

## Lecture 4: Cramér-Rao Information Bound

Lecturer: Tony Sit

Scribe: Cheukhei Chan and Yang Boyu

**Disclaimer:** These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

## 4.1 UMVUE

Q: Suppose we have  $\delta_1$  and  $\delta_2$  as UMVUEs for  $g_1(\theta)$  and  $g_2(\theta)$  respectively. Is  $\delta_1 + \delta_2$  an UMVUE for  $g_1(\theta) + g_2(\theta)$ ?

**Theorem 4.1 (Characterization of UMVUEs; TPE 2.1.7)** Let  $\Delta = \{\delta : E(\delta^2) < \infty\}$ . Then  $\delta_0 \in \Delta$  is UMVU for  $g(\theta) = E(\delta_0)$  if and only if  $E(\delta_0(\theta)U) = 0$  for every  $U \in \mathcal{U} = \{U : E(U) = 0\}$

Proof:

If  $\delta_0$  is a UMVUE, let's consider  $\delta_\lambda = \delta_0 + \lambda U$  for  $\lambda \in \mathbb{R}$  and  $U \in \mathcal{U}$ . Since  $\delta_0$  has minimal variance,

$$\begin{aligned} \text{Var}(\delta_\lambda) &= \text{Var}(\delta_0) + \lambda^2 \text{Var}(U) + 2\lambda \text{cov}(\delta_0, U) \\ &\geq \text{Var}(\delta_0) \end{aligned}$$

Consider the quadratic form  $q(\lambda) = \lambda^2 \text{Var}(U) + 2\lambda \text{cov}(\delta_0, U)$ .

The form  $q$  has the roots  $\lambda = 0$  and  $2\text{cov}(\delta_0, U)/\text{var}(U)$ .

If the roots are distinct, then the form must be negative at some point, which would violate the inequality above.

Hence,  $2\text{cov}(\delta_0, U)/\text{var}(U) = 0$  in which case,  $E(U\delta_0) = \text{cov}(\delta_0, U) = 0$ .

To prove the converse result, we assume that  $E(U\delta_0) = 0 \forall U \in \mathcal{U}$  and consider any  $\delta$  unbiased for  $g(\theta)$ . It follows that  $\delta - \delta_0 \in \mathcal{U}$ , so  $E(\delta_0(\delta - \delta_0)) = 0$ . This implies that  $E(\delta_0\delta) = E(\delta_0^2)$  and subtracting  $E(\delta_0)E(\delta)$  on both sides, we obtain

$$\text{Var}(\delta_0) = \text{cov}(\delta_0, \delta) \leq \sqrt{\text{Var}(\delta_0)\text{Var}(\delta)} \quad \text{by Cauchy-Schwarz inequality. Hence, } \text{Var}(\delta_0) \leq \text{Var}(\delta) \text{ for any arbitrary unbiased estimator } \delta \text{ and } \delta_0.$$

Hence,  $\delta_0$  is an UMVUE for  $g(\theta)$ .

To answer the question,  $\forall U \in \mathcal{U}, E((\delta_1 + \delta_2)U) = E(\delta_1 U) + E(\delta_2 U) = 0$  ( $\delta_1, \delta_2$ : UMVUEs)  
 $\Rightarrow \delta_1 + \delta_2$  is a UMVUE for  $g_1(\theta) + g_2(\theta)$

## 4.2 Variance Bound and Information

Recall  $\text{Cov}(X, Y) \leq \sqrt{\text{Var}(X)\text{Var}(Y)}$

Given this inequality, if  $\delta$  is an unbiased estimator for  $g(\theta)$  and  $\psi$  is an arbitrary random variable, then

$$\text{Var}_\theta(\delta) \geq \text{Cov}^2_\theta(\delta, \psi) / \text{Var}_\theta(\psi) \quad \{*\}$$

If we manage to find a suitable  $\psi$  so that the bound is meaningful in the sense that  $\text{Cov}_\theta(\delta, \psi)$  is the same for all  $\delta$  that are unbiased for  $g(\theta)$ .

Let  $P = \{P_\theta : \theta \in \Theta\}$  be a dominated family with densities  $P_\theta : \theta \in \Theta \in \mathbb{R}$ . To begin,  $E_{\theta+\Delta}(\delta) - E_\theta(\delta)$  gives the same value  $g(\theta + \Delta) - g(\theta)$  for any unbiased  $\delta$ . Hence,  $\Delta$  must be chosen so that  $\theta + \Delta \in \Theta$ .

