

Problem 1

Let $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$, then the mean value of this random vector is $\bar{\mathbf{X}} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n)^T$.
The covariance matrix of \mathbf{X} is calculated by:

$$\mathbf{C}_X = \left\langle (\mathbf{X} - \bar{\mathbf{X}})(\mathbf{X} - \bar{\mathbf{X}})^T \right\rangle$$

Let \mathbf{M} be the constant matrix, then :

$$\mathbf{Y} = \mathbf{M} \cdot \mathbf{X}$$

Then, the covariance of \mathbf{Y} is calculated by:

$$\begin{aligned} \mathbf{C}_Y &= \left\langle (\mathbf{Y} - \bar{\mathbf{Y}})(\mathbf{Y} - \bar{\mathbf{Y}})^T \right\rangle \\ &= \left\langle (\mathbf{M}\mathbf{X} - \mathbf{M}\bar{\mathbf{X}})(\mathbf{M}\mathbf{X} - \mathbf{M}\bar{\mathbf{X}})^T \right\rangle \\ &= \left\langle \mathbf{M}(\mathbf{X} - \bar{\mathbf{X}})(\mathbf{X} - \bar{\mathbf{X}})^T \mathbf{M}^T \right\rangle \\ &= \mathbf{M} \left\langle (\mathbf{X} - \bar{\mathbf{X}})(\mathbf{X} - \bar{\mathbf{X}})^T \right\rangle \mathbf{M}^T \\ &= \mathbf{M}\mathbf{C}_X\mathbf{M}^T \end{aligned}$$

Problem 2

Method 1

The logarithm of the PDF is :

$$\ln p = -\frac{n}{2} \ln |2\pi\mathbf{C}| - \frac{1}{2} \left[(\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true}))^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true})) \right]$$

Then, the first order derivative w.r.t θ is:

$$\begin{aligned} \frac{\partial \ln p}{\partial \theta_i} &= \frac{1}{2} \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right)^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true})) \right] \\ &\quad + \frac{1}{2} \left[(\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true}))^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right) \right] \end{aligned}$$

Then, the second derivative is :

$$\begin{aligned}\frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} = & \frac{1}{2} \left[\left(\frac{\partial^2 \boldsymbol{\mu}}{\partial \theta_i \partial \theta_j} \right)^T \cdot \mathbf{C}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true})) \right] \\ & - \frac{1}{2} \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_j} \right) \right] \\ & - \frac{1}{2} \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_j} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right) \right] \\ & + \frac{1}{2} \left[(\mathbf{x} - \boldsymbol{\mu}(\boldsymbol{\theta}_{true}))^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial^2 \boldsymbol{\mu}}{\partial \theta_i \partial \theta_j} \right) \right]\end{aligned}$$

Then, the $i, j - th$ element in Fisher Matrix is the expectation value of the above equation:

$$F_{ij} = - \left\langle \frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} \right\rangle = 0 - \frac{1}{2} \left\langle \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_j} \right) \right] \right\rangle - \frac{1}{2} \left\langle \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_j} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right) \right] \right\rangle + 0$$

Since this two terms are symmetric, we then have:

$$F_{ij} = - \left\langle \frac{\partial^2 \ln p}{\partial \theta_i \partial \theta_j} \right\rangle = - \left\langle \left[\left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_i} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}}{\partial \theta_j} \right) \right] \right\rangle$$

Method 2

Some properties that would be used in calculating the gradient of matrix:

$$\begin{aligned}x &= \text{tr}(x) \\ \text{tr}(c_1 \mathbf{A} + c_2 \mathbf{B}) &= c_1 \text{tr}(\mathbf{A}) + c_2 \text{tr}(\mathbf{B}) \\ \text{tr}(\mathbf{A}) &= \text{tr}(\mathbf{A}^T) \\ \langle \mathbf{A}, \mathbf{B} \rangle &= \text{tr}(\mathbf{A}^T \mathbf{B}) \\ \text{tr}(\mathbf{AB}) &= \text{tr}(\mathbf{BA}) \\ \text{tr}(\mathbf{ABC}) &= \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB}) \\ d(\mathbf{F}(\mathbf{X})\mathbf{G}(\mathbf{X})) &= d(\mathbf{F}(\mathbf{X}))\mathbf{G}(\mathbf{X}) + \mathbf{F}(\mathbf{X})d(\mathbf{G}(\mathbf{X})) \\ d\mathbf{F}_{p \times q}^T(\mathbf{X}) &= (d\mathbf{F}_{p \times q}(\mathbf{X}))^T\end{aligned}$$

