Saheli Bhattacharya
Stuid = 20999867

CSIC 5011 HW (3rd week)

**8)1)** $X_i \in \mathbb{R}^p \sim N(\mu, \varepsilon)$, $i = 1, 2, \cdots n$

**a)** $\ell n(\mu, \varepsilon) = -\frac{n}{2} \text{trace}(\varepsilon^{-1} S_n) - \frac{n}{2} \log|\varepsilon| + c$ —— Ⓐ

using the P.d.f of multivariable gaussian:

$$f(x_i | \mu, \varepsilon) = \frac{1}{(2\pi)^{1/2} |\varepsilon|^{1/2}} \exp\left\{-\frac{1}{2}(x_i - \mu)^T \varepsilon^{-1}(x_i - \mu)\right\}$$

since $x_i$ are i.i.d, taking the products for $i = 1, \cdots n$ &

then taking the log gives:

$$\ell n(\mu, \varepsilon) = \underbrace{-\frac{p}{2} \times n \log(2\pi)}_{C = \text{independent of } \mu \text{ \& } \varepsilon} - \frac{n}{2} \log|\varepsilon| - \frac{1}{2} \sum_{i=1}^{n}(x_i - \mu)^T \varepsilon^{-1}(x_i - \mu)$$

$$\Rightarrow \ell n(\mu, \varepsilon) = -\frac{n}{2}\log|\varepsilon| - \frac{n}{2}\text{trace}(\varepsilon^{-1} S_n) + C \quad \left\{ \begin{array}{l} \text{using trace}(AB) \\ = \text{trace}(BA) \end{array} \right.$$

**(b)** $f(x) = \text{trace}(Ax^{-1})$, $A, x \geq 0$.

To show: $f(x + \Delta) = f(x) - \text{trace}(x^{-1} A x^{-1} \Delta)$

Let $z = x + \Delta$

$$f(z) = \text{trace}(A(x + \Delta)^{-1}) = \text{trace}(A x^{-1}(I + x^{-1}\Delta)^{-1})$$

$$= \text{trace}(A(x(I + x^{-1}\Delta))^{-1})$$

$$= \text{trace}(A(I + x^{-1}\Delta)^{-1} x^{-1}) \qquad \text{using } (I + x^{-1}\Delta)^{-1} \approx$$

$$= \text{trace}(A(I - x^{-1}\Delta) x^{-1}) \qquad \qquad I - x^{-1}A$$

$$= \text{trace}(A x^{-1} - A x^{-1}\Delta x^{-1})$$

$$f(x + \Delta) = \text{trace}(Ax^{-1}) - \text{tr}(x^{-1} A x^{-1} A) \quad \left\{ \begin{array}{l} \text{using trace}(AB) \\ = \text{trace}(BA) \end{array} \right.$$

$$\Rightarrow \frac{df(x)}{dx} = -x^{-1} A x^{-1} \qquad \left[ f(z) = f(x) + \nabla f(x)^T(z - x) \right]$$

(c)  $g(x) = \log \det(x)$

let $z = x + \delta$

$$g(z) = \log |x + \delta|$$

$$= \log \left| x^{1/2} \left( I + x^{-1/2} \delta x^{-1/2} \right) x^{1/2} \right|$$

$$= \log \left| x \left( I + x^{-1/2} \delta x^{-1/2} \right) \right| \qquad \left[ \text{using } |AB| = |BA| \right]$$

$$= \log |x| + \log \left| I + x^{-1/2} \delta x^{-1/2} \right|$$

$$= \log |x| + \sum_{i=1}^{n} \log(1 + \lambda_i) \quad, \text{ where } \lambda_i \to i^{th}$$

$$\text{eigen value of}$$
$$x^{-1/2} \delta x^{-1/2}$$

since $\delta$ is small $\Rightarrow \lambda_i$ are small

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

$$\Rightarrow \log|z| = \log|x| + \sum_{i=1}^{n} \lambda_i = \log|x| + \text{trace}\left( x^{-1/2} \delta x^{-1/2} \right)$$

$$= \log|x| + \underbrace{\text{trace}\left( x^{-1} \delta \right)}_{\perp} \qquad \left[ \text{using } \text{trace}(AB) = \text{trace}(BA) \right]$$

