

## Lecture 8: Properties of Maximum Likelihood Estimation (MLE)

(LaTeX prepared by Haiguang Wen)  
April 27, 2015

This lecture note is based on ECE 645(Spring 2015) by Prof. Stanley H. Chan in the School of Electrical and Computer Engineering at Purdue University.

### 1 Efficiency of MLE

*Maximum Likelihood Estimation* (MLE) is a widely used statistical estimation method. In this lecture, we will study its properties: efficiency, consistency and asymptotic normality.

MLE is a method for estimating parameters of a statistical model. Given the distribution of a statistical model  $f(y; \theta)$  with unknown deterministic parameter  $\theta$ , MLE is to estimate the parameter  $\theta$  by maximizing the probability  $f(y; \theta)$  with observations  $y$ .

$$\hat{\theta}(y) = \arg \min_{\theta} f(y; \theta) \quad (1)$$

Please see the previous lecture note **Lecture 7** for details.

#### 1.1 Cramér–Rao Lower Bound (CRLB)

*Cramér–Rao Lower Bound* (CRLB) is introduced in **Lecture 7**. Briefly, CRLB describes a lower bound on the variance of estimators of the deterministic parameter  $\theta$ . That is

$$\text{Var} \left( \hat{\theta}(Y) \right) \geq \frac{\left( \frac{\partial}{\partial \theta} \mathbb{E}[\hat{\theta}(Y)] \right)^2}{I(\theta)}, \quad (2)$$

where  $I(\theta)$  is the *Fisher information* that measures the information carried by the observable random variable  $Y$  about the unknown parameter  $\theta$ . For unbiased estimator  $\hat{\theta}(Y)$ , Equation 2 can be simplified as

$$\text{Var} \left( \hat{\theta}(Y) \right) \geq \frac{1}{I(\theta)}, \quad (3)$$

which means the variance of any unbiased estimator is at least as the inverse of the Fisher information.

#### 1.2 Efficient Estimator

From section 1.1, we know that the variance of estimator  $\hat{\theta}(y)$  cannot be lower than the CRLB. So any estimator whose variance is equal to the lower bound is considered as an efficient estimator.

<b>Definition 1.</b> EFFICIENT ESTIMATOR
--

An estimator $\hat{\theta}(y)$ is efficient if it achieves equality in CRLB.
--

---

**Example 1.**

**Question:**  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_n\}$  are i.i.d. Gaussian random variables with distribution  $N(\theta, \sigma^2)$ . Determine the maximum likelihood estimator of  $\theta$ . Is the estimator efficient?

**Solution:** Let  $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$  be the observation, then

$$\begin{aligned} f(\mathbf{y}; \theta) &= \prod_{k=1}^n f(y_k; \theta) \\ &= \prod_{k=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_k - \theta)^2}{2\sigma^2}\right\} \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left\{-\frac{\sum_{k=1}^n (y_k - \theta)^2}{2\sigma^2}\right\}. \end{aligned}$$

Take the log of both sides of the above equation, we have

$$\log f(\mathbf{y}; \theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{\sum_{k=1}^n (y_k - \theta)^2}{2\sigma^2}.$$

Since  $\log f(\mathbf{y}; \theta)$  is a quadratic concave function of  $\theta$ , we can obtain the MLE by solving the following equation.

$$\frac{\partial \log f(\mathbf{y}; \theta)}{\partial \theta} = \frac{2 \sum_{k=1}^n (y_k - \theta)}{2\sigma^2} = 0.$$

Therefore, the MLE is  $\hat{\theta}_{MLE}(\mathbf{y}) = \frac{1}{n} \sum_{k=1}^n y_k$ .

Now let us check whether the estimator is efficient or not. It is easy to check that the MLE is an unbiased estimator ( $\mathbb{E}[\hat{\theta}_{MLE}(\mathbf{y})] = \theta$ ). To determine the CRLB, we need to calculate the Fisher information of the model.

$$I(\theta) = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f(\mathbf{y}; \theta) \right] = \frac{n}{\sigma^2} \quad (4)$$

According to Equation 3, we have

$$\text{Var}(\hat{\theta}_{MLE}(\mathbf{Y})) \geq \frac{1}{I(\theta)} = \frac{\sigma^2}{n}. \quad (5)$$

And the variance of the MLE is

$$\text{Var}(\hat{\theta}_{MLE}(\mathbf{Y})) = \text{Var}\left(\frac{1}{n} \sum_{k=1}^n Y_k\right) = \frac{\sigma^2}{n}. \quad (6)$$

So CRLB equality is achieved, thus the MLE is efficient.

