

CRLB \neq UMVU

Example: $x_i \sim N_0(\vartheta, \sigma^2)$ σ^2 : known, fixed
 $i=1, \dots, n$

$$g(\vartheta) = \vartheta$$

Joint probability distribution:

$$P_\vartheta(x_1, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \vartheta)^2}{2\sigma^2}\right)$$

log likelihood

$$\log P_\vartheta(x_1, \dots, x_n) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \vartheta)^2$$

→ calculate CRLB.

$$\frac{\partial}{\partial \vartheta} \log P_\vartheta(x) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \vartheta) = \frac{n(\bar{x} - \vartheta)}{\sigma^2}$$

(how much the loglikelihood changes when ϑ changes)

Fisher Information Matrix:

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

$$= \frac{n^2}{\sigma^4} \text{Var}(\bar{x}) = \frac{n^2}{\sigma^4} \cdot \frac{\sigma^2}{n} = \frac{n}{\sigma^2}$$

$$I_n(\vartheta) = n/\sigma^2$$

The UMVU estimator for θ

$$T = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{sufficient, complete statistic}$$

$$\text{Var}_{\theta}(T) = \frac{\sigma^2}{n}$$

T attains the CRLB \Leftrightarrow UMVU for θ .

An alternative estimator for θ that attains CRLB but it is not UMVU

$$T(x_1, \dots, x_n) = x_1 \quad (\text{first sample}) \text{ for } g(\theta) = \theta$$

$$E[T] = E[x_1] = \theta \quad (\text{unbiased})$$

$$\text{Var}_{\theta}[T] = \text{Var}_{\theta}(x_1) = \sigma^2$$

$$\text{CRLB} = \frac{1}{n} (I(\theta))^{-1} = \frac{\sigma^2}{n}$$

$$\text{Var}_{\theta}(x_1) = \sigma^2 > \frac{\sigma^2}{n}$$

Estimator unbiased but does not attain CRLB.

Alternative estimator for θ : $\hat{\theta} = 0$
randomized.

$$T(x_1, \dots, x_n, u) = \begin{cases} x_1 & u < \frac{1}{n} \\ x_2 & \frac{1}{n} \leq u < \frac{2}{n} \\ \vdots & \vdots \\ x_n & \frac{n-1}{n} \leq u < 1 \end{cases}$$

We pick one sample at random (uniformly)

from $\{x_1, \dots, x_n\}$. This estimator is not
deterministic in data.

x_i are i.i.d. and we randomly pick one:

$$E[T] = E_u[E[T|u]] = \frac{1}{n} \sum_{i=1}^n E(x_i) = \theta. \quad (\text{unbiased})$$

$$\text{Var}_\theta(T) = \text{Var}_\theta(\text{randomly pick from } \{x_i\}) =$$

$$= \frac{1}{n} \sum \text{var}(x_i) = \sigma^2$$

$$\text{CRLB} := \frac{\sigma^2}{n} \Rightarrow \text{Var}_\theta(T) > \frac{\sigma^2}{n} \text{ CRLB}$$

Estimator T does not attain the CRLB
Not UMVU because T is not a function
of the complete sufficient statistic.

CRLB \Rightarrow UMVU estimator.

Exponential families & CR LB

We consider 1-dimensional exponential family:

$$p_\theta(x) = c \exp(C(\theta)T(x) - d(\theta)) h(x)$$

where: $\theta \in \mathbb{R}$: scalar parameter

$T(x)$: sufficient statistic

$c(\theta)$: normal parameter

$d(\theta)$: function of parameter

$h(x)$: function of data

Theorem: Let $X_i \stackrel{i.i.d}{\sim} p_\theta(x)$ form the exponential family. Then the sufficient statistic T is an efficient estimator for $g(\theta) = E_\theta(T)$

This means that $T(x)$ is unbiased and its variance achieves CRLB. \square

Proof: $S_{\theta}(x) = \frac{\partial}{\partial \theta} \log p_{\theta}(x) = -d'(\theta) + c'(\theta)T(x)$

The score is linear in the sufficient statistic

$$\begin{aligned} T. \\ \Rightarrow E_{\theta} \left(S_{\theta}(x) (T(x) - E_{\theta}[T(x)]) \right) &= \frac{d}{d\theta} E_{\theta}[T(x)] = g'(\theta) \\ \uparrow \end{aligned}$$

By C-R proof: Covariance between the score and the estimators gives the derivative of the expectation

By C-S inequality $(E[X \cdot Y])^2 \leq E(X^2) \cdot E(Y^2)$

$$x = S(\theta)$$

$$y = T(x) - E_{\theta}[T(x)]$$

$$(g'(\theta))^2 \leq E_{\theta}(S_{\theta}(x)^2) \cdot \text{Var}_{\theta}(T) \Rightarrow$$

$$\text{Var}_{\theta}(T) \geq \frac{(g'(\theta))^2}{E_{\theta}(S_{\theta}(x))^2} = \frac{(g'(\theta))^2}{I(\theta)}$$

$$\Rightarrow \text{Var}_{\theta}(T) = \frac{(g'(\theta))^2}{I(\theta)}$$

Because the sufficient statistic $T(x)$ is linearly related to the score function.