And the most important one:

$$\text{tr}(d\mathbf{F}_{p \times q}) = \text{tr} \left(\left(\frac{d\mathbf{F}_{p \times q}}{d\mathbf{X}} \right)^T d\mathbf{X} \right)$$

So, the derivative of the matrix is then the transpose of the left part.

$$\begin{aligned}
\text{tr}(\text{d} \ln p) &= \text{tr} \left(\frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right)^T \cdot \mathbf{C}^{-1} \cdot [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})] + \frac{1}{2} [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})]^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right) \right) \\
&= \text{tr} \left(\frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right)^T \cdot \mathbf{C}^{-1} \cdot [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})] \right) + \text{tr} \left(\frac{1}{2} [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})]^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right) \right) \\
&= \frac{1}{2} \text{tr} \left([\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})]^T \cdot (\mathbf{C}^{-1})^T \cdot \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right) \right) + \frac{1}{2} \text{tr} \left([\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})]^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right) \right)
\end{aligned}$$

So, the gradient of $\ln p$ is:

$$\frac{\text{d} \ln p}{\text{d}\boldsymbol{\theta}} = \frac{1}{2} \left(\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot ((\mathbf{C}^{-1})^T + \mathbf{C}^{-1}) \cdot [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})] \right)$$

then, take the second derivative of this:

$$\begin{aligned}
\text{d} \left(\frac{\text{d} \ln p}{\text{d}\boldsymbol{\theta}} \right) &= \frac{1}{2} \text{d} \left(\left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \right) \cdot [(\mathbf{C}^{-1})^T + \mathbf{C}^{-1}] \cdot [\mathbf{X} - \boldsymbol{\mu}(\boldsymbol{\theta})] \\
&\quad - \frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot [(\mathbf{C}^{-1})^T + \mathbf{C}^{-1}] \cdot \text{d}\boldsymbol{\mu}(\boldsymbol{\theta})
\end{aligned}$$

Then take the expectation value of this equation, the first term is 0, so, we have:

$$\begin{aligned}
\left\langle \text{d} \left(\frac{\text{d} \ln p}{\text{d}\boldsymbol{\theta}} \right) \right\rangle &= \left\langle -\frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot [(\mathbf{C}^{-1})^T + \mathbf{C}^{-1}] \cdot \text{d}\boldsymbol{\mu}(\boldsymbol{\theta}) \right\rangle \\
&= \left\langle -\frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot [(\mathbf{C}^{-1})^T + \mathbf{C}^{-1}] \cdot \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \text{d}\boldsymbol{\theta} \right\rangle
\end{aligned}$$

Therefore:

$$\left\langle \text{d} \left(\frac{\text{d}^2 \ln p}{\text{d}\boldsymbol{\theta}^2} \right) \right\rangle = \left\langle -\frac{1}{2} \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot [(\mathbf{C}^{-1})^T + \mathbf{C}^{-1}] \cdot \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\rangle$$

Since the covariance matrix is symmetric, the final result is then:

$$\left\langle \frac{\text{d}^2 \ln p}{\text{d}\boldsymbol{\theta}^2} \right\rangle = - \left\langle \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^T \cdot \mathbf{C}^{-1} \cdot \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\rangle$$

Then, we have:

$$\mathbf{F}_{ij} = \left(\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_i} \right)^T \cdot \mathbf{C}^{-1} \cdot \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})}{\partial \theta_j}$$

When $n = 1 = m$,

$$F = \left(\frac{\partial \mu}{\partial \theta} \right)^2 \frac{1}{\sigma^2}$$

So, higher value of F means that it contains more information. There are two explanations:

1. when the derivative is large, that means μ is sensitive to θ , a small deviation from the true value could lead to a big difference.

