

Consistency and Efficiency of Estimators

March 6, 2025

The bias and unbiased estimators

In the previous lecture, we have seen an approach to estimation that is based on the likelihood of observed results. Next, we study general theory of estimation that is used to compare between different estimators and to decide on the most efficient one.

- Suppose that we are going to observe a value of a random vector \mathbf{X} . Let \mathcal{X} denote the set of possible values \mathbf{X} can take and, for $\mathbf{x} \in \mathcal{X}$, let $f(\mathbf{x}|\theta)$ denote the probability (or density) that \mathbf{X} takes the value \mathbf{x} where the parameter θ is some unknown element of the set Θ .
- An estimator $\hat{\theta}$ is a procedure that for each possible value $\mathbf{x} \in \mathcal{X}$ specifies which element of Θ we should report as an estimate of θ . When we observe $\mathbf{X} = \mathbf{x}$, we quote $\hat{\theta}(\mathbf{x})$ as our estimate of θ . Thus $\hat{\theta}$ is a function of the random variable \mathbf{X} . Sometimes we write $\hat{\theta}(\mathbf{X})$ to emphasise this point.

Bias and unbiased estimators

To evaluate the usefulness of an estimator $\hat{\theta} = \hat{\theta}(\mathbf{x})$ of θ , examine the properties of the random variable $\hat{\theta} = \hat{\theta}(\mathbf{X})$.

- Given an estimator $\hat{\theta}$, we can calculate its expected value for each possible value of $\theta \in \Theta$.
- An estimator is said to be **unbiased** if this expected value is equal to θ .
- If an estimator is unbiased then, by the Law of Large Numbers, we can conclude that by repeating the experiment an infinite number of times (never happens, not even twice:) with θ fixed and calculate the value of the estimator each time then the average of the estimator values will be exactly equal to θ .

Definition (Unbiased estimators)

An estimator $\hat{\theta} = \hat{\theta}(\mathbf{X})$ is said to be unbiased for a parameter θ if it equals θ in expectation

$$E_{\theta}[\hat{\theta}(\mathbf{X})] = E_{\theta}(\hat{\theta}) = \theta.$$

Intuitively, an unbiased estimator is ‘right on target’. In general, $Bias_{\theta}(\hat{\theta}) = E_{\theta}(\hat{\theta}) - \theta$.

Examples

The most well-known estimators are the sample mean and the sample variance

$$\bar{X} = \sum_{i=1}^n X_i/n, \quad S^2 = \frac{n}{n-1} \overline{(X - \bar{X})^2} = \frac{n}{n-1} (\bar{X^2} - \bar{X}^2)$$

The strange factor $\frac{n}{n-1}$ is to force the unbiasedness of S^2 (Why?).

- Note that even if $\hat{\theta}$ is an unbiased estimator of θ , $g(\hat{\theta})$ will generally not be an unbiased estimator of $g(\theta)$ unless g is linear or affine.
- This limits the importance of the notion of unbiasedness. It might be at least as important that an estimator is accurate in the sense that its distribution is highly concentrated around θ .
- For an arbitrary distribution, the estimator S^2 is an unbiased estimator of the variance of this distribution.
- Consider the estimator S^2 of variance σ^2 in the case of the normal distribution. Although S^2 is an unbiased estimator of σ^2 , S is not an unbiased estimator of σ .

Homework Find the most explicit form of the bias of S as an estimator of σ . Can you simply argue if the bias is positive or negative?

Mean Square Error

Definition (Mean squared error)

The mean squared error of the estimator $\hat{\theta}$ is defined as $MSE_{\theta}(\hat{\theta}) = E_{\theta}(\hat{\theta} - \theta)^2$. Given the same set of data, $\hat{\theta}_1$ is “better” than $\hat{\theta}_2$ if $MSE_{\theta}(\hat{\theta}_1) \leq MSE_{\theta}(\hat{\theta}_2)$ (uniformly better if true for all θ). \square

Lemma (The MSE variance-bias tradeoff)

The MSE decomposes as $MSE_{\theta}(\hat{\theta}) = Var_{\theta}(\hat{\theta}) + Bias_{\theta}(\hat{\theta})^2$.