$$\text{inner product b/w } x^{-1} \text{ \& } \delta.$$

$$f(z) \approx f(x) + \text{trace}\left( x^{-1} (z - x) \right)$$

$$\Rightarrow \frac{d g(x)}{d z} = x^{-1}$$

d)  $$\ln(\mu, z) = -\frac{n}{2} \text{trace}\left( \Sigma^{-1} S_n \right) - \frac{n}{2} \log|\Sigma| + C$$

$$= -\frac{n}{2} \text{trace}\left( \Sigma^{-1} \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^T (x_i - \mu) \right)$$

$$- \frac{n}{2} \log|\Sigma| + C$$

$$= \frac{-1}{2} \sum_{i=1}^{n} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) - \frac{n}{2} \log |\Sigma| + C .$$

$$\frac{\partial \ln(\mu, \Sigma)}{\partial \mu} = \sum_{i=1}^{n} \Sigma^{-1} (x_i - \mu) = 0$$

$$\Rightarrow \mu^{-1}_{m_L} = \frac{1}{n} \sum_{i=1}^{n} x_i .$$

for $\bar{\Sigma}^{-1}_{m_L}$

$$\frac{\partial}{\partial \Sigma} \ln(\mu, \Sigma) = -\frac{n}{2} \left( -\Sigma^{-1} S_n \Sigma^{-1} \right) = -\frac{n}{2} \Sigma^{-1} > 0$$

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

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

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

$$\bar{\Sigma}^{-1}_{m_L} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu}_{m_L})^T (x_i - \hat{\mu}^{-1}_{m_L})$$

**q2)** $y \sim N(\mu, I_p)$

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

Setting the derivative to 0

$$\frac{1}{2} \cdot 2 (y - \mu)(-1) + \frac{\lambda}{2} \times 2 y = 0 .$$

$$\mu (\lambda + 1) = y \Rightarrow \mu = \frac{y}{\lambda + 1}$$

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

$$\text{Var}(\hat{u}^{-1}) = \left[ \left( \frac{1}{1+\lambda} y - \frac{1}{1+\lambda} u \right)^T \left( \frac{y}{1+\lambda} - \frac{u}{1+\lambda} \right) \right]$$

$$= \left( \frac{1}{1+\lambda} \right)^2 \times \left[ (y-u)^T (y-u) \right] = \frac{P}{(1+\lambda)^2}$$

$$\text{Bias}(\hat{u}) = \frac{u - u}{1+\lambda} = \frac{\lambda}{1+\lambda} u.$$

$$\text{Risk} = \text{Var}(\hat{u}) + \text{Bias}(\hat{u}))^2 = \frac{P}{(1+\lambda)^2} + \frac{\lambda^2}{(1+\lambda^2)} \times \|u\|^2$$

(b) $\quad \min_u \frac{1}{2} \|y - u\|^2 + \lambda \|u\|_1 = f(u)$

for minimizing $\Rightarrow 0 \in \partial f(u)$.

$$\partial f(u) = (u-y) + \lambda \, \text{sign}(u)$$

if $u_i > 0$

$\Rightarrow u_i - y_i + \lambda = 0 \Rightarrow u_i = y_i - \lambda$ if $y_i > \lambda_i$

if $u_i < 0$

$\Rightarrow u_i - y_i - \lambda = 0 \Rightarrow u_i = y_i + \lambda$ if $y_i + \lambda < 0$

$\qquad\qquad\qquad\qquad\qquad\qquad y_i < -\lambda$.

if $u_i = 0$

$\Rightarrow u_i - y_i - \lambda = 0 \Rightarrow u_i = y_i + \lambda$ if $y_i + \lambda < 0$

$\qquad\qquad\qquad\qquad\qquad\qquad y_i < \lambda$

if $u$ $\Rightarrow y_i \in [-\lambda, \lambda]$

$$\Rightarrow \mu_i = \begin{cases} y_i - \lambda & , \text{if } y_i > \lambda \\ 0 & \text{if } -\lambda \le y_i \le \lambda \\ y_i + \lambda & \text{if } y_i < -\lambda \end{cases}$$

$$\Rightarrow \mu_i^{soft} = sgn(y_i)\, (|y_i| - \lambda)_+$$

$$\oint \|\hat{\mu}^{soft}(y) - \mu\|^2 \qquad \hat{\mu}^{soft} = y + g(y)$$

$$\lambda = \sqrt{2 \log p}$$

soft thresolding.

