)\|^2$$

By Stein's unbiased risk estimate,

$$\text{the risk } \mathbb{E}\|\hat{\mu}^{JS}(y) - \mu\|^2 = \mathbb{E}U_\alpha(y)$$

$$\text{Let } \frac{\partial U_\alpha(y)}{\partial \alpha} = \frac{-2(p-2) + 2\alpha}{\|y\|^2} = 0, \quad \text{then } \alpha^* = \arg \min_\alpha U_\alpha(y) = p-2.$$

$$U_{\alpha^*}(y) = p - \frac{(p-2)^2}{\|y\|^2}$$

$$\text{For } p > 2, \quad \mathbb{E}\|\hat{\mu}^{JS}(y) - \mu\|^2 = p - \mathbb{E}\frac{(p-2)^2}{\|y\|^2} < p = R(\hat{\mu}^{\text{MLE}}, \mu)$$

(e) In general, an odd monotone unbounded function  $\Theta : \mathbb{R} \rightarrow \mathbb{R}$  defined by  $\Theta_\lambda(t)$  with parameter  $\lambda \geq 0$  is called *shrinkage rule*, if it satisfies

[shrinkage]  $0 \leq \Theta_\lambda(|t|) \leq |t|$ ;

[odd]  $\Theta_\lambda(-t) = -\Theta_\lambda(t)$ ;

[monotone]  $\Theta_\lambda(t) \leq \Theta_\lambda(t')$  for  $t \leq t'$ ;

[unbounded]  $\lim_{t \rightarrow \infty} \Theta_\lambda(t) = \infty$ .

Which rules above are shrinkage rules?

$$(e) \quad \text{For (a), } \theta_\lambda(t) = \frac{1}{1+\lambda}t, \quad 0 \leq \theta_\lambda(|t|) \leq |t|; \quad \theta_\lambda(-t) = -\theta_\lambda(t);$$

$$\theta_\lambda(t) \leq \theta_\lambda(t') \text{ for } t \leq t'; \quad \lim_{t \rightarrow \infty} \theta_\lambda(t) = \infty.$$

$$\text{For (b), } \theta_\lambda(t) = \text{sign}(t)(|t| - \lambda)_+, \quad 0 \leq \theta_\lambda(|t|) \leq |t| - \lambda \leq |t|; \quad \theta_\lambda(-t) = -\theta_\lambda(t);$$

$$\theta_\lambda(t) \leq \theta_\lambda(t') \text{ for } t \leq t'; \quad \lim_{t \rightarrow \infty} \theta_\lambda(t) = \infty$$

$$\text{For (c), } \theta_\lambda(t) = t I(|t| > \lambda), \quad 0 \leq \theta_\lambda(|t|) \leq |t|; \quad \theta_\lambda(-t) = -\theta_\lambda(t);$$

$$\theta_\lambda(t) \leq \theta_\lambda(t') \text{ for } t \leq t'; \quad \lim_{t \rightarrow \infty} \theta_\lambda(t) = \infty$$

$$\text{For (d), } \theta_\lambda(t) = (1 - \frac{\alpha}{\|t\|^2})t. \quad \theta_\lambda'(t) = 1 + \frac{\alpha}{\|t\|^2} > 0, \quad 0 \leq \theta_\lambda(|t|) \leq |t|; \quad \theta_\lambda(-t) = -\theta_\lambda(t);$$

$$\theta_\lambda(t) \leq \theta_\lambda(t') \text{ for } t \leq t'; \quad \lim_{t \rightarrow \infty} \theta_\lambda(t) = \infty$$

So (a) (b) (c) (d) are all shrinkage rules.

3. Necessary Condition for Admissibility of Linear Estimators. Consider linear estimator for  $y \sim \mathcal{N}(\mu, \sigma^2 I_p)$

$$\hat{\mu}_C(y) = Cy.$$

Show that  $\hat{\mu}_C$  is admissible only if

- (a)  $C$  is symmetric;
- (b)  $0 \leq \rho_i(C) \leq 1$  (where  $\rho_i(C)$  are eigenvalues of  $C$ );
- (c)  $\rho_i(C) = 1$  for at most two  $i$ .

These conditions are satisfied for MLE estimator when  $p = 1$  and  $p = 2$ .

*Proof.* In order to show the necessity of (a)(b)(c), we show that if the conditions fail, then there exists a dominating estimator.

(a). Let  $\tilde{A} = (A^T A)^{\frac{1}{2}}$ , then  $\text{tr}(A) \leq \text{tr}(\tilde{A})$ , with equality only if  $A^T = A$ .  
Let  $D = I - \widetilde{(I-C)}$ , then  $D^T = D$ .

Since  $(I-D)^T(I-D) = (\widetilde{I-C})^2 = (I-C)^T(I-C)$ ,  
the two estimators have the same squared bias.

$$\begin{aligned} \text{For the variance terms, } \text{tr} D^T D &= \text{tr} I - 2\text{tr}(I-D) + \text{tr}(I-D)^T(I-D) \\ &\leq \text{tr} I - 2\text{tr}(I-C) + \text{tr}(I-C)^T(I-C) \\ &= \text{tr}(C^T C). \end{aligned}$$

If  $C$  is not symmetric,  $\text{tr} D^T D < \text{tr}(C^T C)$ .

From the variance-bias decomposition, we know  $\hat{\mu}_C$  is not admissible,  
it can be dominated by  $\hat{\mu}_D$ . So (a) is necessary.

(b) Now assume  $C$  is symmetric, consider the eigenvalue decomposition  $C = U \Lambda U^T$ ,  
where  $\Lambda = \text{diag}(\lambda_{ii})$ ,  $UU^T = U^T U = I_p$ .  $\lambda_{ii} = \rho_i(C)$ .

$$\text{Let } \eta = U^T \mu, \quad x = U^T y \sim N(\eta, \sigma^2 I_p)$$

$$\text{Since } E \|Cy - \mu\|^2 = E \|\Lambda - \eta\|^2,$$

$$\text{we have } R(\hat{\mu}_C, \mu) = R(\hat{\eta}_{\Lambda}, \eta) = \sum_{i=1}^p \sigma^2 \lambda_{ii}^2 + (1-\lambda_{ii})^2 \eta_i^2 = \sum_{i=1}^p R(\lambda_{ii}, \eta_i)$$

If  $\exists \lambda_{ii} \notin [0, 1]$ , replace  $\lambda_{ii}$  by 1 if  $\lambda_{ii} > 1$  and by 0 if  $\lambda_{ii} < 0$ ,  
then the new estimator dominates  $\hat{\mu}_C$ .

So (b) is also necessary

(c) Suppose  $\lambda_1 = \dots = \lambda_d = 1 > \lambda_{d+1}, \dots, \lambda_p$ . Let  $x^d = (x_1, \dots, x_d)$

Note that the positive part of J-S estimator is everywhere better than  $\hat{\eta}_J(x^d) = x^d$ .

Define a new estimator  $\hat{\eta}$  to use  $\hat{\eta}^{JS}$  on  $x^d$ , use  $\lambda_i x_i$  for  $i > d$

then  $R(\hat{\eta}, \eta) = R(\hat{\eta}^{JS}, \eta^d) + \sum_{i>d} r(\lambda_i, \eta_i) < r(\lambda, \eta)$   
so  $\hat{\eta}$  dominates  $\hat{\eta}_\lambda$  and hence  $\hat{\theta}_c$ .