### 1.3 Minimum Variance Unbiased Estimator (MVUE)

Recall that a *Minimum Variance Unbiased Estimator* (MVUE) is an unbiased estimator whose variance is lower than any other unbiased estimator for all possible values of parameter  $\theta$ . That is

$$\text{Var}(\hat{\theta}_{MVUE}(\mathbf{Y})) \leq \text{Var}(\hat{\theta}(\mathbf{Y})) \quad (7)$$

for any unbiased  $\hat{\theta}(\mathbf{Y})$  of any  $\theta$ .

**Proposition 1.** UNBIASED AND EFFICIENT ESTIMATORS  
If an estimator  $\hat{\theta}(y)$  is unbiased and efficient, then it must be MVUE.

**Proof.**

Since  $\hat{\theta}(y)$  is efficient, according to CRLB, we have

$$\text{Var}(\hat{\theta}(Y)) \leq \text{Var}(\tilde{\theta}(Y)) \quad (8)$$

for any  $\tilde{\theta}(Y)$ . Therefore,  $\hat{\theta}(Y)$  must be minimum variance (MV). Since  $\tilde{\theta}(Y)$  is also unbiased, it is a MVUE.

□

**Remark:** The converse of the proposition is not true in general. That is, MVUE does NOT need to be efficient. Here is a counter example.

---

**Example 2. COUNTER EXAMPLE**

Suppose that  $\mathbf{y} = \{Y_1, Y_2, \dots, Y_n\}$  are i.i.d. exponential random variables with unknown mean  $\frac{1}{\theta}$ . Find the MLE and MVUE of  $\theta$ . Are these estimators efficient?

**Solution:** Let  $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$  be the observation, then

$$\begin{aligned} f(\mathbf{y}; \theta) &= \prod_{k=1}^n f(y_k; \theta) \\ &= \prod_{k=1}^n \theta \exp\{-\theta y_k\} \\ &= \theta^n \exp\left\{-\theta \sum_{k=1}^n y_k\right\}. \end{aligned} \quad (9)$$

Take the log of both sides of the above equation, we have

$$\log f(\mathbf{y}; \theta) = n \log(\theta) - \theta \sum_{k=1}^n y_k.$$

Since  $\log f(\mathbf{y}; \theta)$  is a concave function of  $\theta$ , we can obtain the MLE by solving the following equation.

$$\frac{\partial \log f(\mathbf{y}; \theta)}{\partial \theta} = \frac{n}{\theta} - \sum_{k=1}^n y_k = 0$$

So the MLE is

$$\hat{\theta}_{MLE}(\mathbf{y}) = \frac{n}{\sum_{k=1}^n y_k}. \quad (10)$$

To calculate the CRLB, we need to calculate  $\mathbb{E}[\hat{\theta}_{MLE}(\mathbf{Y})]$  and  $\text{Var}(\hat{\theta}_{MLE}(\mathbf{Y}))$ . Let  $T(\mathbf{y}) = \sum_{k=1}^n y_k$ , then by moment generating function, we can show that the distribution of  $T(\mathbf{y})$  is the Erlange distribution:

$$f_T(t) = \frac{\theta^n t^{n-1}}{(n-1)!} e^{-\theta t}. \quad (11)$$

So we have

$$\begin{aligned} \mathbb{E}[\hat{\theta}_{MLE}(T(\mathbf{Y}))] &= \int_0^\infty \frac{n}{t} \frac{\theta^n t^{n-1}}{(n-1)!} e^{-\theta t} dt \\ &= \frac{n\theta}{n-1} \int_0^\infty \frac{(\theta t)^{n-2}}{(n-2)!} e^{-\theta t} d\theta t \\ &= \frac{n}{n-1} \theta. \end{aligned} \quad (12)$$

Therefore the MLE is a biased estimator of  $\theta$ .

Similarly, we can calculate the variance of MLE as follows.

$$\begin{aligned}\text{Var}\left(\hat{\theta}_{MLE}(T(\mathbf{Y}))\right) &= \mathbb{E}\left[\hat{\theta}_{MLE}^2(T(\mathbf{Y}))\right] - \mathbb{E}\left[\hat{\theta}_{MLE}(T(\mathbf{Y}))\right]^2 \\ &= \frac{\theta^2 n^2}{(n-1)^2(n-2)}\end{aligned}$$

The Fisher information is

$$I(\theta) = -\frac{\partial^2}{\partial \theta^2} \log f(\mathbf{y}|\theta) = \frac{n}{\theta^2}.$$

So the CRLB is

$$\begin{aligned}\text{Var}\left(\hat{\theta}_{MLE}(T(\mathbf{Y}))\right) &\geq \frac{\left(\frac{\partial}{\partial \theta} \mathbb{E}\left[\hat{\theta}_{MLE}(T(\mathbf{Y}))\right]\right)^2}{I(\theta)} \\ &= \frac{n^2}{(n-1)^2} \bigg/ \frac{n}{\theta^2} \\ &= \frac{n}{(n-1)^2} \theta^2.\end{aligned}$$

The CRLB equality does **NOT** hold, so  $\hat{\theta}_{MLE}$  is not efficient.