Next, we write  $E_{\theta+\Delta}(\delta) - E_\theta(\delta)$  as a covariance under  $P_\theta$ . This step involves the use of likelihood ratio. We assume here that  $P_{\theta+\Delta}(x) = 0$  whenever  $P_\theta(x) = 0$

Define  $L(x) = P_{\theta+\Delta}(x)/P_\theta(x)$ , where  $P_\theta(x) > 0$  and  $L(x) = 1$  otherwise.

We have

$$L(x)P_\theta(x) = \frac{P_{\theta+\Delta}(x)}{P_\theta(x)}P_\theta(x) = P_{\theta+\Delta}(x) \quad a.s. \quad P.$$

and so, for any function  $h$  integrable under  $P_{\theta+\Delta}$ , we have

$$E_{\theta+\Delta}h(X) = \int hP_{\theta+\Delta}d\mu = \int hLP_\theta d\mu = E_\theta(L(X)h(X))$$

Take  $h = 1$ ,  $E_\theta(L(X)) = 1$  because  $E_\theta(L(X)) = \int \frac{P_{\theta+\Delta}(x)}{P_\theta(x)}P_\theta(x)dx = \int P_{\theta+\Delta}(x)dx = 1$

Take  $h = \delta$ ,  $E_{\theta+\Delta}(\delta) = E_\theta(L\delta)$ . So if we define  $\psi(x) = L(x) - 1$ , we can see that

$$E_\theta(\psi(X)) = E_\theta(L(X) - 1) = 1 - 1 = 0$$

and

$$E_{\theta+\Delta}(\delta) - E_\theta(\delta) = E_\theta(L\delta) - E_\theta(\delta) = E_\theta(\psi\delta)$$

As a result,

$$\text{Cov}_\theta(\delta, \psi) = g(\theta + \Delta) - g(\theta)$$

for any unbiased estimator  $\delta$ . With this particular choice of  $\psi$ , the inequality (\*) can be rewritten as

$$\text{Var}_\theta(\delta) = \frac{\{g(\theta + \Delta) - g(\theta)\}^2}{\text{Var}_\theta(\psi)} = \frac{\{g(\theta + \Delta) - g(\theta)\}^2}{E_\theta\left(\frac{P_{\theta+\Delta}(x)}{P_\theta(x)} - 1\right)^2}, \quad (**)$$

which is known as the Hammersley-Chapman-Robbins Inequality.

Under suitable regularity conditions, we can show that

$$\frac{\lim_{\Delta \rightarrow 0} \left\{ \frac{g(\theta + \Delta) - g(\theta)}{\Delta} \right\}^2}{\lim_{\Delta \rightarrow 0} E_\theta \left( \frac{\{P_{\theta+\Delta}(x) - P_\theta(x)\}/\Delta}{P_\theta(x)} \right)^2} \rightarrow \frac{(g'(\theta))^2}{E_\theta \left( \frac{\frac{\partial}{\partial \theta} P_\theta(x)}{P_\theta(x)} \right)^2} \quad (***)$$

The denominator on the RHS of (\*\*\*) is known as Fisher information, denoted as  $I(\theta)$  and is given by

$$I(\theta) = E_{\theta}\left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right)^2$$

With enough regularity conditions to interchange integration and differentiation,

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta}(1) = \frac{\partial}{\partial \theta} \int P_{\theta}(x) d\mu(x) = \int \frac{\partial}{\partial \theta} P_{\theta}(x) d\mu(x) \\ &= \int \frac{\partial \log P_{\theta}(x)}{\partial \theta} P_{\theta}(x) d\mu(x) = E_{\theta}\left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right) \end{aligned}$$

and so,

$$I(\theta) = E_{\theta}\left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right)^2 - \left(E_{\theta}\left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right)\right)^2 = \text{Var}_{\theta}\left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right).$$

Furthermore, since

$$\int \frac{\partial^2 \log P_{\theta}(x)}{\partial \theta^2} d\mu(x) = E_{\theta}\left(\frac{\partial^2 P_{\theta}(x) / \partial \theta^2}{P_{\theta}(x)}\right) = 0$$

