

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 μ and Σ ;

$$(a) f(x) = \frac{1}{|2\pi|^{p/2} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}.$$

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

$$\text{i.e. } l_n(\mu, \Sigma) = -\frac{np}{2} \log(2\pi) - \frac{n}{2} \log|\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

By cyclicity, $(x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \text{trace}(\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T) = n \cdot \text{trace}(\Sigma^{-1} S_n)$

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

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

$$\begin{aligned} (b) f(X + \Delta) &= \text{trace}(A(X + \Delta)^{-1}) = \text{trace}(A[(I + \Delta X^{-1})X]^{-1}) \\ &= \text{trace}(AX^{-1}(I + \Delta X^{-1})) \approx \text{trace}(AX^{-1}(I - \Delta X^{-1})) \\ &= \text{trace}(AX^{-1}) - \text{trace}(AX^{-1}\Delta X^{-1}) \end{aligned}$$

$$\begin{aligned} \text{trace}(AX^{-1}\Delta X^{-1}) &= \text{trace}(X^{-1}\Delta' X^{-1} A') = \text{trace}(X^{-1}\Delta X^{-1} A') \\ &= \text{trace}(X^{-1} A' X^{-1} \Delta) \end{aligned}$$

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

$$\text{Thus, } \frac{df(X)}{dX} = -X^{-1} A X^{-1}.$$

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

$$\begin{aligned} (c) g(X + \Delta) &= \log(\det(X + \Delta)) = \log \left\{ \det \left(X^{\frac{1}{2}} (I + X^{\frac{1}{2}} \Delta X^{-\frac{1}{2}}) X^{\frac{1}{2}} \right) \right\} \\ &= \log \det(X) + \log \det(I + X^{\frac{1}{2}} \Delta X^{-\frac{1}{2}}) \\ &= \log \det(X) + \sum_{i=1}^n \log(1 + \lambda_i), \text{ where } \lambda_i \text{ are eigenvalues of } X^{\frac{1}{2}} \Delta X^{-\frac{1}{2}}. \end{aligned}$$

Since Δ is small, so are its eigenvalues λ_i 's,

thus $\log(1+\lambda_i) \approx \lambda_i$, and

$$\begin{aligned} g(X+\Delta) &= g(X) + \sum_{i=1}^n \log(1+\lambda_i) = g(X) + \sum_{i=1}^n \lambda_i = g(X) + \text{trace}(X^{-\frac{1}{2}} \Delta X^{-\frac{1}{2}}) \\ &= g(X) + \text{trace}(X^{-1} \Delta) \end{aligned}$$

Hence, $\frac{d g(X)}{d X} = X^{-1}$

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

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

A reference for (b) and (c) can be found in Convex Optimization, by Boyd and Vandenberghe, examples in Appendix A.4.1 and A.4.3:

https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf

$$\begin{aligned} (\text{d}), \quad \frac{\partial \ln(\mu \cdot \Sigma)}{\partial \Sigma} &= \frac{\partial}{\partial \Sigma} (-\frac{1}{2} \text{trace}(\Sigma^{-1} S_n) - \frac{1}{2} \cdot \log \det(\Sigma)) \\ &= -\frac{1}{2} \cdot (-\Sigma^{-1} S_n \Sigma^{-1}) - \frac{1}{2} \Sigma^{-1} = 0 \end{aligned}$$

Thus, $\hat{\Sigma} = S_n$

2. Shrinkage : Suppose $y \sim N(\mu, I_p)$

- (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 $C = I$.

$$\begin{aligned} (\alpha) \text{ let } R(\mu) &= \frac{1}{2} (y - \mu)^T (y - \mu) + \frac{\lambda}{2} \mu^T \mu \\ &= \frac{1}{2} y^T y - y^T \mu + \mu^T \mu (1 + \lambda) \end{aligned}$$

$$0 = \frac{\partial R}{\partial \mu} = -y + (1 + \lambda) \mu \quad \text{gives} \quad \mu = \frac{1}{1+\lambda} y$$

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

$$\hat{\mu} = \frac{1}{1+\lambda} C y \quad \text{where } C = \frac{1}{1+\lambda} I.$$

$$MSE = \text{Var}(\hat{\mu}) + \text{bias}^2(\hat{\mu}), \text{ where } \text{bias}^2(\hat{\mu}) = \|(\mathbb{I} - (I + \lambda)^{-1})y\|_2^2 = 0$$

$$\text{Var}(\hat{\mu}) = \sigma^2 \text{tr}(C^T C) = \frac{p\sigma^2}{(1+\lambda)^2} = \frac{p}{(1+\lambda)^2}.$$

(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 μ , such a risk is smaller than that of MLE? Note: see Gaussian Estimation by Iain Johnstone, Lemma 2.9 and the reasoning before it.