The distribution in Equation 9 belongs to exponential family and  $T(\mathbf{y}) = \sum_{k=1}^n y_k$  is a complete sufficient statistic. So the MLE can be expressed as  $\hat{\theta}_{MLE}(T(\mathbf{y})) = \frac{n}{T(\mathbf{y})}$ , which is a function of  $T(\mathbf{y})$ . However, the MLE is a biased estimator (Equation 12). But we can construct an unbiased estimator based on the MLE. That is

$$\begin{aligned}\tilde{\theta}(T(\mathbf{y})) &= \frac{n-1}{n} \hat{\theta}_{MLE}(T(\mathbf{y})) \\ &= \frac{n-1}{T(\mathbf{y})}.\end{aligned}$$

It is easy to check  $\mathbb{E}\left[\tilde{\theta}(T(\mathbf{Y}))\right] = \mathbb{E}\left[\frac{n-1}{n} \hat{\theta}_{MLE}(T(\mathbf{Y}))\right] = \frac{n-1}{n} \frac{n}{n-1} \theta = \theta$ . Since  $\tilde{\theta}(T(\mathbf{y}))$  is an unbiased estimator and it is a function of complete sufficient statistic,  $\tilde{\theta}(T(\mathbf{y}))$  is MVUE. So

$$\hat{\theta}_{MVUE}(T(\mathbf{y})) = \frac{n-1}{T(\mathbf{y})}. \quad (13)$$

The variance of MVUE is

$$\begin{aligned}\text{Var}\left(\hat{\theta}_{MVUE}(T(\mathbf{Y}))\right) &= \text{Var}\left(\frac{n-1}{n} \hat{\theta}_{MLE}(T(\mathbf{Y}))\right) \\ &= \frac{(n-1)^2}{n^2} \frac{n^2 \theta^2}{(n-1)^2(n-2)} \\ &= \frac{\theta^2}{n-2}.\end{aligned}$$

So the CRLB is

$$\text{Var}\left(\hat{\theta}_{MVUE}(T(\mathbf{Y}))\right) \geq \frac{1}{I(\theta)} = \frac{\theta^2}{n}.$$

Therefore, the MVUE is **NOT** an efficient estimator.

## 2 Consistency of MLE

### Definition 2. CONSISTENCY

Let  $\{Y_1, \dots, Y_n\}$  be a sequence of observations. Let  $\hat{\theta}_n$  be the estimator using  $\{Y_1, \dots, Y_n\}$ . We say that  $\hat{\theta}_n$  is consistent if  $\hat{\theta}_n \xrightarrow{P} \theta$ , i.e.,

$$\mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0, \text{ as } n \rightarrow \infty \quad (14)$$

**Remark:** A sufficient condition to have Equation 14 is that

$$\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] \rightarrow 0, \text{ as } n \rightarrow \infty. \quad (15)$$

### Proof.

According to Chebyshev's inequality, we have

$$\mathbb{P}(|\hat{\theta}_n - \theta| \geq \varepsilon) \leq \frac{\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right]}{\varepsilon^2} \quad (16)$$

Since  $\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] \rightarrow 0$ , then we have

$$0 \leq \mathbb{P}(|\hat{\theta}_n - \theta| \geq \varepsilon) \leq \frac{\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right]}{\varepsilon^2} \rightarrow 0.$$

Therefore,  $\mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0$ , as  $n \rightarrow \infty$ . □

### Example 3.

$\{Y_1, Y_2, \dots, Y_n\}$  are i.i.d. Gaussian random variables with distribution  $N(\theta, \sigma^2)$ . Is the MLE using  $\{Y_1, Y_2, \dots, Y_n\}$  consistent?

### Solution:

From Example 1., we know that the MLE is

$$\hat{\theta}_n = \frac{1}{n} \sum_{k=1}^n Y_k.$$

Since

$$\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] = \text{Var}(\hat{\theta}_n) = \frac{\sigma^2}{n}, \text{ (see Equation 6),}$$

so  $\mathbb{E} \left[ \left( \hat{\theta}_n - \theta \right)^2 \right] \rightarrow 0$ . Therefore  $\hat{\theta}_n \xrightarrow{P} \theta$ , i.e.  $\hat{\theta}_n$  is consistent.

In fact, the result of the example above it holds for any distribution. The following proposition states this result:

### Proposition 2.

(MLE of i.i.d observation is consistent) Let  $\{Y_1, \dots, Y_n\}$  be a sequence of i.i.d. observations where

$$Y_k \stackrel{iid}{\sim} f_\theta(y).$$

Then the MLE of  $\theta$  is consistent.