$$g(y) = \begin{cases} -\lambda & \text{if } y_i > \lambda \\ -y & \text{if } -\lambda \le y_i \le \lambda \\ \lambda & \text{if } y_i < -\lambda \end{cases}$$

from SURE Lemma 2.2 of the notes →

$$R(\hat{\mu}^{soft}, \mu) = \mathbb{E}_\mu \left( p + 2 \nabla^T g(y) + \| g(y) \|^2 \right)$$

$$\nabla^T g(y) = \sum_{i=1}^{p} \frac{\partial}{\partial y_i} g(y)$$

$$\frac{\partial g}{\partial y_i}(y) = -I\{ |y_i| \le \lambda \}$$

$$\Rightarrow R(\hat{\mu}_{soft}, \mu) = \mathbb{E}_\mu \left( p - 2 \sum_{i=1}^{p} I\{|y_i| \le \lambda\} \right. $$
$$\left. + \sum_{i=1}^{p} \min(\lambda^2, y^4) \right)$$

minimizing the RHS w.r.t $\lambda$ to get $\hat{\lambda}$ SURE

for univariate case, $y = \mu + z \sim N(\mu, 1)$ from [1]

$$\mu(\lambda, \mu) \le \mu(\lambda, 0) + \min(\mu^2, 1 + \lambda^2)$$

for $\lambda = \sqrt{2 \log p}$

$$\mu(\lambda, \mu) \le \frac{1}{p} + (2 \log p + 1) \min(\mu^2, 1)$$

for $p$-variate distribution, summing over the element

using lemma 2.11 of [1]:

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

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

If $\mu$ is sparse then $R(\hat{\mu}^{soft}, \mu) < R(\hat{\mu}^{MLE}, \mu)$

(c) $$\min_{\mu} \|y - \mu\|_2^2 + \lambda^2 \|\mu\|_0 = \min_{\mu} \sum_{i=1}^{n}\left((y_i - \mu_i)^2 + \lambda^2 \mathbb{1}\{\mu_i \neq 0\}\right)$$

Solving componentwise:

$$\min_{\mu_i} (y_i - \mu_i)^2 + \lambda^2 \mathbb{1}(\mu_i \neq 0)$$

if $\mu_i = 0 \Rightarrow$ cost $= y_i^2$

if $y_i^2 \leq \lambda^2$ then set $\mu_i = 0$

if $\mu_i \neq 0 \Rightarrow$ min. cost is $\lambda^2$ when $\mu_i = y_i$

$$\Rightarrow \mu_i = \begin{cases} y_i & \text{if } y_i^2 \geq \lambda^2 \\ 0 & \text{if } y_i^2 \leq \lambda^2 \end{cases}$$

$$\Rightarrow \mu_i = \begin{cases} y_i, & \text{if } y_i > \lambda \\ 0, & \text{if } -\lambda \leq y_i \leq \lambda \\ y_i, & \text{if } y_i < -\lambda \end{cases}$$

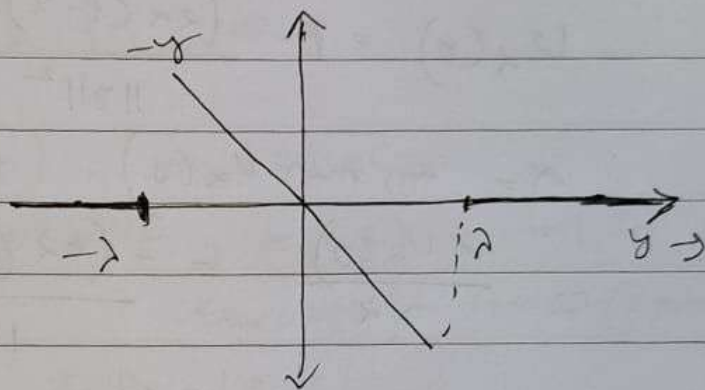$$\Rightarrow \hat{\mu_i}^{hard} = y_i \mathbb{1}(|y_i| > \lambda)$$

$\vec{u}^{hard} = y + g(y)$ .

$$\Rightarrow g(y) = \begin{cases} 0 & \text{if } y >> \\ -y & \text{if } -\lambda \le y \le \lambda \\ 0 & \text{if } y <-> . \end{cases}$$

$g(y) =$



$g(y)$ is not
weakly differentiable
due to sudden
jumps at $-\lambda$ & $\lambda$.

