

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 ;

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

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

$$\begin{aligned} f(X + \Delta) &= \text{trace}(A(X + \Delta)^{-1}) = \text{trace}(A[(I + \Delta X^{-1})X]^{-1}) \\ &= \text{trace}(A X^{-1} (I + \Delta X^{-1})^{-1}) \quad (I + \Delta X^{-1})^{-1} \approx I - \Delta X^{-1} \\ &\approx \text{trace}(A X^{-1}) - \text{trace}(A X^{-1} \Delta X^{-1}) \quad \text{trace}(a) = \text{trace}(a') \\ &= \text{trace}(A X^{-1}) - \text{trace}((X^{-1})^T \Delta^T (X^{-1})^T A^T) \\ &= \text{trace}(A X^{-1}) - \text{trace}(X^{-1} \Delta X^{-1} A') \\ &= \text{trace}(A X^{-1}) - \text{trace}(A' X^{-1} \Delta X^{-1}) \quad \text{cyclicity} \\ &= \text{trace}(A X^{-1}) - \text{trace}(X^{-1} A' X^{-1} A) \\ &= f(X) - \text{trace}(X^{-1} A' X^{-1} \Delta) \end{aligned}$$

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

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

where λ_i is the i -th eigenvalue of $X^{-1/2} \Delta X^{-1/2}$

Since Δ is small λ_i 's are small

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

$$\begin{aligned} \text{Hence } \log \det(I + \Delta) &\approx \log \det(X) + \sum_{i=1}^n \lambda_i \\ &= \log \det(X) + \text{trace}(X^{-1/2} \Delta X^{-1/2}) \\ &= \log \det(X) + \text{trace}(X^{-1} \Delta) \\ \Rightarrow g(X + \Delta) &= g(X) + \text{trace}(X^{-1} \Delta) \end{aligned}$$

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

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

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

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

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

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

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

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

$$\begin{aligned}
 \text{RSS}(\mu) &= \frac{1}{2}(y - \mu)^T(y - \mu) + \frac{\lambda}{2}\mu^T\mu \\
 &= \frac{1}{2}y^Ty - 2y^T\mu + \mu^T\mu + \frac{\lambda}{2}\mu^T\mu \\
 0 = \frac{\partial \text{RSS}(\mu)}{\partial \mu} &= -\frac{1}{2}y + \frac{1+\lambda}{2}\mu
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow \hat{\mu}^{\text{ridge}} &= \frac{1}{1+\lambda}y \\
 \Rightarrow \hat{\mu}_i^{\text{ridge}} &= \frac{1}{1+\lambda}y_i
 \end{aligned}$$

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

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

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

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

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

$$\text{Var}(\hat{\mu}_c) = b^2 \text{trace}(C^TC) = 1 \cdot 1^p = p$$

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

(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

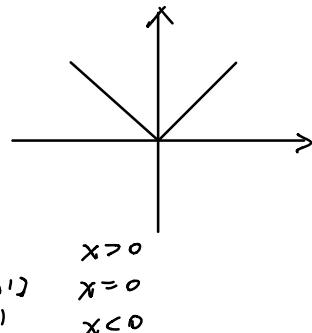
$$\mathbb{E}\|\hat{\mu}^{\text{soft}}(y) - \mu\|^2 \leq 1 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, 1).$$

Under what conditions on 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.

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

for $\|\mu\|_1$ refers to $\|x\|_1$,

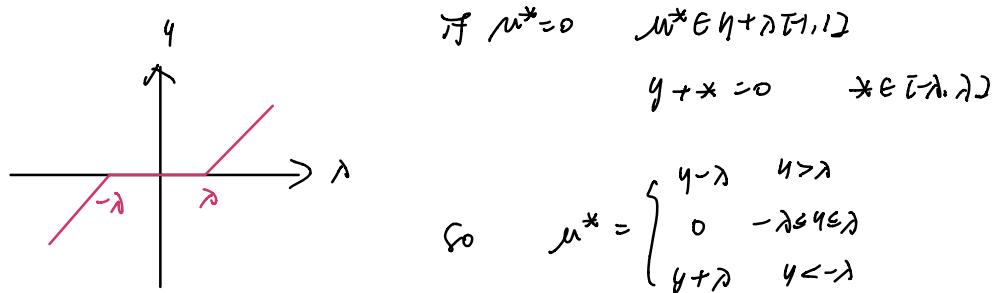
$$\partial\|x\|_1 = \text{sign}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$



