

Homework 3

1. (a) The joint probability density function is

$$f_X(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

The likelihood function

$$\begin{aligned} L(\mu, \Sigma) &= \prod_{j=1}^n f_X(x_j, \mu, \Sigma) \\ &= \prod_{j=1}^n \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)\right) \\ &= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} (x_j - \mu)^T \Sigma^{-1} (x_j - \mu)\right) \end{aligned}$$

$$\ln(L(\mu, \Sigma))$$

$$\begin{aligned} &= -\frac{np}{2} \ln(2\pi) - \frac{n}{2} \ln \det(\Sigma) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \\ &= -\frac{n}{2} \text{trace}(\Sigma^{-1} S_n) - \frac{n}{2} \ln \det(\Sigma) + C \end{aligned}$$

where $S_n = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$ and C does not depend on μ and Σ .

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

$$\frac{df(x)}{dx} = \lim_{\Delta \rightarrow 0} \frac{f(x+\Delta) - f(x)}{\Delta} = -X^T A X^{-1}$$

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

where λ_i is the i th eigenvalue of $X^{-\frac{1}{2}}\Delta X^{-\frac{1}{2}}$. Since Δ is small, which implies λ_i is small. $\log(1 + \lambda_i) \approx \lambda_i$.

$$\begin{aligned} \log \det(X + \Delta) &\approx \log \det X + \sum_{i=1}^n \lambda_i \\ &= \log \det X + \text{trace}(X^{-\frac{1}{2}}\Delta X^{-\frac{1}{2}}) \\ &= \log \det X + \text{trace}(X^{-1}\Delta) \\ \frac{dg(x)}{dx} &= \lim_{\Delta \rightarrow 0} \frac{g(X + \Delta) - g(X)}{\Delta} = X^{-1} \end{aligned}$$

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

$$\begin{aligned} \frac{\partial \ln}{\partial n} &= \nabla n \left(-\frac{1}{2} \sum_{i=1}^n (x_i - m) \Sigma^{-1} (x_i - m) \right) \\ &= \sum_{i=1}^n \Sigma^{-1} (x_i - m) \\ &= \Sigma^{-1} \sum_{i=1}^n (x_i - m) \end{aligned}$$

which is equal to 0 only if

$$\sum_{i=1}^n x_i - nm = 0 \Rightarrow m = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\begin{aligned}
\frac{\partial \ln}{\partial \Sigma} &= \nabla_{\Sigma} \left[-\frac{n}{2} \text{trace}(\Sigma^{-1} S_n) - \frac{n}{2} \log \det(\Sigma) + C \right] \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \nabla_{\Sigma} \left(\sum_{i=1}^n \text{trace}(\Sigma^{-1} S_n) \right) \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \nabla_{\Sigma} \left(\text{trace} \left(\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \right) \Sigma^{-1} \right) \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \left[\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T \right]^T \\
&= \frac{n}{2} \Sigma^T - \frac{S_n}{2}
\end{aligned}$$

Set the equation to 0, we get

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

2. (a) let $L(\mu) = \frac{1}{2} \|y - \mu\|_2^2 + \frac{\lambda}{2} \|\mu\|_2^2$

$$\frac{\partial L}{\partial \mu} = (\mu - y) + \lambda \mu$$

Set the equation to 0, we get

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

$$\text{MSE}(\lambda) = E\{(W\lambda\hat{\beta} - \beta)^T (W\lambda\hat{\beta} - \beta)\}$$

$$= r^2 \text{trace}\{W\lambda(X^T X)^{-1} W\lambda^T\} + \beta^T (W\lambda - I)^T (W\lambda - I) \beta$$

$$\text{Since } X^T X = I = (X^T X)^T$$

$$\text{MSE}(\hat{\mu}_i^{\text{ridge}}) = \frac{\beta}{(1+\lambda)^2} + \frac{\lambda^2}{(1+\lambda)^2} y_i^T y_i$$

$$(b) \quad \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 + \lambda \|\mathbf{u}\|_1$$

Expanding out the first term we get $\frac{1}{2} \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{u} + \frac{1}{2} \mathbf{u}^T \mathbf{u}$
 and since $\mathbf{u}^T \mathbf{u}$ does not contain \mathbf{y} , we can consider another equivalent problem

$$\min_{\mathbf{u}} (-\mathbf{y}^T \mathbf{u} + \frac{1}{2} \|\mathbf{u}\|^2) + \lambda \|\mathbf{u}\|_1.$$