Proof.

$$\begin{aligned}
 MSE(\hat{\theta}) &= E(\hat{\theta} - \theta)^2 = E\{ [\hat{\theta} - E(\hat{\theta})] + [E(\hat{\theta}) - \theta] \}^2 = E[\hat{\theta} - E(\hat{\theta})]^2 + E[E(\hat{\theta}) - \theta]^2 \\
 &\quad + 2 \underbrace{E\{ [\hat{\theta} - E(\hat{\theta})][E(\hat{\theta}) - \theta] \}}_{=0} \\
 &= E[\hat{\theta} - E(\hat{\theta})]^2 + E[E(\hat{\theta}) - \theta]^2 = Var(\hat{\theta}) + \underbrace{[E(\hat{\theta}) - \theta]^2}_{Bias(\hat{\theta})^2}.
 \end{aligned}$$

NOTE: This lemma implies that the mean squared error of an unbiased estimator is equal to the variance of the estimator.

Example: Standard deviation estimation

- For the estimation of σ^2 in the two-parameter normal distribution, we had two explicit estimators: sample variance and $n - 1$ divisor sample variance.
- One can argue that the unbiased estimator is better, since it has the smaller MSE. Why?
- The estimation of σ can be also performed by the square roots of these two estimators. However the MSE of these estimators are not derivable directly from the MSE of the variance estimators.

Homework S is the biased estimator of σ . Argue that it is not the MLE by explicitly identifying the MLE. The MLE is also biased. Investigate the bias and MSE of these two estimators. Summarize conclusions.

The minimal variance linear mean estimator

Let X_1, \dots, X_n be independent random variables with means $\mathbb{E}(X_i) = \mu$ and variances $\text{Var}(X_i) = \sigma_i^2$. Consider pooling the estimators of μ into a common estimator using the linear combination

$\hat{\mu} = w_1 X_1 + w_2 X_2 + \dots + w_n X_n$. We will see that the following is true

- (i) The estimator $\hat{\mu}$ is unbiased if and only if $\sum w_i = 1$.
- (ii) The estimator $\hat{\mu}$ has minimum variance among this class of estimators when the weights are inversely proportional to the variances σ_i^2 .
- (iii) The variance of $\hat{\mu}$ for optimal weights w_i is $\text{Var}(\hat{\mu}) = 1 / \sum_i \sigma_i^{-2}$.

Minimum-Variance Unbiased Estimation

Getting a small MSE often involves a tradeoff between variance and bias. For unbiased estimators, the MSE obviously equals the variance, $MSE(\hat{\theta}) = Var(\hat{\theta})$, so no tradeoff can be made. One approach is to restrict ourselves to the subclass of estimators that are *unbiased* and *minimum variance*.

Definition (Minimum-variance unbiased estimator)

If an unbiased estimator of $g(\theta)$ has minimum variance (for all possible values of θ) among all unbiased estimators of $g(\theta)$ it is called a uniformly minimum variance unbiased estimator (UMVUE or MVUE).

We will develop a method of finding the MVUE when it exists. When such an estimator does not exist we will be able to find a lower bound for the variance of an unbiased estimator in the class of unbiased estimators, and compare the variance of our unbiased estimator with this lower bound.

Assumptions

- We suppose throughout that we have a regular parametric model and further that Θ is an open subset of the line.
- The support set of the distributions does not depend on $\theta \in \Theta$.
- $\partial \log p(x, \theta) / \partial \theta$ exists.
- Then, for a statistics T , we have the following change of integration with derivation

$$\frac{\partial}{\partial \theta} E_\theta(T) = \frac{\partial}{\partial \theta} \int_{\mathcal{X}} T(x)p(x|\theta) dx = \int_{\mathcal{X}} \frac{\partial}{\partial \theta} T(x)p(x|\theta) dx$$

In our discussion, we will work under these assumptions. We do not discuss in any detail when these assumptions are satisfied. It should be however remembered that

- these assumptions are satisfied by exponential families in the one- and multiparameter cases,
- they are frequently valid both for the discrete and continuous distributions.

Homework Discuss the issue of changing the order of differentiation and integration by providing examples when this can and cannot be done. Provide some general conditions that are sufficient for this to hold.