We can see that

$$\begin{aligned} \frac{\partial^2 \log P_{\theta}(x)}{\partial \theta^2} &= \frac{\partial^2 P_{\theta}(x) / \partial \theta^2}{P_{\theta}(x)} - \left(\frac{\partial \log P_{\theta}(x)}{\partial \theta}\right)^2 \\ \implies I(\theta) &= -E_{\theta}\left(\frac{\partial^2 \log P_{\theta}(x)}{\partial \theta^2}\right) \end{aligned}$$

Therefore,

$$\text{Var}_{\theta}(\delta) \geq \frac{\{g'(\theta)\}^2}{I(\theta)}, \quad \theta \in \Theta$$

**Theorem 4.2** Let  $P = \{P_{\theta} : \theta \in \Theta\}$  be a dominated family with  $\Theta$  and open set in  $\mathbb{R}$  and densities  $P_{\theta}$  differentiable with respect to  $\theta$ . If  $E_{\theta}(\psi) = 0$  and  $E_{\theta}(\delta^2) < \infty$ , then

$$\text{Var}_{\theta}(\delta) \geq \frac{\{g'(\theta)\}^2}{I(\theta)}, \quad \theta \in \Theta.$$

This result is called the Cramér-Rao / Information Bound.

**Example 4.2.1** Let  $\mathcal{P}$  be a one-parameter exponential family in canonical form and density  $p_{\eta}$  given by

$$p_{\eta} = \exp(\eta T(x) - A(\eta)) h(x),$$

then

$$\frac{\partial \log p_{\eta}(x)}{\partial \eta} = T(x) - A'(\eta).$$

By the previous results, we have

$$I(\eta) = \text{Var}_\eta(T(x) - A'(\eta)) = \text{Var}_\eta(T(x)) = A''(\eta).$$

If the family is parameterised instead by  $\mu = A'(\eta) = \eta\mathbb{T}(\mathbb{X})$ , then

$$A''(\eta) = I(\mu)(A'(\eta))^2,$$

and so, because  $A''(\eta) = \text{Var}(T)$ , we have  $I(\mu) = \frac{1}{\text{Var}_\eta(T)}$ . Observe also that because  $T$  is UMVUE for  $\mu$ , the lower bound variance  $\text{Var}(\delta) \geq (I(\mu))^{-1}$  for unbiased estimation  $\delta$  of  $\mu$  is

**Example 4.2.2** Suppose  $X$  is an absolutely continuous random variable with density  $f$ . The family of distributions  $\mathcal{P} = \{P_\theta : \theta \in \mathbb{R}\}$  with  $P_\theta$  the distribution of  $\theta + \epsilon$  is called a location family.

$$\begin{aligned} \int g(x) dP_\theta(x) &= \mathbb{E}\theta(g(X)) \\ &= \mathbb{E}_\theta(g(\theta + \epsilon)) \\ &= \int g(\theta + \epsilon) f(\epsilon) d\epsilon \\ &= \int g(x) f(x - \theta) dx. \end{aligned}$$

So  $P_\theta$  has the density  $p_\theta = f(x - \theta)$ . The corresponding Fisher information for the family is

$$\begin{aligned} I(\theta) &= \mathbb{E}_\theta \left( \frac{\partial \log f(X - \theta)}{\partial \theta} \right)^2 \\ &= \mathbb{E}_\theta \left( -\frac{f'(X - \theta)}{f(X - \theta)} \right)^2 \\ &= \mathbb{E} \left( \frac{f'(\epsilon)}{f(\epsilon)} \right)^2 \\ &= \int \frac{\{f'(x)\}^2}{f(x)} dx \perp \theta. \end{aligned}$$

So, for the location families,  $I(\theta)$  is constant with respect to  $\theta$ .

If two (or more) independent vectors are observed, then the total Fisher information is the sum of the Fisher information provided by the individual observations.