$$\text{thus } \frac{\partial J(\mu)}{\partial \mu} = \mu^* y + \lambda \operatorname{sgn}(\mu^*)$$

$$0 \in \partial J(\mu) \text{ discards } \begin{cases} \text{if } \mu^* > 0 & \mu^* - y + \lambda = 0 \\ & \mu^* = y - \lambda > 0 \end{cases} \quad y > \lambda$$

$$\begin{cases} \text{if } \mu^* < 0 & \mu^* - y - \lambda = 0 \\ & \mu^* = y + \lambda \end{cases} \quad y + \lambda < 0$$



$$\mu^{\text{soft}} = \operatorname{sgn}(y) / (|y| - \lambda)^+$$

Then we found the upperbound of the risk.

Suppose $\hat{\mu}_\lambda(y) = y + g_\lambda(y)$ were $g_\lambda(y)$ for soft thresholding

$$\text{Consider } y = \mu + z \sim N(0, 1)$$

$$\text{the risk function } r_\delta(\lambda, \mu) = \int [\hat{\mu}_\lambda(\mu + z) - \mu]^2 \phi(z) dz$$

$$[\hat{\mu}_\lambda(\mu + z) - \mu]^2 = \begin{cases} (z + \lambda)^2 & \mu < -\lambda - z \\ \mu^2 & \text{otherwise.} \\ (z - \mu)^2 & \mu > \lambda - z \end{cases}$$

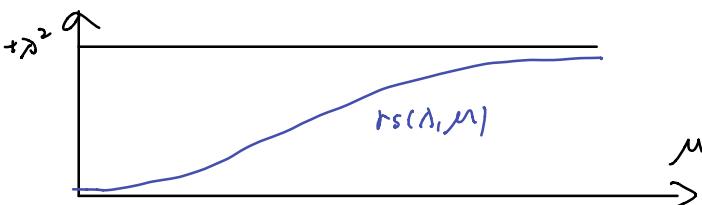
$$r_\delta(\lambda, 0) \leq 2\lambda^{-1} \phi(\lambda) \leq e^{-\lambda^2/2}$$

$$r_\delta(\lambda, \infty) = 1 + \lambda^2$$

$$\frac{\partial r_\delta(\lambda, \mu)}{\partial \mu} = 2\mu P(|\mu + z| \leq \lambda) \leq 2\mu$$

$$\frac{\partial}{\partial \mu}$$

Thus,



$$RS(\lambda, \mu) - RS(\lambda, 0) \leq \mu^2$$

$$RS(\lambda, \mu) \leq RS(\lambda, 0) + \min(\mu^2, 1/\lambda^2) \quad \text{choose } \lambda = \sqrt{2 \log p}$$

$$RS(\lambda, \mu) \leq RS(\lambda, 0) + (2 \log p + 1) \min(\mu^2, 1)$$

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

$$Risk(\hat{\mu}, \mu) \leq \sigma^2 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, \sigma^2)$$

for $\epsilon = 1$ we get the conclusion for this exercise.

$$\| \hat{\mu}^{soft}(\mu) - \mu \|_2^2 \leq 1 + (2 \log p + 1) \sum_{i=1}^p \min(\mu_i^2, 1)$$

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

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

(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

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

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

$$\begin{aligned} RSS(\mu) &= (y - \mu)^T (y - \mu) + \lambda^2 \left(\sum_{i=1}^p I(\mu_i \neq 0) \right) \\ &= y^T y - y^T \mu - \mu^T y + \mu^T \mu + \lambda^2 \left(\sum_{i=1}^p I(\mu_i \neq 0) \right) \end{aligned}$$

for $\mu_i \neq 0$

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

When is $y_i^2 < (y_i - \mu_i)^2 + \lambda^2$ is always true?

