

Foundations of Statistical Inference

Julien Berestycki

Department of Statistics
University of Oxford

MT 2016

Lecture 3 - Estimators, Minimum Variance Unbiased Estimators and the Cramér-Rao Lower Bound.

Estimators

Definition

A **point estimate** for θ is a statistic of the data.

$$\hat{\theta} = \hat{\theta}(x) = t(x_1, \dots, x_n).$$

An **interval estimate** is a set valued function $C(X) \subseteq \Theta$ such that $\theta \in C(X)$ with a specified probability.

Definition (Maximum likelihood estimation)

If $L(\theta)$ is differentiable and there is a unique maximum in the interior of $\theta \in \Theta$, then the MLE $\hat{\theta}$ is the solution of

$$\frac{\partial}{\partial \theta} L(\theta; x) = 0 \text{ or } \frac{\partial}{\partial \theta} \ell(\theta) = 0,$$

where $\ell(\theta) = \log L(\theta; x)$.

Lemma 2 : MLEs and exponential families

Consider a k -dimensional exponential family in canonical form

$$L(\theta; x) = \exp \left\{ \sum_{j=1}^k \phi_j \left(\sum_{i=1}^n B_j(x_i) \right) + nD(\phi) + \sum_{i=1}^n C(x_i) \right\}.$$

Let $T_j(X) = \sum_{i=1}^n B_j(X_i)$, $j = 1, \dots, k$. If the realized data are $X = x$, then the statistics evaluated on the data are $T_j(x) = t_j$.

Theorem

The MLEs of ϕ_1, \dots, ϕ_k are the solution of

$$t_j = \mathbb{E}_X(T_j), j = 1, \dots, k.$$

i.e. set the expected values of the sufficient statistics equal to their realised values and solve for ϕ_j . [If the family is not in canonical form, there is a similar slightly more complicated matrix equation]

Proof

$$\ell = \log L = \text{const} + \sum_{j=1}^k \phi_j t_j + nD(\phi)$$

$$\Rightarrow \frac{\partial}{\partial \phi_j} \ell = t_j + n \frac{\partial}{\partial \phi_j} D(\phi)$$

However, since $\mathbb{E}_X[B_i(X)] = -\frac{\partial}{\partial \phi_i} D(\phi)$ and $T_j(X) = \sum_{i=1}^n B_j(X_i)$ we know that

$$\mathbb{E}_X[T_j] = -n \frac{\partial}{\partial \phi_j} D(\phi), \text{ so}$$

$$\frac{\partial}{\partial \phi_j} \ell = t_j - \mathbb{E}_X(T_j) = 0$$

is equivalent to $t_j = \mathbb{E}_X(T_j)$.

Bias, Variance, Mean Squared Error

Definition

A statistic $T_n = T(X_1, \dots, X_n)$ is **unbiased** for a function $g(\theta)$ if

$$\mathbb{E}_X(T_n) = \int_{\mathcal{X}} t_n(x) f(x; \theta) dx = g(\theta), \text{ for all } \theta \in \Theta.$$

The **bias** of an estimator T_n is $\text{bias}(T_n) = \mathbb{E}_X [T_n - g(\theta)]$

T_n is a **consistent** estimator if

$$\forall \epsilon > 0, P(|T_n - \theta| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

The **Mean Squared Error** (MSE) of T_n is

$$\text{MSE}(T_n) = \mathbb{E}_X [T_n - g(\theta)]^2 = V_X(T_n) + [\text{bias}(T_n)]^2$$

Example 10 $N(\mu, \sigma^2)$. $\hat{\mu} = \bar{X}$ and $S^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$ are unbiased estimates of μ and σ^2 .

Minimum Variance Unbiased Estimators (MVUE)

- If we want to find a good estimator then one obvious strategy is to try to find estimators that minimise MSE. This is often difficult.
- For example, if we choose the estimator $\hat{\theta} = \theta_0$ then this has $\text{MSE}=0$ when $\theta = \theta_0$, so no other estimator can be uniformly best unless it has zero MSE everywhere.
- If we restrict attention to unbiased estimators then the situation becomes more tractable. In this case, MSE reduces to the variance of the estimator and we can focus on minimising the variance of estimators. That is, we search for minimum variance unbiased estimators (MVUE).