$$(b) \text{ Let } R(\mu) = \frac{1}{2} (y - \mu)^T (y - \mu) + \lambda \|\mu\|_1.$$

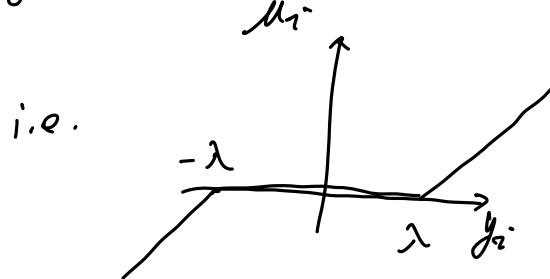
$$\text{Let } 0 = \frac{\partial R(\mu)}{\partial \mu_i} = -y_i + \mu_i + \lambda \cdot \text{sgn}(\mu_i)$$

$$\textcircled{1} \quad \mu_i > 0, \quad \mu_i + \lambda = y_i$$

$$\textcircled{2} \quad \mu_i = 0, \quad y_i = \mu_i = 0$$

$$\textcircled{3} \quad \mu_i < 0, \quad \mu_i = y_i + \lambda$$

$$\text{So } \hat{\mu}_i^{soft} = \text{sgn}(y_i) \cdot (|y_i| - \lambda)_+$$



Let $r(\hat{\mu}_i, \mu_i) = \mathbb{E} \|\hat{\mu}_i - \mu_i\|^2 = \int \|\hat{\mu}_i - \mu_i\|^2 \phi(z) dz$, where $z \sim N(0, 1)$

$$\|\hat{\mu}_i - \mu_i\|^2 = \begin{cases} (y_i - \lambda - \mu_i)^2, & y_i > \lambda \\ \mu_i^2, & y_i = \lambda \\ (y_i + \lambda - \mu_i)^2, & y_i < \lambda \end{cases} \xrightarrow{y_i - \mu_i = z} \begin{cases} (z - \lambda)^2, & z > \lambda - \mu_i \\ \mu^2, & z = \lambda - \mu_i \\ (z + \lambda)^2, & z < \lambda - \mu_i \end{cases}$$

$$\text{Let } \mu_i \rightarrow \infty. \quad r(\hat{\mu}_i, \mu_i) = 1 + \lambda^2.$$

$$\frac{\partial}{\partial \mu_i} r(\hat{\mu}_i, \mu_i) = 2\mu_i \cdot \mathbb{P}(|\mu_i + z| \leq \lambda) \leq 2\mu_i, \text{ thus } r(\hat{\mu}_i, \mu_i) \leq 1 + \lambda^2.$$

On the other hand,

$$r(\hat{\mu}_i, \mu_i) - r(\hat{\mu}_i, 0) \leq \int \min(1 + \lambda^2, \mu_i^2) \phi(z) dz = \min(1 + \lambda^2, \mu_i^2) \leq (1 + \lambda^2) \min(\mu_i^2, 1)$$

$$\text{and } r(\hat{\mu}_i, 0) \leq 2\lambda^{-1} \phi(\lambda) \leq e^{-\lambda^2/2} \leq 1$$

$$\text{Thus. } r(\hat{\mu}, \mu) = \sum_{i=1}^p r(\hat{\mu}_i, \mu_i) \leq \sum_{i=1}^p (1 + (1 + \lambda^2) \min(\mu_i^2, 1))$$

$$\leq 1 + \sum_{i=1}^p (1+\lambda^2) \min(\mu_i^2, 1)$$

Let $\lambda = \sqrt{2 \log p}$, then $r(\hat{\mu}, \mu) \leq 1 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, 1)$.

The condition of $r(\hat{\mu}^{\text{soft}}, \mu^{\text{MLE}})$ is

(c) Consider the l_0 regularization

$$\min_{\mu} \|y - \mu\|_2^2 + \lambda^2 \|\mu\|_0, \quad 1 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, 1) < p.$$

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?

$$(c) \text{ Let } R(\mu) = \|y - \mu\|_2^2 + \lambda^2 \|\mu\|_0 = (y - \mu)^T (y - \mu) + \lambda^2 \left(\sum_{i=1}^p I\{\mu_i \neq 0\} \right)$$

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

$$\text{Let } \sigma = \frac{\partial R(\mu_i)}{\partial \mu_i} = \begin{cases} 2(y_i - \mu_i) & \mu_i \neq 0 \\ 0 & \mu_i = 0 \end{cases}$$

When $y_i^2 < (y_i - \mu_i)^2 + \lambda^2$ is always true, $\Delta = 4y_i^2 - 4\lambda^2 < 0$, or $|y_i| < \lambda$

this time, set $\mu_i = 0$. Otherwise, let $\mu_i = y_i$.