The score function

Definition (Score function)

For the (possibly vector valued) observation $X = x$ to be informative about θ , the density must vary with θ . If $p(x|\theta)$ is smooth and differentiable, then for finding MLE we have used zeros of the score function

$$S(\theta) = S(\theta|x) = \frac{\partial}{\partial \theta} \log p(x|\theta) \equiv \frac{\partial p(x|\theta)/\partial \theta}{p(x|\theta)}.$$

If differentiation wrt θ and integration wrt x can be interchanged, as we assume, we have for X distributed according to $p(x|\theta)$:

$$E_{\theta} \{ S(\theta|X) \} = \int \frac{\partial p(x|\theta)/\partial \theta}{p(x|\theta)} p(x|\theta) dx = \int \partial p(x|\theta)/\partial \theta dx = \frac{\partial}{\partial \theta} \left\{ \int p(x|\theta) dx \right\} = \frac{\partial}{\partial \theta} 1 = 0.$$

- Thus the score function has expectation zero.
- The score function $S(\theta|X)$, for a fixed θ , is a random variable (but it is not a statistic, why?).
- We often drop explicit dependence on X from the notation by simply writing $S(\theta)$.
- The negative derivative of the score function $I_{obs}(\theta) = -\partial S(\theta)/\partial \theta$ measures how concave down is the likelihood around value θ .

The Fisher Information

Definition (Fisher information)

The Fisher information is defined as the average value of the minus derivative of the score function

$$I(\theta) \equiv -E_\theta \left(\frac{\partial}{\partial \theta} S(\theta) \right) = E_\theta I_{obs}(\theta).$$

The negative derivative of the score function $I_{obs}(\theta)$, which is a random variable dependent on X , is referred to as empirical or observed information about θ . □

Lemma

$$I(\theta) = \text{Var}_\theta S(\theta)$$

We note that $\text{Var}_\theta S = E_\theta S^2$ and

$$\frac{\partial S}{\partial \theta} = \frac{\partial^2 \log p}{\partial \theta^2} = \frac{\partial}{\partial \theta} \left[\frac{1}{p} \frac{\partial p}{\partial \theta} \right] = -\frac{1}{p^2} \left[\frac{\partial p}{\partial \theta} \right]^2 + \frac{1}{p} \frac{\partial^2 p}{\partial \theta^2} = - \left[\frac{\partial \log p}{\partial \theta} \right]^2 + \frac{1}{p} \frac{\partial^2 p}{\partial \theta^2} = -S^2 + \frac{1}{p} \frac{\partial^2 p}{\partial \theta^2}$$

If integration and differentiation can be interchanged

$$E_\theta \left[\frac{1}{p} \frac{\partial^2 p}{\partial \theta^2} \right] = \int_X \frac{\partial^2 p}{\partial \theta^2} dx = \frac{\partial^2}{\partial \theta^2} \int_X p dx = \frac{\partial^2}{\partial \theta^2} 1 = 0.$$

The Cramér-Rao lower bound (CRLB)

Theorem (Information Inequality)

Let $\hat{\theta}$ be an unbiased estimator of θ . Then $\text{Var}_\theta \hat{\theta} \geq 1/I(\theta)$.

Proof.

Unbiasedness means that $\mathbb{E}_\theta(\hat{\theta}) = \int \hat{\theta}(x)p(x|\theta)dx = \theta$. Assume we can differentiate wrt θ under the integral, then $\int \frac{\partial}{\partial \theta} \left\{ \hat{\theta}(x)p(x|\theta) \right\} dx = \int \hat{\theta}(x) \frac{\partial}{\partial \theta} p(x|\theta) dx = 1$ since the estimator $\hat{\theta}(x)$ can't depend on θ .