(d) $\quad u^{-1 js}(y) = \left(1 - \dfrac{\alpha}{\|y\|^2}\right) y$

$\mathbb{E} \; \|\hat{u}^{js}(y) - u\|^2 = \mathbb{E} \left\| y - \dfrac{\alpha y}{\|y\|^2} - u \right\|^2$

$\qquad = \mathbb{E} \left( \|y - u\|^2 + 2(y-u)^T \dfrac{\alpha y}{\|y\|^2} + \dfrac{\alpha^2 \|y\|^2}{\|y\|^4} \right)$

$\qquad = p + 2\mathbb{E}_u \left[ \dfrac{\alpha (y-u)^T y}{\|y\|^2} \right] + \mathbb{E}\left[ \dfrac{\alpha^2}{\|y\|^2} \right]$

$\qquad = p + \mathbb{E}_u \left[ \dfrac{2 \alpha \|y\|^2 - 2 \alpha u^T y + \alpha^2}{\|y\|^2} \right]$

$\qquad = \mathbb{E}_y \; 2\alpha \left( \|y\|^2 - u^T y \right)$

$\qquad = \mathbb{E}_y \; 2\alpha \left( y^T y - u^T y \right) = 2\alpha \left( (y-u)^T y \right)$

$$= \mathbb{E}_{y} \, 2\alpha \left( ||y - \mu||^2 - ||\mu||^2 - \mu^{\top} y \right) \big]$$

$$= 2\alpha (p-2)$$

$$\Rightarrow \mathbb{E} \, || \hat{\mu}^{JS}(y) - \mu ||^2 = \mathbb{E} \left[ p - \frac{(2\alpha(p-2) - \alpha^2)}{||y||^2} \right]$$

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

$$\alpha = \text{arg min } U_{\alpha}(y)$$

$$\frac{\partial U_{\alpha}(y)}{\partial \alpha} = - \frac{(2(p-2) - 2\alpha)}{||y||^2} = 0$$

$$\Rightarrow (p-2) - \alpha = 0 \Rightarrow \alpha^{*} = p-2$$

$$R_{NLE}^{(\hat{\mu}_{MLE})} = \mathbb{E} \, || \hat{\mu}_{MLE} - \mu ||^2$$

$$= \mathbb{E} \, || y - \mu ||^2 = p$$

$$R_{JS} = p - \mathbb{E} \, \frac{\left(2\alpha(p-2) - \alpha^2\right)}{||y||^2}$$

$$\text{for } p > 2 \Rightarrow \mathbb{E} \left( \frac{2\alpha(p-2) - \alpha^2}{||y||^2} \right) > 0$$

$$\Rightarrow R_{JS} < p = R_{MLE}$$

$$\Rightarrow R_{JS} < R_{MLE}.$$

(c)    $\alpha-1$ norm minimization, i.e. soft thresholding

& James Stein are shrinkage rule which satisfy

(a)    $\theta_\lambda((t)) \leq |t|$

(b)    $\theta_\lambda(-t) = \theta_\lambda(t)$

(c)    $\theta_\lambda(t) \leq \theta_\lambda(t')$ for $t \leq t'$

(d)    $\lim_{t \to \infty} \theta_\lambda(t) = \infty$

3) 3)    $y \sim N(\mu, \sigma^2 I_p)$, $\hat{\mu}_c(y) = cy$

(i)    let $|A| = (A^T A)^{1/2}$, tr $(A) \leq $ tr $|A|$

equality if $A = A^T$ (symmetric)

let $D$ be such that $I - D = |I - c|$

$\Rightarrow$ $D$ is symmetric

NSE, $\Rightarrow$ $\mathbb{E}[||\hat{\mu} - \mu||^2] = \mathbb{E}||\hat{\mu} - \mathbb{E}\hat{\mu}||^2 +$

$$||\mathbb{E}\hat{\mu} - \mu||^2 = \text{var}(\hat{\mu}) + \text{bias}^2(\hat{\mu}).$$

for linear estimators, $\text{Var}(\hat{\mu}) = \text{tr}(\text{Cov}(\hat{\mu})) =$

$$\text{tr}(\sigma^2 cc^T)$$

$\Rightarrow$ $\text{Var}(\hat{\mu}) = \sigma^2 \text{tr}(c^T c)$

$\text{Bias} = \mathbb{E}\hat{\mu} - \mu = |c - I| \mu$

$\Rightarrow$ $MSE = \sigma^2 \text{tr}(c^T c) + ||(I - c)\mu||^2$ ——— Ⓐ

Claim    MSE of $\hat{\mu}_D$ is everywhere better than

$\mu_c$ if $c$ is not symmetric.

