

# Chapter 9. Properties of Point Estimators and Methods of Estimation

**9.1** Introduction

**9.2** Relative Efficiency

**9.3** Consistency

**9.4** Sufficiency

**9.5** The Rao-Blackwell Theorem and Minimum-Variance Unbiased Estimation

**9.6** The Method of Moments

**9.7** The Method of Maximum Likelihood

## 9.1 Introduction

- Estimator  $\hat{\theta} = \hat{\theta}_n = \hat{\theta}(Y_1, \dots, Y_n)$  for  $\theta$  : a function of  $n$  random samples,  $Y_1, \dots, Y_n$ .
- Sampling distribution of  $\hat{\theta}$  : a probability distribution of  $\hat{\theta}$
- Unbiased estimator,  $\hat{\theta}$  :  $E(\hat{\theta}) - \theta = 0$ .
- Properties of  $\hat{\theta}$  : efficiency, consistency, sufficiency
- Rao-Blackwell theorem : an unbiased estimator with small variance is a function of a sufficient statistic
- Estimation method
  - Minimum-Variance Unbiased Estimation
  - Method of Moments
  - Method of Maximum Likelihood

## 9.2 Relative Efficiency

- We would like to have an estimator with smaller bias and smaller variance : if one can find several unbiased estimators, we want to use an estimator with smaller variance.
- Relative efficiency

(Def 9.1) Suppose  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators for  $\theta$ , with variances,  $V(\hat{\theta}_1)$  and  $V(\hat{\theta}_2)$ , respectively.

Then *relative efficiency* of  $\hat{\theta}_1$  relative to  $\hat{\theta}_2$ , denoted  $\text{eff}(\hat{\theta}_1, \hat{\theta}_2)$ , is defined to be the ratio  $\text{eff}(\hat{\theta}_1, \hat{\theta}_2) = \frac{V(\hat{\theta}_2)}{V(\hat{\theta}_1)}$

- $\text{eff}(\hat{\theta}_1, \hat{\theta}_2) > 1$  :  $V(\hat{\theta}_2) > V(\hat{\theta}_1)$ , and  $\hat{\theta}_1$  is *relatively more efficient* than  $\hat{\theta}_2$ .
- $\text{eff}(\hat{\theta}_1, \hat{\theta}_2) < 1$  :  $V(\hat{\theta}_2) < V(\hat{\theta}_1)$ , and  $\hat{\theta}_2$  is *relatively more efficient* than  $\hat{\theta}_1$ .

(Example) want to estimate the mean of a normal distribution. Let  $\hat{\theta}_{med}$  be the sample median and  $\hat{\theta}_{mean}$  be the sample mean. Then  $\hat{\theta}_{mean}$  is better than  $\hat{\theta}_{med}$ . Why?

$$eff(\hat{\theta}_{med}, \hat{\theta}_{mean}) = \frac{V(\hat{\theta}_{mean})}{V(\hat{\theta}_{med})} = \frac{\sigma^2/n}{(1.2533)^2 \sigma^2/n} = .6366.$$

(Exercise) Let  $Y_1, \dots, Y_n$  denote a random sample from a population with mean  $\mu$  and  $\sigma^2$ . Consider the following three estimates for  $\mu$  :

$$\hat{\mu}_1 = \frac{1}{2}(Y_1 + Y_2), \quad \hat{\mu}_2 = \frac{1}{4}Y_1 + \frac{Y_2 + \dots + Y_{n-1}}{2(n-2)} + \frac{1}{4}Y_n, \quad \hat{\mu}_3 = \bar{Y}.$$

(a) Show that they are unbiased.

(b) Find the efficiency of  $\hat{\mu}_3$  relative to  $\hat{\mu}_2$  and  $\hat{\mu}_1$ , respectively.

## 9.3 Consistency

- Motivation - (Tossing a coin)

A coin is tossed  $n$  times independently, and  $p$  is the (unknown) probability of resulting in heads. Suppose we are interested in the number of heads among  $n$  tosses,  $Y$ .