Noting that $\hat{\mathbf{u}}^{LS} = \mathbf{y}$, the previous problem can be written as

$$\min_{\mathbf{u}} \sum_{i=1}^p -\hat{u}_i^{LS} u_i + \frac{1}{2} u_i^2 + \lambda |u_i|.$$

Fix a certain i . Then, we want to minimize

$$L_i = -\hat{u}_i^{LS} u_i + \frac{1}{2} u_i^2 + \lambda |u_i|.$$

If $\hat{u}_i^{LS} > 0$, then we must have $u_i \geq 0$ since otherwise we could flip its sign and get a lower value for the objective function. Likewise if $\hat{u}_i^{LS} < 0$, then we must choose $u_i \leq 0$.

(Case 1). $\hat{u}_i^{LS} > 0$. Since $u_i \geq 0$,

$$L_i = -\hat{u}_i^{LS} u_i + \frac{1}{2} u_i^2 + \lambda u_i,$$

and differentiating this with respect to u_i and setting equal to 0, we get $u_i = \hat{u}_i^{LS} - \lambda$ and this is only feasible if the right-hand side is nonnegative, the actual solution is

$$u_i^{(1)} = \text{sign}(y_i) (|y_i| - \lambda)^+.$$

(Case 2). $\hat{u}_i^{LS} \leq 0$. This implies we must have $u_i \leq 0$ and so

$$L_i = -\hat{u}_i^{LS} u_i + \frac{1}{2} u_i^2 - \lambda u_i.$$

Differentiating with respect to u_i and setting equal to 0,

$u_i = \hat{u}_i^{LS} + \lambda = \text{sign}(y_i) (|y_i| - \lambda)$. To ensure this is feasible, we need $u_i \leq 0$, which is achieved by taking

$$\hat{u}_i^{(asso)} = \text{sign}(y_i)(|y_i| - \lambda)^+$$

Thus $\hat{u}_i^{\text{soft}} = u_{\text{soft}}(y_i; \lambda) = \text{sign}(y_i)(|y_i| - \lambda)^+$.

The risk function

$$r_s(\lambda, n) = \int [u_n(m+2) - n]^2 \phi(z) dz.$$

Differentiation under the expectation yields

$$0 \leq \frac{\partial}{\partial n} r_s(\lambda, n) = 2n P(|m+2| \leq \lambda) \leq 2n$$

The risk at $n=0$ has a simple form

$$r_s(\lambda, 0) = 2 \int_{-\infty}^{\infty} (z - \lambda)^2 \phi(z) dz$$

$$= 2(\lambda^2 + 1) \bar{\Phi}(\lambda) - 2\lambda \phi(\lambda)$$

and, using the bound for Mills ratio $\bar{\Phi}(u) \leq \lambda^{-1} \phi(u) \leq e^{-\frac{\lambda^2}{2}}$

$$r_s(\lambda, 0) \leq 2\lambda^{-1} \phi(\lambda) \leq e^{-\frac{\lambda^2}{2}}$$

We have $r_s(\lambda, n) - r_s(\lambda, 0) \leq n^2$. Using the bound at ∞ , we get

$$r_s(\lambda, n) \leq r_s(\lambda, 0) + \min(n^2, 1 + \lambda^2).$$

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

$$r_s(\lambda, n) \leq 1 + (2 \log p + 1) \sum_{i=1}^n \min(u_i^2, 1).$$

(c) The problem can be written as

$$\min_m \|y - m\|_2^2 + \lambda^2 \|u\|_0 = \min_m \frac{1}{2} \sum_{i=1}^n (u_i - y_i)^2 + \lambda \mathbb{1}_{\{u_i \neq 0\}}$$

Since it has component wise property, we can solve the scalar problem,

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

In case $u_i = 0$ the cost is $\frac{1}{2} y_i^2$. Hence as long $\frac{1}{2} y_i^2 \leq \lambda$ it is worth to set $u_i = 0$. In the other case the minimum value

of the cost function is λ where $u_i = y_i$.