Suppose  $X$  and  $Y$  are independent, and that  $X$  has density  $p_\theta$  and  $Y$  has density  $q_\theta$ . The Fisher information from  $X$  and  $Y$  are respectively

$$I_X(\theta) = \text{Var}_\theta \left( \frac{\partial \log p_\theta(X)}{\partial \theta} \right),$$

and

$$I_Y(\theta) = \text{Var}_\theta \left( \frac{\partial \log q_\theta(Y)}{\partial \theta} \right).$$

$$\begin{aligned} I_{X,Y}(\theta) &= \text{Var}_\theta \left( \frac{\partial \log \{p_\theta(X)q_\theta(Y)\}}{\partial \theta} \right) \\ &= \text{Var}_\theta \left( \frac{\partial \log p_\theta(X)}{\partial \theta} + \frac{\partial \log q_\theta(Y)}{\partial \theta} \right) \\ &= \text{Var}_\theta \left( \frac{\partial \log p_\theta(X)}{\partial \theta} \right) + \text{Var}_\theta \left( \frac{\partial \log q_\theta(Y)}{\partial \theta} \right) \\ &= I_X(\theta) + I_Y(\theta). \end{aligned}$$

$\Rightarrow$  If  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P_\theta$ ,  $I_X(\theta) = I_{X_1}(\theta) + \dots + I_{X_n}(\theta) = nI_{X_1}(\theta)$ .

$\Rightarrow \text{Var}_\theta(\delta) \geq \frac{g'(\theta)}{nI(\theta)}$

### 4.3 Multi-parameter Cramer-Rao Inequality

Suppose the following conditions hold:

1.  $\Theta \subseteq \mathbb{R}^k$  is an open set,
2.  $\{P_\theta : \theta \in \Theta\}$  have common support  $I$ ,
3.  $\frac{\partial p_\theta(X)}{\partial \theta_i}$  exists,  $\forall i = 1, \dots, k, x \in I$  and is finite,
4.  $\frac{\partial \int_x p_\theta(x) d\mu}{\partial \theta_i} = \int_x \frac{\partial p_\theta(x)}{\partial \theta_i} d\mu, \forall i = 1, \dots, k$ ,
5.  $\frac{\partial \int_x \delta(x) p_\theta(x) d\mu}{\partial \theta_i} = \int_x \delta(x) \frac{\partial p_\theta(x)}{\partial \theta_i} d\mu, \forall i = 1, \dots, k$

Define the  $k \times k$  information matrix  $I(\theta)$  by

$$I_{ij}(\theta)_{i,j=1,\dots,k}, \text{ with } I_{ij}(\theta) = \mathbb{E}_\theta \left( \frac{\partial \log p_\theta(X)}{\partial \theta_i} \frac{\partial \log p_\theta(X)}{\partial \theta_j} \right)$$

Specially, if  $k = 1$ ,  $I_{11} = \mathbb{E}_\theta \left( \frac{\partial \log p_\theta(X)}{\partial \theta} \right)^2$ . Assume that  $I(\theta)$  is finite and positive definite, then

$$\text{Var}_\theta(\delta(X)) \geq \alpha^T I(\theta)^{-1} \alpha, \text{ with } \alpha_i = \frac{\partial g(\theta)}{\partial \theta_i}.$$

**Proof:** Let  $\psi_i(X) = \frac{\partial \log p_\theta(X)}{\partial \theta_i}$ , then

$$\begin{aligned} \mathbb{E}_\theta(\psi_i(X)) &= \int_x \left( \frac{\partial \log p_\theta(x)}{\partial \theta_i} \right) p_\theta(x) d\mu(x) \\ &= \int_x \frac{\frac{\partial p_\theta(x)}{\partial \theta_i}}{p_\theta(x)} p_\theta(x) d\mu \\ &= \int_x \frac{\partial p_\theta(x)}{\partial \theta_i} d\mu \\ &= \frac{\partial \int_x p_\theta(x) d\mu}{\partial \theta_i} = 0 \end{aligned}$$

Fix a non-zero vector  $(a_1, \dots, a_k)$ . Then  $\mathbb{E}_\theta(\sum_{i=1}^k a_i \psi_i(X)) = 0$ . Claim:  $\text{Var}(\sum_{i=1}^k a_i \psi_i(X)) = a^T I(\theta) a$ .

$$\begin{aligned} \text{Var}\left(\sum_{i=1}^k a_i \psi_i(X)\right) &= \sum_{i,j} a_i a_j \text{Cov}(\psi_i(X), \psi_j(X)) \\ &= \sum_{i,j} a_i a_j \mathbb{E}(\psi_i(X) \psi_j(X)) \\ &= \sum_{i,j} a_i a_j I_{ij}(\theta) \\ &= a^T I(\theta) a. \end{aligned}$$