(Q1) what is the distribution of  $Y$ ?

Since  $p$  is unknown, consider  $\hat{p} = \hat{p}_n = Y/n$ . As  $n$  increases, the amount of information in the sample (here, the quality of  $\hat{p}_n$ ) also increases in the sense that  $\hat{p}_n = Y/n$  should be very close to  $p$  as  $n \rightarrow \infty$ .

(Q2) How one can express “closeness” of  $\hat{p}_n$  to  $p$ ?

Since  $\hat{p}_n = Y/n$  is a statistic, consider the probability that  $|\hat{p}_n - p|$  will be less than some arbitrary positive number as  $n$  increases: what is the value of  $P(|Y/n - p| \leq \epsilon)$  as  $n \rightarrow \infty$ ?

- Consistency and related theorems

(Def 9.2)  $\hat{\theta}_n = \hat{\theta}(Y_1, \dots, Y_n)$  is said to be a *consistent estimator* of  $\theta$  (i.e.,  $\hat{\theta}_n \xrightarrow{P} \theta$ )

if, for any  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| \leq \epsilon) = 1 \text{ or } \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0$$

(note) “ $\hat{\theta}_n$  is a *consistent estimator* of  $\theta$ ” means “ $\hat{\theta}_n$  converges in probability to  $\theta$ ”

(Thm 9.1) An unbiased  $\hat{\theta}_n$  for  $\theta$  is a consistent estimator of  $\theta$  if  $\lim_{n \rightarrow \infty} V(\hat{\theta}_n) = 0$ .

(Example 9.2) Let  $Y_1, \dots, Y_n$  denote a random sample from a distribution with mean  $\mu$  and variance  $\sigma^2 < \infty$ . Show that  $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$  is a consistent estimator of  $\mu$ .

How about consistency of  $\bar{Y}_1 - \bar{Y}_2$  for  $\mu_1 - \mu_2$ ?

(Thm 9.2) Suppose that  $\hat{\theta}_n$  and  $\hat{\theta}'_n$  are consistent estimators of  $\theta$  and  $\theta'$ , respectively. Then,

- (a)  $\hat{\theta}_n + \hat{\theta}'_n$  is a consistent estimator of  $\theta + \theta'$
- (b)  $\hat{\theta}_n \times \hat{\theta}'_n$  is a consistent estimator of  $\theta \times \theta'$
- (c) If  $\theta' \neq 0$ ,  $\hat{\theta}_n / \hat{\theta}'_n$  is a consistent estimator of  $\theta / \theta'$
- (d) If  $g(\cdot)$  is a real-valued function that is continuous at  $\theta$ , then  $g(\hat{\theta}_n)$  is a consistent estimator of  $g(\theta)$ .

(Example 9.3) Let  $Y_1, \dots, Y_n$  denote a random sample from a distribution with finite  $E(Y_i) = \mu$ ,  $E(Y_i^2) = \mu_2'$ , and  $E(Y_i^4) = \mu_4'$ . Show that  $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2$  is a consistent estimator of  $\sigma^2 = V(Y_i)$ .

(Question) Suppose  $Y_1, \dots, Y_n$  is a random sample from any distribution with mean  $\mu$  and known variance  $\sigma^2$ . Then, in Section 8.6, we derived a large-sample confidence interval for  $\mu$  with confidence coefficient approximately equal to  $1 - \alpha$ ,

$$[\bar{Y} - z_{\alpha/2}(\sigma/\sqrt{n}), \bar{Y} + z_{\alpha/2}(\sigma/\sqrt{n})].$$

If  $\sigma^2$  is unknown and  $n$  is large, can one replace  $\sigma$  with  $S_n$ ? Why?