2. when the σ variance is small, that means the distribution is more localized around the mean value.

Problem 3

Take the derivative of $\chi^2(a)$:

$$\begin{aligned}\chi^2(a) &= [\mathbf{x} - a\mathbf{v}]^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}] \\ d(\chi^2(a)) &= [-\mathbf{v}da]^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}] + [\mathbf{x} - a\mathbf{v}]^T \mathbf{C}^{-1} [-\mathbf{v}da] \\ &= \left\{ [-\mathbf{v}]^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}] da + [\mathbf{x} - a\mathbf{v}]^T \mathbf{C}^{-1} [-\mathbf{v}] da \right\} \\ &= -2 [\mathbf{v}^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}]] da \\ \frac{d(\chi^2(a))}{da} &= -2 [\mathbf{v}^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}]]\end{aligned}$$

Since the transpose of a scalar is still itself, we can combine the two term in the third line.

To get the a_{best} that makes the χ^2 minimized, make the derivative 0:

$$\begin{aligned}\frac{d(\chi^2(a))}{da} &= -2 [\mathbf{v}^T \mathbf{C}^{-1} [\mathbf{x} - a\mathbf{v}]] = 0 \\ a &= \frac{\mathbf{v}^T \mathbf{C}^{-1} \mathbf{x}}{\mathbf{v}^T \mathbf{C}^{-1} \mathbf{v}}\end{aligned}$$

For a random variable A_{best} , its variance is:

$$\begin{aligned}A_{best} &= \frac{\mathbf{X}^T \mathbf{C}^{-1} \mathbf{V}}{\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V}} \\ \text{Cov}(A_{best}) &= \text{Cov}\left(\frac{\mathbf{X}^T \mathbf{C}^{-1} \mathbf{V}}{\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V}}\right) \\ &= \frac{1}{(\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V})^2} \text{Cov}(\mathbf{X}^T \mathbf{C}^{-1} \mathbf{V}) \\ &= \frac{1}{(\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V})^2} \text{Cov}(\mathbf{V}^T \mathbf{C}^{-1} \mathbf{X}) \\ &= \frac{1}{(\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V})^2} (\mathbf{V}^T \mathbf{C}^{-1}) \text{Cov}(\mathbf{X}) (\mathbf{V}^T \mathbf{C}^{-1})^T \\ &= \frac{1}{(\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V})^2} (\mathbf{V}^T \mathbf{C}^{-1}) \mathbf{C} \mathbf{C}^{-1} \mathbf{V} \\ &= \frac{1}{\mathbf{V}^T \mathbf{C}^{-1} \mathbf{V}}\end{aligned}$$

So, the variance of A_{best} is then $\text{Var}(A_{best}) = \mathbf{V}^{-1} \mathbf{C} (\mathbf{V}^T)^{-1}$.

Fisher matrix for the gaussian distribution is :

$$\mathbf{F} = \left(\frac{\partial \boldsymbol{\mu}(\vec{\theta})}{\partial \vec{\theta}} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial \boldsymbol{\mu}(\vec{\theta})}{\partial \vec{\theta}} \right)$$

Here, the variable is a_{best} , so the Fisher matrix becomes:

$$\begin{aligned} \mathbf{F} &= \left(\frac{\partial a \mathbf{v}}{\partial a} \right)^T \cdot \mathbf{C}^{-1} \cdot \left(\frac{\partial a \mathbf{v}}{\partial a} \right) \\ &= \mathbf{v}^T \cdot \mathbf{C}^{-1} \cdot \mathbf{v} \end{aligned}$$

So, the inverse of Fisher matrix is :

$$\mathbf{F}^{-1} = \mathbf{V}^{-1} \cdot \mathbf{C} \cdot (\mathbf{V}^T)^{-1}$$

So, this is the same as the inverse of Fisher matrix.

Cramer-Rao-bound :

$$\text{var}(\hat{\theta}) \geq \frac{1}{F(\theta)}$$

The precision of any unbiased estimation is at most the Fisher information matrix.

The estimator is statistically efficient, making the best use of the available information in estimating the parameter of interest.