$$T(x) - E_\theta(T(x)) = \lambda \cdot S_\theta(x) \quad \lambda: \text{constant}$$

$$\bar{T}(\cdot) = E_\theta(T(\cdot)) + \lambda \cdot S_\theta(\cdot)$$

$$S_\theta(\cdot) = C'(\theta) T(\cdot) - d'(\theta)$$

$$T(x) = \frac{1}{C'(\theta)} \cdot S_\theta(x) + \frac{d'(\theta)}{C'(\theta)}$$

Lemma: Assume A1-A3 hold, and that $T: X \rightarrow \mathbb{R}$ is unbiased estimator for θ . Then there exists $c(\theta), d(\theta)$ differentiable such that

$$P_\theta(x) = \exp(c(\theta)T(x) - d(\theta) h(x)) \neq 0 \forall \theta \in \Theta.$$

By Cramér - Rao lower bound: Let T be an unbiased estimator of $g(\theta) \in \mathbb{R}^k$ with $E_\theta(T) < \infty$, $\forall \theta \in \Theta$. Then assuming I

is invertible:

$$\text{Cov}(\tau) \geq D_{\theta}g(\theta)(I(\theta))^{-1}D_{\theta}g(\theta)^T.$$

Proof: T_j unbiased $\Rightarrow \text{bias}_\theta(T_j) = 0 \forall j \in \{1, \dots, k\}$

$$\Rightarrow \text{bias}_\theta\left(\sum_j \alpha_j T_j\right) = 0 \quad \sum \alpha_j \frac{\partial}{\partial \theta_j} g(\theta)$$

$$\Rightarrow a^T \text{Cov}(\tau) a = \text{Cov}(a\tau) \\ \geq a^T D_{\theta}g(\theta)(I(\theta))^{-1}D_{\theta}g(\theta)^T a$$

□

Fisher Information Matrix:

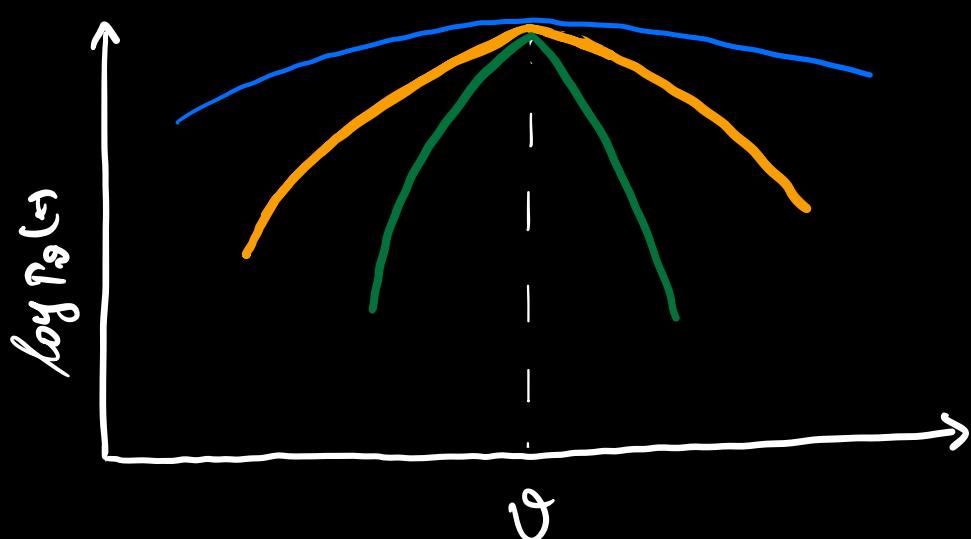
Definition: We define the Fisher Information Matrix as:

$$I(\theta)_{ij} = -\mathbb{E}_{\theta}\left[\frac{\partial}{\partial \theta_i} \log p_{\theta}(x) \frac{\partial}{\partial \theta_j} \log p_{\theta}(x)\right]$$

for $i, j \in \{1, \dots, k\}$

- Note:
- Fisher information matrix is a covariance matrix.
 - FIM shows us how sensitive the model is to changes of θ .

Curvature of log likelihood function



- $G^2 = 2.0$ sharp peak
- $G^2 = 1.0$ moderate curvature
- $G^2 = 0.2$ flat curve

$$\text{Curvature} = - \frac{d^2}{d\theta^2} \log P_\theta(x) = \frac{1}{G^2}$$

more curvature \Rightarrow sharper likelihood

\Leftrightarrow smaller estimator variance

high fisher information matrix
 \Rightarrow low uncertainty

low Fisher information matrix
→ higher uncertainty.

Kullback-Leibler divergence:

Q: If I observe some data, can I tell whether it came from $P_{\theta_1}(x)$ or $P_{\theta_2}(x)$? Is it easy to distinguish between two parameter choices θ_1, θ_2 in the same statistical model?