$$(I - D)^T (I - D) = |I - c|^2 = (I - c)^T (I - c)$$

$\Rightarrow$ the bias squared is same for both estimators (from Ⓐ)

Now, for the variance term is

$$tr\left(D^T D\right) = tr[\cdot\; tr\; I - 2\, tr\left(I - D\right) + tr\left(I - D\right)^T \left(I - D\right)$$

Now, $tr\left(D^T D\right) < tr\left(C^T C\right)$ iff

$$tr\left(I - D\right) = tr\; tr\, \left| I - c \right| > tr\left(I - C\right)$$

$\Rightarrow$ It occurs only if $C$ is not symmetric

ii)  EVD of symmetric $C$ (proved in (i))

$\Rightarrow\quad C = U \Lambda U^{-1}$   All the eigenvalues are real

let $\eta = U^T \mu$ & $x = U^T y \sim N\left(\eta,\, \sigma^2 I_p\right)$ since

$$U^T U = I$$

Now, $\mathcal{E} \left\| (y - \mu) \right\|^2 = \mathcal{E} \left\| U \Lambda U^T y - \mu \right\|^2$

$$= \mathcal{E} \left\| \Lambda x - \eta \right\|^2$$

$\Rightarrow M\left(\hat{\mu},\, \mu\right) = M\left(\hat{\eta}_\Lambda,\, \eta\right) = \sum_{i=1}^{P} \sigma^2 \lambda_i +$

$$\left(1 - \lambda_i\right)^2 \eta_i$$

If $\lambda_i \in [0, 1]$, a strictly better MSE can be

obtained by replacing $\lambda_i$ by

1 If $\lambda_i > 1$  for 0  If $\lambda_i < 0$.

iii) Let $\lambda_1 = \lambda_2 = \cdots \cdots = \lambda_d = 1 > \lambda_i$, $i > d \geq 3$ &

let $x^d = (x_1, \cdots x_d)$

Positive part of JS estimator is everywhere better than

$$\hat{\eta}_1 (x^d) = x^d$$

If a new estimator $\hat{\eta}$ is defined to use $\hat{\eta}^{JS}$ on $x^d$
& to continue to use $\lambda_i x_i$ for $i > d$ then

$$r(\hat{\eta}, \eta) = r(\hat{\eta}^{JS}, \eta^d) + \sum_{i > d} r(\lambda_i, \eta_i) \quad \omega(\eta, \eta)$$

$\Rightarrow$ So, $\hat{\eta}$ dominates $\hat{\eta}_\lambda$ & hence $\hat{\mu}_c$.

8) 4)     for $P = 1$     $\gamma \sim N(\mu, \sigma^2)$

$$R(\hat{\mu}^{JS}, \mu) = P - E\mu \frac{(P-2)^2}{\|\gamma\|^2} \qquad\qquad —①$$

$\|\gamma\|^2 \Rightarrow$ follows non central Chi Squared
distribution with non - centality param $\|\mu\|^2$

Non-central distn can be realized as a mixture of
central Chi Squared distn $x^2_{p+2N}$, where $N$ is a
Poisson variable with mean $\|\mu\|^2/2$

& $E\left[1/x^2_p\right] = \frac{1}{p-2} \qquad — Ⓐ$

for $P = 1$ & $2$

for $P = 1 \Rightarrow R(\hat{\mu}_{JS}, \mu) = 1 - \frac{1}{(-1)} = 2 > 1 = MLE$

for $P = 2$, $R(\hat{\mu}_{JS}, \mu) = 2 = R(\hat{\mu}_{MLE}, \mu)$

Now, conditioning on $N$ & using $A$

$$\mathbb{E}\left(\frac{1}{\|\gamma\|^2}\right) = \mathbb{E}\left(\frac{1}{x^2_{p+2N}}\right) = \mathbb{E}\left[\frac{1}{p+2N-2}\right]$$

$$\geq \frac{1}{p-2+\|\mu\|^2}$$

Substituting in ①

$$R\left(\hat{\mu}^{JS}, \mu\right) - p - \mathbb{E}_\mu\left(\frac{(p-2)^2}{\|\gamma\|^2}\right) \leq p - \frac{(p-2)^2}{p-2+\|\mu\|^2}$$

$$= 2 + (p-2) - \frac{(p-2)^2}{(p-2)+\|\mu\|^2}$$

$$= 2 + \frac{(p-2)^2\|\mu\|^2}{(p-2)+\|\mu\|^2}.$$