Now, $\frac{\partial p}{\partial \theta} = p \frac{\partial}{\partial \theta} (\log p)$, so that

$$\int \hat{\theta}(x)p \frac{\partial}{\partial \theta} (\log p) dx = E_\theta \left[\hat{\theta}(x) \frac{\partial}{\partial \theta} (\log p) \right] = E_\theta(\hat{\theta}S) = 1.$$

We already know that the score function has expectation zero, $E_\theta(S) = 0$. Consequently $\text{Cov}_\theta(\hat{\theta}, S) = E_\theta(\hat{\theta}S) - E_\theta(\hat{\theta})E_\theta(S) = E_\theta(\hat{\theta}S) = 1$. By Schwartz's inequality

$$1 = \text{Cov}_\theta(\hat{\theta}, S)^2 \leq \text{Var}_\theta(\hat{\theta}) \text{Var}_\theta(S)$$

This implies the Rao-Cramér inequality. □

Comments

- Why ‘**information**’? – Variance measures the lack of knowledge. The reciprocal of the variance could be defined as the amount of information carried by the (possibly vector valued) random observation X about θ .
- We require that the ranges of the integrals (supports of the densities) do not depend on θ . That is, the range of x , here the support of $p(x|\theta)$, cannot depend on θ . This second condition is violated for some density functions, i.e. the CRLB is not valid for the uniform distribution.
- We can have absolute assessment for unbiased estimators by comparing their variances to the CRLB. We can also assess biased estimators. If its variance is lower than CRLB then it can be indeed a very good estimate, although it is biased.
- In the iid case, i.e. $p(x|\theta) = p_1(x_1|\theta) \dots p_1(x_n|\theta)$, then $I(\theta) = nI_1(\theta)$, where the $I_1(\theta)$ is based on $p_1(x|\theta)$.

Example

Consider IID case with

$$p_1(x_i|\mu) = \frac{1}{\mu} \exp\left(-\frac{1}{\mu}x_i\right).$$

and $p(x|\mu) = \left(\frac{1}{\mu}\right)^n \exp\left(-\frac{1}{\mu} \sum_{i=1}^n x_i\right)$. The score function, which is the partial derivative of $\log f$ wrt the unknown parameter μ , is

$$S(\mu) = \frac{\partial}{\partial \mu} \log f = -\frac{n}{\mu} + \frac{1}{\mu^2} \sum_{i=1}^n x_i.$$

For $X \sim \text{Exp}(1/\mu)$, we have $E_\mu(X) = \mu$ implying

$$I(\mu) = -E_\mu \left\{ \frac{\partial}{\partial \mu} \left(-\frac{n}{\mu} + \frac{1}{\mu^2} \sum_{i=1}^n X_i \right) \right\} = -\frac{n}{\mu^2} + \frac{2}{\mu^3} E_\mu \left\{ \sum_{i=1}^n X_i \right\} = -\frac{n}{\mu^2} + \frac{2n\mu}{\mu^3}$$

Hence $CRLB = \frac{\mu^2}{n}$.

Let us propose $\hat{\mu} = \bar{X}$ as an unbiased estimator of μ . For $X \sim \text{Exp}(1/\mu)$, we have $E(X) = \mu = \sqrt{\text{Var}(X)}$, implying the unbiased estimator $\hat{\mu} = \bar{x}$ achieves its CRLB

$$\text{Var}(\hat{\mu}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{n\mu^2}{n^2} = \frac{\mu^2}{n}.$$

Example: Uniform distribution

- Let U be a random variable uniformly distributed over $(0, \theta)$.
- Then $\hat{\theta} = U$ is the MLE of θ .
- One can argue that $\hat{\theta}_1 = 2U$ is the unbiased estimator. Why?

Homework Discuss difficulties with the score function and the Fisher information in this case. Of these two estimators which has smaller MSE?

- Let U_1, \dots, U_n be a sample from the uniform distribution over $(0, \theta)$.
- Then the MLE of θ is

$$\hat{\theta} = \max_{i=1, \dots, n} U_i$$

Homework Argue that the MLE is not unbiased. Propose an unbiased estimator based on the MLE. Discuss the MSE and Fisher information for these two estimators.

Efficiency

Definition (Efficiency)

Define the efficiency of an unbiased estimator $\hat{\theta}$ as

$$\text{eff}(\hat{\theta}) = \frac{\text{CRLB}}{\text{Var}(\hat{\theta})},$$

where $\text{CRLB} = 1/I(\theta)$. Clearly $0 < \text{eff}(\hat{\theta}) \leq 1$. An unbiased estimator $\hat{\theta}$ is said to be efficient if $\text{eff}(\hat{\theta}) = 1$. □

Definition (Asymptotic efficiency)

The *asymptotic efficiency* of an unbiased estimator $\hat{\theta}$ is the limit of the efficiency as $n \rightarrow \infty$. An unbiased estimator $\hat{\theta}$ is said to be asymptotically efficient if its asymptotic efficiency is equal to 1. □

Examples

The following are examples where the efficiency or asymptotical efficiency can be easily demonstrated

- the MLE $\hat{\theta} = r/n$ for the binomial distribution that was considered.
- the MLE for the Poisson distribution is 100% efficient.
- the MLE $\hat{\theta}$ for the exponential distribution with parameter θ is $\tilde{\theta}$ asymptotically efficient.
- the MLE estimator of variance in the normal distribution is asymptotically efficient.