Finally,

$$\begin{aligned}
\text{Cov}(\delta(X), \sum_{i=1}^n k a_i \psi_i(X)) &= \sum_{i=1}^n k a_i \text{Cov}(\delta(X), \psi_i(X)) \\
&= \sum_{i=1}^n k a_i \mathbb{E}(\delta(X) \psi_i(X)) \\
&= \sum_{i=1}^n k a_i \int_x \delta(x) \frac{\partial \log p_\theta(x)}{\partial \theta_i} p_\theta(x) d\mu \\
&= \sum_{i=1}^n k a_i \int_x \delta(x) \frac{\partial p_\theta(x)}{\partial \theta_i} dx \\
&= \sum_{i=1}^n k a_i \frac{\partial \int_x \delta(x) p_\theta(x) d\mu}{\partial \theta_i} \\
&= \sum_{i=1}^n k a_i \alpha_i(\theta).
\end{aligned}$$

By Cauchy-Schwarz inequality,

$$\text{Var}(\delta(X)) \text{Var}(\sum_{i=1}^n k a_i \psi_i(X)) \geq \text{Cov}(\sum_{i=1}^n k a_i \psi_i(X), \delta(X)),$$

$\Rightarrow$

$$\text{Var}(\delta(X)) \geq \sup_{a \neq 0} \frac{(\sum_{i=1}^n k a_i \alpha_i(\theta))^2}{a^T I(\theta) a} = \alpha^T I(\theta)^{-1} \alpha.$$

■

**Example 4.3.1** Suppose  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$ ,  $\mu \in \mathbb{R}$ , and  $\sigma^2 > 0$ . We want to estimate  $g_1(\mu, \sigma^2) = \mu$  and  $g_2(\mu, \sigma^2) = \sigma^2$ . Look at only unbiased estimators.

*Claim 1:*  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is UMVUE for  $\sigma^2$ ,  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

**Proof:**

1.  $(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$  is complete and sufficient.
2.  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  is a function of  $(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ . This is because  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n) = \sum_{i=1}^n X_i^2 - n\bar{X}_n^2$ .
3.  $\frac{1}{\sigma^2} \sim \chi_{n-1}^2 \Rightarrow \mathbb{E}(\frac{1}{n-1} \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2) = \frac{\sigma^2}{n-1} (n-1) = \sigma^2$

■

*Claim 2:*  $\bar{X}_n$  is UMVNE for  $\mu$ .  $\mathbb{E}(\bar{X}_n - \mu)^2 = \sigma^2/n$ . Note that  $\mathbb{E}(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 / (n-1) - \sigma^2) = \text{Var}(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 / (n-1)) = \frac{2(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{(n-1)}$ .

$$\log p_{\mu, \sigma^2}(X) = -\frac{n}{2} \log \sigma^2 - \sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2} + C_n,$$

$$\frac{\partial \log p_\theta(X)}{\partial \mu} = \sum_{i=1}^n \frac{X_i - \mu}{\sigma^2},$$

$$\begin{aligned}\frac{\partial^2 \log p_\theta(X)}{\partial \mu^2} &= -\frac{n}{\sigma^2}, \\ \frac{\partial \log p_\theta(X)}{\partial \sigma^2} &= -\frac{n}{\sigma^2} + \sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^4} \\ \frac{\partial^2 \log p_\theta(X)}{\partial (\sigma^2)^2} &= \frac{n}{\sigma^4} - \sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^6} \\ \frac{\partial^2 \log p_\theta(X)}{\partial \mu \partial \sigma^2} &= -\sum_{i=1}^n \frac{X_i - \mu}{\sigma^4}\end{aligned}$$

Therefore, the Fisher information matrix is

$$I = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}$$

$$I_{22} = -\frac{n}{2\sigma^4} + \mathbb{E}\left(\sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^6}\right) = -\frac{n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6} = \frac{n}{2\sigma^4}$$

$$\Rightarrow \text{CRLB for } \mu = \begin{pmatrix} 1 & 0 \end{pmatrix} I^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \frac{1}{I_{11}} = \frac{\sigma^2}{n}.$$

$$\Rightarrow \text{CRLB for } \sigma^2 = \begin{pmatrix} 0 & 1 \end{pmatrix} I^{-1} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \frac{1}{I_{22}} = \frac{2\sigma^4}{n}.$$