Theorem 2 : Cramér-Rao inequality (and bound).

Theorem

If $\hat{\theta}$ is an unbiased estimator of θ , then subject to certain regularity conditions on $f(x; \theta)$, we have

$$\text{Var}(\hat{\theta}) \geq I_{\theta}^{-1}.$$

where I_{θ} , the Fisher information, is given by

$$I_{\theta} = -\mathbb{E}_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta) \right]$$

Comment If an estimator achieves the bound then it is MVUE. There is no guarantee that the bound will be attainable. In many cases it is attainable asymptotically. Intuitively, the more 'information' we have about θ , the larger I_{θ} will be and lowest possible variance of the estimator will be smaller.

Regularity conditions for CRLB

- We will not be concerned with the details of the required regularity conditions.
- The main reason they are needed is to ensure that it is ok to interchange integration and differentiation during parts of the proof.
- One condition that is often easy to check is that the range of the rv X must not depend on θ . So for example, the result can not be applied when working with the uniform distribution $U[0, \theta]$ and we wish to estimate θ .

In order to prove the CRLB we will need to use a few results.

Lemma (Variance-Covariance inequality)

Let U and V be scalar rv. Then

$$\text{cov}(U, V)^2 \leq \text{var}(U)\text{var}(V)$$

with equality if and only if $U = aV + b$ for constants and $a \neq 0$.

The Fisher Information I_θ , which is used in the Cramér-Rao lower bound, can be expressed in two different forms.

Lemma

Under regularity conditions

$$I_\theta = -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta) \right] = \mathbb{E}_\theta \left[\left(\frac{\partial \ell}{\partial \theta} \right)^2 \right] = \text{Var}[S(X; \theta)],$$

where the **score function** $s(x; \theta)$ is defined as

$$s(x; \theta) = \frac{\partial}{\partial \theta} \ell(\theta) = \frac{f'(x; \theta)}{f(x; \theta)}$$

Lemma 3 - Proof

We need to prove $-\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta) \right] = \mathbb{E} \left[\left(\frac{\partial \ell}{\partial \theta} \right)^2 \right]$.

$$\begin{aligned} \frac{\partial^2 \ell}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \left\{ \frac{1}{L} \frac{\partial L}{\partial \theta} \right\} \quad \left[\text{since } \frac{\partial \ell}{\partial \theta} = \frac{1}{L} \frac{\partial L}{\partial \theta} \right] \\ &= -\frac{1}{L^2} \left(\frac{\partial L}{\partial \theta} \right)^2 + \frac{1}{L} \frac{\partial^2 L}{\partial \theta^2} \\ &= -\left(\frac{\partial \ell}{\partial \theta} \right)^2 + \frac{1}{L} \left(\frac{\partial^2 L}{\partial \theta^2} \right) \end{aligned}$$

The second term has expectation zero because

$$\mathbb{E} \left[\frac{1}{L} \left(\frac{\partial^2 L}{\partial \theta^2} \right) \right] = \int \frac{1}{L} \frac{\partial^2 L}{\partial \theta^2} L dx = \int \frac{\partial^2 L}{\partial \theta^2} dx = \frac{\partial^2}{\partial \theta^2} \int L dx = 0$$

The alternative form $I_\theta = \text{Var}[S(X; \theta)]$ follows from $\mathbb{E} \left[\frac{\partial \ell}{\partial \theta} \right] = 0$.

Proof of the CRLB

We consider only unbiased estimators, so we have

$$\mathbb{E}(\hat{\theta}) = \int_{\mathcal{X}} \hat{\theta}(x) L(\theta; x) dx = \theta$$

Differentiate both sides w.r.t. θ

$$\int_{\mathcal{X}} \hat{\theta} \frac{\partial L}{\partial \theta} dx = 1$$

Now

$$\frac{\partial L}{\partial \theta} = L \frac{\partial \ell}{\partial \theta}$$

so

$$1 = \int_{\mathcal{X}} \hat{\theta} \frac{\partial \ell}{\partial \theta} L dx = \mathbb{E} \left[\hat{\theta} \frac{\partial \ell}{\partial \theta} \right]$$