Homework *Find a formal argument for all these statements.*

Multiparameter extension

The information matrix is the covariance of the score vector
 $\mathbf{S} = \partial \log p(X|\theta) / \partial \theta$ function

$$\mathbf{I}(\theta) = \text{Cov}_{\theta}(\mathbf{S}) = \left[E_{\theta} \left(\frac{\partial}{\partial \theta_i} \log p(X|\theta) \frac{\partial}{\partial \theta_j} \log p(X|\theta) \right) \right]_{i,j=1}^d$$

Theorem

Suppose that $\hat{\theta}$ is an unbiased estimator of θ and $\mathbf{I}(\theta)$ is non-singular.

$$\text{Cov}_{\theta} \hat{\theta} - \mathbf{I}(\theta)^{-1}$$

is a positive definite matrix.

Canonical exponential family

Consider the canonical form of the exponential family

$$p(x|\theta) = \exp \left(\sum_{j=1}^d T_j(x)\theta_j - A(\theta) \right) h(x)$$

Then

$$\mathbf{I}(\theta) = \text{Cov}_{\theta} \mathbf{T}(X) = \ddot{A}(\theta).$$

Asymptotics of estimators

We have already defined asymptotic efficiency of estimators, where asymptotics is with respect to the sample size converging to infinity. There are other properties that are desirable.

- Asymptotically unbiased $Bias_{\theta,n} \rightarrow 0$.
- Consistency $\hat{\theta}_n \rightarrow \theta$, where convergence is either in probability or with probability one (strong consistency).
- Asymptotic normality when $a_n(\hat{\theta} - \theta)$ converges to a standard normal distribution, for properly chosen non-random a_n (can be a matrix in the multiparameter case).
- The asymptotic theory will be discussed later in the course. Here we only consider the consistency condition.

Consistency

Definition

An estimator $\hat{\theta}_n$ of a parameter θ is called consistent if it is convergent in probability to θ as $n \rightarrow \infty$, i.e. for each $\theta \in \Theta$

$$\lim_{n \rightarrow \infty} P_\theta(|\hat{\theta}_n - \theta| > \epsilon) = 0.$$

A sample mean \bar{X}_n is a consistent estimator of the mean $\mu = E_\theta X$ since by the Law of Large Number it is always

$$\bar{X}_n \rightarrow EX.$$

An estimator is called uniformly consistent if

$$\lim_{n \rightarrow \infty} \sup_{\theta \in \Theta} P_\theta(|\hat{\theta}_n - \theta| > \epsilon) = 0.$$

Consistency through Chebyshev's inequality

- Recall that for an unbiased estimator $\hat{\theta}_n$:

$$P_\theta(|\hat{\theta}_n - \theta| \geq \epsilon) \leq \text{Var}_\theta(\hat{\theta}_n)/\epsilon^2.$$

- Thus if one shows that for each θ we have that $\text{Var}_\theta(\hat{\theta}_n)$ converges to zero, then it yields consistency. If $\sup_{\theta \in \Theta} \text{Var}_\theta(\hat{\theta}_n)$ converges to zero, then it yields uniform consistency.
- Any efficient estimator based on an IID sample is consistent since

$$\text{Var}_\theta(\hat{\theta}_n) = \frac{1}{nl_1(\theta)} \rightarrow 0.$$

- Uniform consistency does not need to hold.

Consistency for the exponential families

Theorem

The MLE for an exponential family is consistent.

- One of the reasons why the MLE for the exponential families is a good estimator.
- There will be more reasons given later.

Homework Consider the exponential distribution with the intensity λ and the MLE for λ . Is this estimator uniformly consistent?