Hence

$$\hat{u}_i^{\text{hard}} = u_{\text{hard}}(y_i, \lambda) = y_i I(|y_i| > \lambda).$$

$$\hat{u}_{\text{hard}}(y) = \begin{cases} 0 & \text{if } |y| \leq \lambda \\ y & \text{if } |y| > \lambda \end{cases}$$

$$g(y) = \begin{cases} 1 & \text{if } |y| \leq \lambda \\ 0 & \text{if } |y| > \lambda \end{cases}$$

$g(y)$ is weakly differentiable, since $g(y)$ is not differentiable but integrable.

$$(d) \hat{u}^{\text{JS}}(y) = \left(1 - \frac{\alpha}{\|y\|^2}\right)y$$

$g(y) = -\alpha \|y\|^{-2} y$ is weakly differentiable,

$$Dg_i(y) = -\alpha \left(\frac{1}{\|y\|^2} - \frac{2y_i^2}{\|y\|^4}\right)$$

$\nabla^T g(y) = -\alpha^2 \|y\|^{-2}$ and the unbiased risk estimator

$$V(x) = n - \alpha^2 \|x\|^2$$

Consequently, $V(\hat{u}^{\text{JS}}, n) = n - (n-2)^2 E \|x\|^2$

The estimator $\hat{u}(y) = \left(1 - \frac{\beta}{\|x\|^2}\right)x$ has unbiased risk

estimate $V_\beta(x) = n - \frac{2\beta(n-2) - \beta^2}{\|x\|^2}$, the quantity is minimized for each x by the choice $\beta = n-2$.

Thus $E\|\hat{u}^{\text{JS}}(y) - u\|^2 = E V_\alpha(y)$

where $V_\alpha(y) = p - \frac{2\alpha(p-2) - \alpha^2}{\|y\|^2}$.

$$\alpha^* = 2$$

We begin with the identity

$$\sum_{i=1}^N (\hat{m}_i - m_i)^2 \leq \sum_{i=1}^N [(z_i - \hat{m}_i)^2 - (z_i - m_i)^2 + 2(\hat{m}_i - m_i)(z_i - m_i)]$$

Taking expectations on both sides

$$E(\|\hat{m} - m\|^2) = E(\|z - \hat{m}\|^2) - N + 2 \sum_{i=1}^N \text{cov}(\hat{m}_i, z_i)$$

Using integral by part, we can get

$$\text{cov}(\hat{m}_i, z_i) = E\left(\frac{\partial \hat{m}_i}{\partial z_i}\right)$$

$$E(\|\hat{m} - m\|^2) = E(\|z - \hat{m}\|^2) - N + 2 \sum_{i=1}^N E\left[\frac{\partial \hat{m}_i}{\partial z_i}\right]$$

Plug in $\hat{m}^{(MLE)} = z$

$$E(\|z - \hat{m}^{(MLE)} - m\|^2) = E(\|z - z\|^2) - N + 2N = N$$

Plug in $\hat{m}^{(JS)} = z - \frac{n-2}{\|z\|^2} z$

$$\frac{\partial \hat{m}_i^{(JS)}}{\partial z_i} = 1 - \frac{n-2}{\|z\|^2} + \frac{2(n-2)z_i}{\|z\|^4}$$

$$\begin{aligned} \sum_{i=1}^N E\left(\frac{\partial \hat{m}_i^{(JS)}}{\partial z_i}\right) &= E\left[\sum_{i=1}^N \frac{\partial \hat{m}_i^{(JS)}}{\partial z_i}\right] \\ &= n - E\left(\frac{(n-2)z}{\|z\|^2}\right). \end{aligned}$$

Note that $E(\|z - \hat{m}^{(JS)}\|^2) = E\left(\left\|\frac{n-2}{\|z\|^2} z\right\|^2\right)$

$$\begin{aligned} E(\|\hat{m}^{(JS)} - m\|^2) &= E(\|z - \hat{m}^{(JS)}\|^2) - N + 2 \sum_{i=1}^N E\left[\frac{\partial \hat{m}_i^{(JS)}}{\partial z_i}\right] \\ &= n - E\left(\frac{(n-2)z}{\|z\|^2}\right) \end{aligned}$$

Therefore, $E(\|\hat{m}^{(JS)} - m\|^2) < E(\|\hat{m}^{(MLE)} - m\|^2)$ where $n > 2$.

(e) Ridge regression, LASSO regression and the James-Stein Estimator are shrinkage rules.