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

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

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

$$\Leftrightarrow y_i^2 < \lambda^2$$

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

otherwise $\mu_i = y_i$

$$\begin{aligned} \hat{\mu}_{\text{hard}}(y_i) &= (1 - g(y_i)) y_i \\ &= (1 - \mathbb{I}(|y_i| \leq \lambda)) y_i = y_i \mathbb{I}(|y_i| > \lambda) \end{aligned}$$

$g(y) = \mathbb{I}(|y| \leq \lambda)$ 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 $\mu \in \mathbb{R}^p$.

$$\hat{\mu} \stackrel{\Delta}{=} \hat{\mu}^{TS}(y) = (1 - \frac{\alpha}{\|y\|^2}) y = y - \frac{\alpha}{\|y\|^2} y$$

$$g(y) = -\frac{\alpha}{\|y\|^2} y \text{ satisfies.}$$

1. g is weakly differentiable

$$2. \sum_{i=1}^p \int \left| \frac{\partial g_i(y)}{\partial y_i} \right| dy < \infty$$

Thus $U(Y) \triangleq p + 2D^T g(Y) + \|g(Y)\|^2$ we have

$$\text{Risk } \mathbb{E}\|\hat{\mu}^{TS}(Y) - \mu\|^2 = \mathbb{E}U_\alpha(Y)$$

$$\text{where } D^T g(Y) = \sum_{i=1}^p \frac{\partial}{\partial y_i} g_i(Y) = -\frac{\alpha(p-2)}{\|Y\|^2}$$

$$\|g(Y)\|^2 = \left\| -\frac{\alpha}{\|Y\|^2} Y \right\|^2 = \alpha^2 \frac{1}{\|Y\|^2}$$

$$\text{Thus } U(Y) = p - (2\alpha(p-2) - \alpha^2)/\|Y\|^2$$

$$D = \frac{\partial U(Y)}{\partial \alpha} = -(2\alpha(p-2) - 2\alpha)/\|Y\|^2$$

$$\Rightarrow \text{the optimal } \alpha^* = \arg \min_\alpha U(Y) = p-2$$

$$\text{So } \alpha(y) = p - (2(p-2)^2 - (p-2)^2) / \|y\|^2 \\ = p - (p-2)^2 / \|y\|^2$$

Aim:

$$\text{Risk}(\hat{\mu}_{TS}, \mu) = p - \mathbb{E} \frac{(p-2)^2}{\|y\|^2} < p = \text{Risk}(\hat{\mu}_{MLE}, \mu)$$

Find the upper bound of the risk of James-Stein.

estimator $\|y\|^2 \sim \chi^2(\|y\|^2, p)$

$$T(a, b) \quad f(x; a, b) = \frac{\beta^a}{\Gamma(a)} x^{a-1} e^{-\beta x}$$

$$\chi^2(k) \quad f(x; k) = \frac{1}{(\frac{k}{2})^{\frac{k}{2}} \Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} e^{-x}$$

$$\chi^2(k) \sim T(\frac{k}{2}, \frac{1}{2})$$

$$\|y\|^2 \sim \chi^2(\|y\|^2, p) \sim T(\frac{\|y\|^2}{2}, \frac{p}{2}) \sim T(0, \frac{p}{2} + N)$$

$$N \sim \text{Poisson}(\frac{\|y\|^2}{2})$$

$$\text{a.s.a. } Y \sim \chi^2(0, p+2N)$$

$$\text{So, } \mathbb{E}(\frac{1}{\|y\|^2}) = \mathbb{E}_N \mathbb{E}_Y \left[\frac{1}{\|y\|^2} \mid N \right] = \mathbb{E}(\frac{1}{p+2N-2})$$

$$\begin{aligned} \text{Tensen's ineq} \quad & \geq \frac{1}{p+2N-2} \\ & = \frac{1}{p+\frac{\|y\|^2}{2}-2} \end{aligned}$$

Thus for $p \geq 3$

$$\text{Risk}(\hat{\mu}_{TS}, \mu) < p - \frac{(p-2)^2}{p-2+\frac{\|y\|^2}{2}} = 2 + \frac{(p-2)\|y\|^2}{p-2+\frac{\|y\|^2}{2}} < p$$

- (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?