Proof of the CRLB

Now we use the inequality that for two random variables U, V

$$\text{Cov}[U, V]^2 \leq \text{Var}[U]\text{Var}[V]$$

with $U = \hat{\theta}$, $V = \frac{\partial \ell}{\partial \theta}$. We know $\text{Var}[\frac{\partial \ell}{\partial \theta}] = I_{\theta}$. Must show $\text{Cov}[U, V] = 1$.

$$\text{Cov}[U, V] = \mathbb{E}[UV] - \mathbb{E}[U]\mathbb{E}[V], \quad \mathbb{E}[U] = \theta, \quad \mathbb{E} \left[\hat{\theta} \frac{\partial \ell}{\partial \theta} \right] = 1$$

$$\mathbb{E}[V] = \int_{\mathcal{X}} \frac{\partial \ell}{\partial \theta} L dx = \int_{\mathcal{X}} \frac{\partial L}{\partial \theta} dx = \frac{\partial}{\partial \theta} \left[\int_{\mathcal{X}} L dx \right] = \frac{\partial}{\partial \theta} [1] = 0$$

$$\text{Var}[\hat{\theta}] = \text{Var}[U] \geq \frac{\text{Cov}[U, V]^2}{\text{Var}[V]} = \frac{1^2}{I_{\theta}} = I_{\theta}^{-1}$$

Information in a sample of size n .

If we have n iid observations then

$$f(x; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

and the Fisher information is

$$I_n(\theta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta) \right] = -\int \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(x_i; \theta) f(x; \theta) dx = n i_1(\theta).$$

That is, $i_1(\theta)$ is calculated from the density as

$$i_1(\theta) = -\int \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) f(x; \theta) dx$$

Question Under what conditions will we be able to attain the Cramér-Rao bound and find a MVUE?

Corollary (1)

There exists an unbiased estimator $\hat{\theta}$ which attains the CR lower bound (under regularity conditions) if and only if

$$S(x, \theta) = \frac{\partial \ell}{\partial \theta} = I_{\theta}(\hat{\theta} - \theta)$$

Proof In the CR proof

$$\text{Cov}[U, V]^2 \leq \text{Var}[U]\text{Var}[V]$$

and the lower bound is attained if and only equality is achieved. If $U = \hat{\theta}$, $V = \frac{\partial \ell}{\partial \theta}$, the equality occurs when $\frac{\partial \ell}{\partial \theta} = c + d\hat{\theta}$, where c, d are constants. $\mathbb{E}[V] = 0$ so $c = -d\theta$ and $\frac{\partial \ell}{\partial \theta} = d(\hat{\theta} - \theta)$.

Multiply by $\partial\ell/\partial\theta$ and take expectations.

$$\mathbb{E} \left[\left(\frac{\partial\ell}{\partial\theta} \right)^2 \right] = d \mathbb{E} \left[\frac{\partial\ell}{\partial\theta} \hat{\theta} \right] - d\theta \mathbb{E} \left[\frac{\partial\ell}{\partial\theta} \right] = d \times 1 - 0$$

The LHS is I_θ so we have $d = I_\theta$ and

$$\frac{\partial\ell}{\partial\theta} = I_\theta(\hat{\theta} - \theta)$$

Question What is the relationship between the CRLB and exponential families?

Corollary (2)

If there exists an unbiased estimator $\hat{\theta}(X)$ which attains the CR lower bound (under regularity conditions) it follows that X must be in an exponential family

Proof Taking $n = 1$

$$\frac{\partial \log f(x; \theta)}{\partial\theta} = \frac{\partial\ell}{\partial\theta} = I_\theta(\hat{\theta} - \theta)$$

and

$$\log f(x; \theta) = \hat{\theta}A(\theta) + D(\theta) + C(x)$$

which is an exponential family form. The constant of integration $C(x)$ is a function of x .

Question What is the relationship between the CRLB and MLEs?

Corollary (3)

Suppose $\tilde{\theta}(X)$ is an unbiased estimator that attains the CRLB, and so is a MVUE. Suppose that the MLE $\hat{\theta}$ is a solution to $\partial\ell/\partial\theta = 0$ (so, not on boundary). Then $\tilde{\theta} = \hat{\theta}$.

i.e. if the CRLB is attained then it is generally the MLE that attains it.

Proof $\tilde{\theta}$ must satisfy $\frac{\partial\ell}{\partial\theta} = I_\theta(\tilde{\theta} - \theta)$.