3. (a) We use the notation $|A| = (\text{tr} A^T A)^{\frac{1}{2}}$ and the fact that $\text{tr} A \leq \text{tr} |A|$, with equality only if A is symmetric, $A^T = A$. Let D be defined via the identity $I - D = (I - C)$. D is symmetric and we use the variance-bias decomposition to show that the MSE of $\hat{\theta}_0$ is everywhere better than that of $\hat{\theta}_C$ if C is not symmetric. Since

$$(I - D)^T (I - D) = (I - C)^T (I - C)$$

the two estimators have the same bias. Turning to the variance terms, write

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

Comparing with the corresponding variance term for $\hat{\theta}_C$, we see that $\text{tr } D^T D < \text{tr } C^T C$ iff

$$\text{tr} (I - D) = \text{tr} (I - C) > \text{tr} (I - C)$$

which occurs iff C fails to be symmetric.

- (b) We assume C is symmetric, we can find a decomposition $C = V R V^T$ with V orthogonal and $R = \text{diag}(\xi_i)$ containing the eigenvalues of C . Now change variables to $y = V^T \theta$ and $x = V^T y \sim N(\eta, \epsilon^2 I)$. Orthogonality of V implies that $E \|(\eta - \theta)\|^2 = E \|Rx - \eta\|^2$, so we have

$$r(\hat{\theta}_C, \theta) = r(\hat{\eta}_R, y) = \sum_i \epsilon^2 \xi_i^2 + (1 - \xi_i^2) \eta_i^2$$

$$= \sum_i r_C(\xi_i, \eta_i)$$

if any eigenvalue $\xi_i \notin (0, 1)$, a strictly better MSE results by replacing ξ_i by 1 if $\xi_i > 1$ and by 0 if $\xi_i < 0$.

(c) Suppose that $b_1 = \dots = b_d = 1 > \beta_i$ for $i > d > 3$, and let $x^d = (x_1, \dots, x_d)$. We have note that the James-Stein estimator is everywhere better than $\hat{y}_J(x^u) = x^d$. We define a new estimator \hat{y} to use \hat{y}^{JS} on x^d and to continue to use $\beta_i x_i$ for $i > d$, then

$$r(\hat{y}, y) = r(\hat{y}^{JS}, y^d) + \sum_{i>d} r_i(\beta_i x_i, y_i) < r(\hat{y}_J, y),$$

and so \hat{y} dominates \hat{y}_J and hence $\hat{\theta}_C$.

4. For $p=1$ case, the J-S estimator has higher risk than the MLE for any value of n . This is known as the Stein's paradox, where the J-S estimator shrinks the estimates towards the origin, leading to a higher risk than the MLE.
 For $p=2$ case, the J-S estimator has a lower risk than the MLE for n in a certain range.

let $X \sim N(n, I)$. If $p \geq 3$, then

$$\xi \frac{1}{\|X\|^2} \leq \frac{1}{p-2} \left(\frac{p}{p+||n||^2} \right)$$

$$\frac{R(\hat{y}^{JS}, n)}{R(n, X)} = 1 - \frac{(p-2)^2}{p} \xi \frac{1}{\|X\|^2}$$

$$\leq 1 - \frac{(p-2)^2}{p} \left(\frac{1}{p-2+||n||^2} \right)$$

$$\therefore R(\hat{y}^{JS}, n) \leq p - \frac{(p-2)^2}{p-2+||n||^2} \text{ since } R(n, X) \geq p$$

$$5. (a) \quad \begin{aligned} u_i &\sim N(\mu, A) \quad (i=1, \dots, n) & p(\theta) &= N(\mu, A) \\ x_i &\sim N(u_i, 1) & p(x|\theta) &= N(u_i, 1) \end{aligned}$$

$$\begin{aligned} \ln p(x, \theta) &= \ln p(x|\theta) + \ln p(\theta) \\ &= \ln N(x_i | u_i, 1) + \ln N(u_i | \mu, A) \\ &= -\frac{1}{2} \ln 2\pi - \frac{(x_i - u_i)^2}{2} - \frac{1}{2} \ln 2\pi - \frac{1}{2} \ln A - \frac{u_i^2}{2A} \end{aligned}$$