Example 1: $P_{\theta_1}(x) = \mathbb{I}(x \in \{1, 2\})$
 $P_{\theta_2}(x) = \mathbb{I}(x \in \{2, 3\})$

(a) if $x = 1$ $P_{\theta_1}(x) = 1$ $P_{\theta_2}(x) = 0$

(b) if $x = 2$ $P_{\theta_1}(x) = 1$ $P_{\theta_2}(x) = 1$

(c) if $x = 3$ $P_{\theta_1}(x) = 0$ $P_{\theta_2}(x) = 1$

(a), (c) are distinguishable.

(b) distributions are identical. I cannot tell them apart.

Definition: The discriminating power between θ_1, θ_2 of result x is

$$h(\theta_1, \theta_2)(x) = \log \frac{P_{\theta_1}(x)}{P_{\theta_2}(x)}$$

defined on $\{\sum P_{\theta_1}(x) > 0\} \cap \{\sum P_{\theta_2}(x) > 0\}$

This power tells you how much evidence the data x gives in favor of θ_1 over θ_2 .

To get an overall measure of distinguishing, we average the log likelihood ratio under P_{θ_1} :

$$\begin{aligned} D_{KL}(P_{\theta_1} \| P_{\theta_2}) &= E_{\theta_1} \left[\log \frac{P_{\theta_1}(x)}{P_{\theta_2}(x)} \right] \\ &= \int \log \frac{P_{\theta_1}(x)}{P_{\theta_2}(x)} P_{\theta_1}(x) \mu(dx) \end{aligned}$$

Properties of KL divergence:

$$1. D_{KL}(P_{\theta_1} \| P_{\theta_2}) \geq 0$$

$$2. D_{KL}(P_{\theta_1} \| P_{\theta_2}) = 0 \iff P_{\theta_1}(x) = P_{\theta_2}(x) \text{ } \mu\text{-a.s.}$$

Example: (Poisson)

$$\text{Compare } P_{\theta_1}(x) = \frac{\theta_1^x e^{-\theta_1}}{x!}, \quad P_{\theta_2}(x) = \frac{\theta_2^x e^{-\theta_2}}{x!}$$

KL divergence:

$$D_{KL}(P_{\theta_1} \| P_{\theta_2}) = \mathbb{E}_{\theta_1} \left[\log \frac{P_{\theta_1}(x)}{P_{\theta_2}(x)} \right]$$

$$= \mathbb{E}_{\theta_1} \left[x \log \left(\frac{\theta_1}{\theta_2} \right) - (\theta_1 - \theta_2) \right]$$

Since $\mathbb{E}_{\theta_1}(x) = \theta_1$,

$$= \theta_1 \left(\log \left(\frac{\theta_1}{\theta_2} \right) - (\theta_1 - \theta_2) \right)$$

Remark: KL divergence is not symmetric

$$D_{KL}(P_{\theta_1} \| P_{\theta_2}) \neq D_{KL}(P_{\theta_2} \| P_{\theta_1})$$

Lemma: Let Y r.o. with $E[|Y|] = E[Y^+] + E[Y^-] < \infty$
For any convex function $Q: Y \rightarrow \mathbb{R}$, $E[Q(Y)] < \infty$

Proof: Consider $\gamma = \frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)}$

$$E_{\mathcal{D}_1}(|\gamma|) = \int \frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)} \cdot P_{\mathcal{D}_1}(x) \mu(dx) = 1.$$

$$\Rightarrow E_{\mathcal{D}_1} \left(\left(\log \frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)} \right) \right) = E_{\mathcal{D}_1} \left((-\log \gamma) \right) < \infty$$

\Rightarrow The integral is well defined on $\mathbb{R} \setminus \text{out}$

By Jensen inequality

$$D_{KL}(P_{\mathcal{D}_1} || P_{\mathcal{D}_2}) = E_{\mathcal{D}_1} \left(-\log \frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)} \right) \geq -\log E_{\mathcal{D}_1} \left[\frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)} \right] \\ = -\log 1 = 0$$

To show the reverse direction we need reverse Jensen.

$$-\log \gamma \text{ Strictly Convex} \Rightarrow \gamma = \frac{P_{\mathcal{D}_2}(x)}{P_{\mathcal{D}_1}(x)} \equiv c^{-1}$$

$P_{\mathcal{D}_1}$ -a.s. 

Lemma: Let Y RV with $E[Y]$ well defined on $\text{RV}\{\infty\}$ and ϕ is a convex function Then

$$E[\phi(Y)] \geq \phi(E(Y))$$

with equality if -f $\phi(y) = ay + b$ P -a.s.

In particular if ϕ is strictly convex

then $y = c$ P -a.s.

Theorem: Let A1-A5 hold, and assume

$$\text{that } \frac{\partial^2}{\partial \theta_i \partial \theta_j} \int \log \frac{P_{\theta_i}(x)}{P_{\theta_0}(x)} P_{\theta_0}(x) \mu(dx)$$

$$= \int \frac{\partial^2}{\partial \theta_i \partial \theta_j} \left(\log \frac{P_{\theta_1}(x)}{P_{\theta_2}(x)} P_{\theta_2}(x) \right) \mu(dx)$$

$$\text{then } I(\theta)_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} D_H(\theta | \theta') \Big|_{\theta = \theta'}$$

Fisher Information Matrix