Thus, $\hat{\mu}_i^{\text{hard}} = y_i \cdot \mathbb{1}\{|y_i| > \lambda\}$.

$g(y) = 1 - \mathbb{1}\{|y| > \lambda\} = \mathbb{1}\{|y| \leq \lambda\}$ is not weakly differentiable since $g(y)$ 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 $\mu \in \mathbb{R}^p$.

$$(d) \hat{\mu}^{JS}(y) = (1 - \frac{2}{\|y\|^2})y, \quad g(y) = -\frac{2}{\|y\|^2} \cdot y$$

Thus $U(y) = p + 2\mathbf{g}^T g(y) + \|g(y)\|^2$, and

$$\nabla g^T(y) = \sum_{i=1}^P \frac{\partial}{\partial y_i} g(y) = -2(p-2)/\|y\|^2$$

$$\|g(y)\|^2 = \left\| -\frac{2}{\|y\|^2} \cdot y \right\|^2 = \frac{2^2}{\|y\|^2} \cdot \|y\|^2 = \frac{2^2}{\|y\|^2}$$

thus, $U_2(y) = p - (2(p-2) - 2^2)/\|y\|^2$. Risk $R(\hat{\mu}^{JS}, \mu) = E U_2(y)$.

$$\text{Let } D = \frac{\partial U_2(y)}{\partial \alpha} = \frac{1}{\|y\|^2} (2(p-2) - 2\alpha), \quad \alpha^* = p-2$$

$$\text{Risk } R(\hat{\mu}^{JS}, \mu) = E U_2(y) = E(p - (p-2)^2/\|y\|^2) < p = R(\hat{\mu}^{MLE}, \mu), \quad p > 2.$$

- (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). All of the four estimators are shrinkage rules.

We can check these conditions by drawing the picture of $\Theta(y) - y$.

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 $p = 2$.

Reference: Theorem 2.3 in Gaussian Estimation by Iain Johnstone,
<http://statweb.stanford.edu/~imj/Book100611.pdf>

Proof. \Leftarrow If the condition fails, we show a dominating estimator.

If (1) fails:

Lemma: $|A| = (A^T A)^{\frac{1}{2}}$ and $\text{tr}(A) \leq \text{tr}(|A|)$, with equality only if $A = A^T$.

Let $I - D = |I - C|$, D is symmetric.

Since $(I - D)^T(I - D) = |I - C|^2 = (I - C)^T(I - C)$, so $\hat{\theta}_p$ and $\hat{\theta}_c$ has the same bias.

Write $\text{tr } D^T D = \text{tr } I - 2\text{tr}(I - D) + \text{tr}(I - D)^T(I - D)$.

and $\text{tr } D^T D < \text{tr } C^T C$ iff $\text{tr}(I - D) = \text{tr}|I - C| > \text{tr}(I - C)$, this occurs only if C is not symmetric. So C is inadmissible.

If (2) fails:

Now C is symmetric, let $C = U\Lambda U^T$ with U orthogonal and $\Lambda = \text{diag}(\lambda_i)$

Let $\eta = U^T \mu$, $x = U^T y \sim N(\eta, \sigma^2 I_p)$. and $E\|Cy - x\|^2 = E\|\Lambda - \eta\|^2$.

We have $r(\hat{\mu}_c, \mu) = r(\hat{\eta}_1, \eta) = \sum_i \sigma^2 \lambda_i^2 + (1 - \lambda_i)^2 \eta_i^2 = \sum_i r(\lambda_i, \eta_i)$

If $\lambda_i \notin [0, 1]$, a strictly better LSE results by replacing λ_i by 1 if $\lambda_i > 1$ and by 0 if $\lambda_i < 0$.

If (3) fails:

Suppose that $\lambda_1 = \dots = \lambda_d = 1 > \lambda_{d+1}$ for $i > d \geq 3$. Let $x^d = (x_1, \dots, x_d)$

Noted that the positive part James-Stein estimator is even better than $\hat{\eta}(x^d) = x^d$

So define $\hat{\eta}^{JS}$ on x^d and use $\lambda_i x_i$ for $i > d$, denote by $\hat{\eta}$.

$r(\hat{\eta}, \eta) = r(\hat{\eta}^{JS}, \eta^d) + \sum_{i>d} r(\lambda_i x_i, \eta_i) < r(\lambda, \eta)$. Thus, $\hat{\eta}$ domain $\hat{\eta}_n$. hence, $\hat{\eta}_c$.

Thus, condition (i)-(iii) implied $\hat{\mu}_c$ is admissible.

(\Rightarrow) is similar to (\Leftarrow), just from conclusion to condition.

As for MLE, $C = Z_p$ is symmetric, $\lambda_i = 1$.

As for $p=1, 2$, at most two $\lambda_i = 1$. So $\hat{\mu}_{MLE}$ is admissible.