Consider the second-order term,

$$\begin{aligned} &- \frac{(x_i - u_i)^2}{2} - \frac{u_i^2}{2A} \\ &= -\frac{1}{2} [x_i^2 - 2x_i u_i + (1 + \frac{1}{A}) u_i^2] \\ &= -\frac{1}{2} \begin{pmatrix} x_i \\ u_i \end{pmatrix}^T \begin{pmatrix} 1 & -1 \\ -1 & (1 + \frac{1}{A}) \end{pmatrix} \begin{pmatrix} x_i \\ u_i \end{pmatrix} \end{aligned}$$

We get

$$\text{cov}\left[\begin{pmatrix} x_i \\ u_i \end{pmatrix}\right] = \begin{pmatrix} 1 & -1 \\ -1 & (1 + \frac{1}{A}) \end{pmatrix}^{-1} = \begin{pmatrix} A+1 & A \\ A & A \end{pmatrix}$$

$$\text{we have } E\left[\begin{pmatrix} x_i \\ u_i \end{pmatrix}\right] = \begin{pmatrix} M \\ M \end{pmatrix}$$

$$p(x) = N(M, A+1)$$

Now we consider the conditional distribution $p(\theta|x)$. Pick out all second-order terms of u_i , where x_i is regarded as a constant, we have

$$-\frac{1}{2} \left(1 + \frac{1}{A}\right) u_i^2,$$

We can get the variance of $\theta|x$

$$\text{var}(\theta|x) = \frac{A}{1+A}$$

$$\text{we have } E(\theta|x) = \frac{A}{1+A} x + \frac{M}{1+A}$$

The conditional distribution of $\theta|x$ is given by

$$\theta|x \sim N\left(\frac{A}{1+A}x + \frac{M}{1+A}, \frac{A}{1+A}\right)$$

We have $Q = \sum_{i=1}^N \frac{x_i}{1+A} \sim X_N^2$

Note that $Q \sim X_N^2$, $\frac{1}{Q}$ follows Inverse- X^2 with $df = n$ and $E\left(\frac{1}{Q}\right) = \frac{1}{n-3}$

Therefore, $E\left(\frac{1}{\sum_{i=1}^N \frac{x_i^2}{1+A}}\right) = \frac{1}{n-3}$, and we use

$$\hat{M} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{B} = 1 - \frac{n-3}{S} \text{ with } S = \sum_{i=1}^n x_i^2$$

This leads to the James-Stein estimator

$$\hat{m}_i^{JS} = \bar{x} + \left(1 - \frac{n-3}{S}\right)(x_i - \bar{x}).$$

$$E\|\hat{m}^{Bayes} - m\|^2 = E[\|M + B(x_i - M) - m\|^2]$$

$$= E[\|m - m\|^2] + nB = nB$$

$$E\|\hat{m}^{MLE} - m\|^2 = E(\|m - m\|^2) - n + 2n = n$$

$$\hat{m}^{MLE} = \bar{x}$$

(b) We assume $m \sim g(\cdot)$, where $g(\cdot)$ is an arbitrary PDF.
 $m \sim g(\cdot)$, $z | m \sim N(m, 1)$.

The marginal distribution of z is

$$f(z) = \int_{-\infty}^{\infty} \varphi(z - m) g(m) dm$$

where $\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$

The posterior distribution of m is

$$g(m|z) = \varphi(z - m) g(m) / f(z).$$

Take derivative on both sides

$$-\frac{d\psi(y)}{dy} \exp(-\psi(y)) \int \exp(yx) h_0(x) dx + \exp(-\psi(y)) \int \exp(yx) h_0(x) x dx = 0$$
$$-\frac{d\psi(y)}{dy} + \exp(-\psi(y)) \int \exp(yx) h_0(x) x dx = 0$$

and $\exp(-\psi(y)) \int \exp(yx) h_0(x) x dx = \int x h(x) dx = E(x)$

$$\frac{d\psi(y)}{dy} = E(x)$$

$$g(u|z) = \psi(z-u) g(u) / f(z)$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z-u)^2}{2}\right) g(u) / f(z)$$

$$= [\exp(zu)] \left[\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) / f(z) \right] [\exp\left(-\frac{u^2}{2}\right) g(u)]$$

$$= \exp(zu - \psi(z)) h_0(u).$$

where $\psi(z) = \log \frac{f(z)}{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)}$ is the cumulant

Generating function, and $h_0(u) = \exp\left(-\frac{u^2}{2}\right) g(u)$.

$$E(\ln(z)) = z + \frac{d}{dz} \log f(z).$$