**Proof.**

(This proof is partially correct. See Levy Chapter 4.5 for complete discussion.)  
The MLE of  $\theta$  is

$$\begin{aligned}\hat{\theta}_n &= \arg \max_{\theta} \prod_{k=1}^n f_{\theta}(y_k) \\ &= \arg \max_{\theta} \log \left( \prod_{k=1}^n f_{\theta}(y_k) \right) \\ &= \arg \max_{\theta} \frac{1}{n} \sum_{k=1}^n \log f_{\theta}(y_k) \\ &= \arg \max_{\theta} \varphi_n(\theta),\end{aligned}$$

where  $\varphi_n(\theta) = \frac{1}{n} \sum_{k=1}^n \log f_{\theta}(y_k)$ . Let  $\ell_{\hat{\theta}}(y_k) = \log \frac{f_{\theta}(y_k)}{f_{\hat{\theta}}(y_k)}$ , then we have

$$\begin{aligned}\mathbb{E}_{\hat{\theta}} [\ell_{\hat{\theta}}(Y_k)] &\stackrel{def}{=} \int \log \frac{f_{\theta}(y_k)}{f_{\hat{\theta}}(y_k)} f_{\hat{\theta}}(y_k) dy_k \\ &= D(f_{\theta} \| f_{\hat{\theta}}).\end{aligned}$$

According to the weak law of large numbers (WLLN), we have

$$\frac{1}{n} \sum_{k=1}^n \ell_{\hat{\theta}}(y_k) \xrightarrow{p} D(f_{\theta} \| f_{\hat{\theta}}). \quad (17)$$

Since  $\hat{\theta}_n$  is the MLE which maximizes  $\varphi_n(\theta)$ , then

$$\begin{aligned}0 &\geq \varphi_n(\theta) - \varphi_n(\hat{\theta}) \\ &= \frac{1}{n} \sum_{k=1}^n \log f_{\theta}(y_k) - \frac{1}{n} \sum_{k=1}^n \log f_{\hat{\theta}}(y_k) \\ &= \frac{1}{n} \sum_{k=1}^n \log \frac{f_{\theta}(y_k)}{f_{\hat{\theta}}(y_k)} \\ &= \frac{1}{n} \sum_{k=1}^n \ell_{\hat{\theta}}(y_k) \\ &= \frac{1}{n} \sum_{k=1}^n \ell_{\hat{\theta}}(y_k) - D(f_{\theta} \| f_{\hat{\theta}}) + D(f_{\theta} \| f_{\hat{\theta}}).\end{aligned}$$

Therefore,

$$D(f_{\theta} \| f_{\hat{\theta}}) \leq \left| \frac{1}{n} \sum_{k=1}^n \ell_{\hat{\theta}}(y_k) - D(f_{\theta} \| f_{\hat{\theta}}) \right|.$$

By Equation 17, we have

$$0 \leq D(f_{\theta} \| f_{\hat{\theta}}) \leq \left| \frac{1}{n} \sum_{k=1}^n \ell_{\hat{\theta}}(y_k) - D(f_{\theta} \| f_{\hat{\theta}}) \right| \xrightarrow{p} 0.$$

So we must have  $D(f_{\theta} \| f_{\hat{\theta}}) \xrightarrow{p} 0$ , and then  $\hat{\theta} \xrightarrow{p} \theta$ .

□

### 3 Asymptotic Normality of MLE

The previous proposition only asserts that MLE of i.i.d. observations is consistent. However, it provides no information about the distribution of the MLE.

**Proposition 3.** ASYMPTOTIC NORMALITY

Let  $\{Y_1, \dots, Y_n\}$  be a sequence of i.i.d. observations where

$$Y_k \stackrel{iid}{\sim} f_\theta(y)$$

$\hat{\theta}$  is the MLE of  $\theta$ , then

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N\left(0, \frac{1}{I(\theta)}\right).$$

**Proof.**

See Lehmann, “Elements of Large Sample Theory”, Springer, 1999 for proof. □

---

**Example 4.**

$\{Y_1, Y_2, \dots, Y_n\}$  are i.i.d. Gaussian random variables with distribution  $N(\theta, \sigma^2)$ . Find the asymptotic distribution of  $\hat{\theta}_{ML}$

**Solution:**

Similar to Example 2, we can calculate the Fisher information of  $\theta$ ,

$$I(\theta) = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f_\theta(Y) \right] = \frac{1}{\sigma^2}$$

We know that  $\hat{\theta}_{ML} = \frac{1}{n} \sum_{k=1}^n y_k$ . So if  $\hat{\theta}_n = \frac{1}{n} \sum_{k=1}^n y_k$ , then we have

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \sigma^2).$$

---