(a) Plot the graph of $\mu(q)$

We see clearly that all estimators
above Ridge, LASSO, l_0 -Regularizer, James-Stein.
satisfy shrinkage, odd, monotone, unbounded

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 $p = 2$.

Lemma 2.1 Let C be any matrix and let P be any orthogonal matrix. Then:

$P^T C P$ admissible iff Cy is admissible

Let $L(\psi(\theta), \delta(y))$ be the loss for estimating $\psi(\theta)$ by $\delta(y)$

then Risk $R(\psi, \delta; \theta) = E_\theta L(\psi(\theta), \delta(y))$

Suppose P is orthogonal if ψ_1, δ_1 given

define $\psi_2(\theta) = \psi_1(P\theta)$ $\delta_2(y) = \delta_1(Py)$

$$\begin{aligned} R(\psi_2, \delta_2, \theta) &= E_\theta L(\psi_1(P\theta), \delta_1(Py)) \\ &= R(\psi_1, \delta_1; P\theta) \\ &= R(\psi_1, \delta_1; \theta) \end{aligned}$$

Thus. ψ_2 is admissible for δ_2

If ψ_1 is admissible for δ_1

Set $\psi_1(\theta) = \theta$ $\delta_1(y) = Cy$

so Cy is admissible for $P\theta$ iff

Cy is admissible for θ

$$L(p\psi_1, P\delta) = L(\psi_1, \delta)$$

$\Rightarrow p'Cy$ is admissible for θ iff

Cy is admissible for $p\theta$

$\Rightarrow p'Cy$ is admissible for θ iff

Cy is admissible for θ

$$C = P^T D P \quad D = \text{diag}(d_1, d_2, \dots)$$

Dy is the unique Bayes solution w.r.t. prior distribution

$$dg(\theta) = (2\pi)^{-p/2} \prod_{i=1}^p \lambda_i^{-1/2} e^{-\lambda_i \theta_i^2/2} \prod_{i=1}^p d\theta_i$$

$$\text{where } \lambda_i = \frac{(1-d_i)}{d_i}$$

Following replacement leads to better estimator

$$d_i > 1 \rightarrow d_i > 1$$

$$d_i < 0 \rightarrow -d_i$$

$$\text{three } d_i = 1 \rightarrow (1 - (k-2)) / \sum_{i=1}^p q_i^2$$

thus, it only remains to prove Dy is admissible

if one or two $d_i = 1$ while others $0 < d_i < 1$

Let. $d_1 = d_2 = 1 \quad 0 < d_i < 1 \quad i = 3, 4, \dots$

Suppose such a Dy is inadmissible

Then, there exists an estimate $H(y) = (h_1(y), \dots, h_p(y))$
 s.t. $H(y)$ is better than Dy

$$\rho(D, H(y)) \leq \rho(D, Dy)$$

$$\sum_{i=1}^p \int_{\Omega} (h_i(y) - \theta_i)^2 d\psi(y; \theta, \Sigma) \leq 2 + \sum_{i=1}^p (\alpha_i^2 + \theta_i^2 (\alpha_i - 1)^2)$$

Multiplying both sides by

$$\varphi(\theta) = \prod_{i=1}^p (\lambda_i / 2\pi)^{1/2} e^{-\lambda_i \theta_i^2 / 2}$$

$$\Rightarrow \int_{\Omega} [(\lambda_1 y_1 - \theta_1)^2 + (\lambda_2 y_2 - \theta_2)^2] d\psi(y; \theta, \Sigma) \leq 2$$

$$\text{where } M' = (\theta_1, \theta_2, \dots, \theta_p) \quad \Sigma = \text{diag}(1, 1, \frac{\lambda_1^2}{2}, \dots)$$

Thus, Dy is admissible \Rightarrow dy is admissible