Setting $\frac{\partial\ell}{\partial\theta} = 0$ and solving will give the MLE $\hat{\theta}$.

Since $I_\theta > 0$ (in all but exceptional circumstances), this gives $\tilde{\theta} = \hat{\theta}$.

Question Do all MLEs attain the CRLB?

No, because not all MLEs are unbiased.

Example 11

Let X_1, \dots, X_n be a random sample from $N(\mu, \sigma^2)$.

Then we know the MLEs are $\hat{\mu} = \bar{X}$, $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2$.

Exercise $\hat{\mu}$ is unbiased, but $\hat{\sigma}^2$ is biased. CRLBs are $1/I_\mu = \sigma^2$ and $1/I_{\sigma^2} = 2\sigma^4/n$.

$\text{Var}(\hat{\mu}) = \sigma^2/n$ which equals the CRLB so is MVUE.

$\text{Var}(\hat{\sigma}^2) = 2(n-1)\sigma^4/n^2$ is less than the CRLB. But $\hat{\sigma}^2$ is biased.

The sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ is unbiased and has variance $2\sigma^4/(n-1)$ which is larger than the CRLB.

Question Is S^2 a MVUE?

Efficiency

Definition

The (Bahadur) **efficiency** of an estimator $\tilde{\theta}$ is defined as a comparison of the variance of $\tilde{\theta}$ with the CR bound I_{θ}^{-1} . That is

$$\text{Efficiency of } \tilde{\theta} = \frac{I_{\theta}^{-1}}{\text{Var}[\tilde{\theta}]} = \frac{1}{I_{\theta} \text{Var}[\tilde{\theta}]}$$

The **asymptotic efficiency** is the limit as $n \rightarrow \infty$.

There are similar definitions for the relative efficiency of two estimators.

Asymptotic normality of MLE

Revision from Part A Statistics As the sample size $n \rightarrow \infty$, the MLE

$$\hat{\theta} \approx N(\theta, I_{\theta}^{-1}).$$

This is a powerful and general result. Assuming the usual regularity conditions hold then it tells us that the MLE has the following properties

- 1 it is asymptotically unbiased
- 2 it is asymptotically efficient i.e. it attains the CRLB asymptotically.
- 3 it has a normal distribution asymptotically.

Extensions to the Cramér-Rao inequality

1. If $\hat{\theta}$ is an estimator with bias $b(\theta) = \text{bias}(\hat{\theta})$, then

$$\text{Var}[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 I_{\theta}^{-1}$$

2. If $\hat{g}(x)$ is an unbiased estimator for $g(\theta)$, then

$$\text{Var}[\hat{g}(X)] \geq \left(\frac{\partial g}{\partial \theta}\right)^2 I_{\theta}^{-1}.$$

Proof Begin with $\mathbb{E}_{\theta}(\hat{\theta}(X)) = \theta + b(\theta)$ (in 1.) and $\mathbb{E}_{\theta}(\hat{g}(X)) = g(\theta)$ (in 2.). Differentiate both sides and proceed as above to find $\text{Cov}[U, V] = (1 + \partial b / \partial \theta)$ (in 1.) and $\text{Cov}[U, V] = \partial g / \partial \theta$ (in 2., with $U = \hat{g}$). The bound is against $\text{Cov}[U, V]^2$ which leads to the results above.

Fisher Information for a d -dimensional parameter

Information matrix:

$$I_{ij} = \mathbb{E} \left[\frac{\partial \ell}{\partial \theta_i} \frac{\partial \ell}{\partial \theta_j} \right] = -\mathbb{E} \left[\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]$$

under regularity conditions. The CR inequality is

$$\text{Var}(\hat{\theta}_i) \geq [I^{-1}]_{ii}, \quad i = 1, \dots, d.$$

Exercise: verify that we have already proved $\text{Var}(\hat{\theta}_i) \geq [I_{ii}]^{-1}$. Note that $[I^{-1}]_{ii} \geq [I_{ii}]^{-1}$ (GJJ) so bound above is stronger.

Exercise For an Exponential family in canonical form,

$$I_{ij} = -\frac{\partial^2}{\partial \phi_i \partial \phi_j} nD(\phi).$$