(Thm 9.3) Suppose that  $U_n$  has a distribution function that converges to a standard normal distribution function as  $n \rightarrow \infty$ . If  $W_n$  converges in probability to 1, then the distribution function of  $U_n/W_n$  converges to a standard normal distribution.

(Example 9.4)



(Note)

- Example 9.4 implies that when  $n$  is large,  $\sqrt{n}(\bar{Y}_n - \mu)/S_n$  has approximately a standard normal distribution whatever is the form of the distribution from which the sample is taken.
- We know from Chapter 7(p.359) that if the sample is taken from a normal distribution and  $n$  is finite,  $\sqrt{n}(\bar{Y}_n - \mu)/S_n$  has a  $t$  distribution with  $n - 1$  degrees of freedom
- This implies that if a large sample is taken from a normal distribution, the distribution of  $\sqrt{n}(\bar{Y}_n - \mu)/S_n$  can be approximated by a standard normal distribution.

In other words, if  $n$  gets large, then the number of degrees of freedom also gets large, and the  $t$ -distribution can be approximated by a standard normal distribution (see Table 4 and 5 in pp.848-849).

(Exercise) Suppose that  $Y \sim b(n, p)$ . Show that

- 1)  $\hat{p}_n = Y/n$  is an unbiased estimator of  $p$
- 2)  $(\hat{p}_n - p)/\sqrt{\hat{p}_n\hat{q}_n/n}$  converges to a standard normal distribution.
- 3) Derive a large-sample confidence interval for  $p$  with confidence coefficient  $1 - \alpha$

## 9.4 Sufficiency

- We summarize the information in the sample by using a statistic,  $g(Y_1, \dots, Y_n)$  for the parameters of interest.

(Example)  $\bar{Y}$  and  $S^2$  as the unbiased estimator for the population mean  $\mu$  and variance  $\sigma^2$

(Question) Does the process of summarizing the original set of  $n$  sample observations to a few statistics,  $g_1(Y_1, \dots, Y_n), \dots, g_m(Y_1, \dots, Y_n)$  retain all the information about the parameters of interest in  $n$  sample observations?

(Answer) There are methods to find statistics  $g_1(Y_1, \dots, Y_n), \dots, g_m(Y_1, \dots, Y_n)$  summarizing all the information in a sample about the parameters of interest : we call such statistics *sufficient* statistics.

- Sufficient statistic

(Def 9.3) Let  $Y_1, \dots, Y_n$  denote a random sample from a probability distribution with unknown parameter  $\theta$ . Then the statistic  $U = g(Y_1, \dots, Y_n)$  is said to be *sufficient* for  $\theta$  if the conditional distribution of  $Y_1, \dots, Y_n$ , given  $U$ , does not depend on  $\theta$ .

(Example) Consider the outcomes of  $n$  independent trials of a binomial experiment,  $X_1, \dots, X_n$  where  $X_i \sim b(1, p)$  and  $p = P(\text{the } i\text{-th trial is a success})$ . Suppose we are given a value of  $Y = \sum_{i=1}^n X_i$  (i.e., # of successes among  $n$  trials). If we know the value of  $Y$ , do we have to look at  $X_1, \dots, X_n$  or other functions of  $X_1, \dots, X_n$  in order to gain further information about  $p$ ?

Why?  $P(X_1 = x_1, \dots, X_n = x_n \mid Y = y)$  does not depend on  $p$  (i.e., as long as  $Y$  is known, there is no further information about  $p$  from other functions of  $X_1, \dots, X_n$ ). We call  $Y$  a sufficient statistic for  $p$ .

- How to find a sufficient statistic

- (Def 9.3) tells us how to check if a statistic is sufficient or not.

- *likelihood of the sample*(Def 9.4) and *factorization criterion*(Thm 9.4) help find a sufficient statistic for  $\theta$

(Def 9.4) Let  $y_1, \dots, y_n$  be sample observation taken on corresponding random variables  $Y_1, \dots, Y_n$  whose distribution depends on  $\theta$ . Then the *likelihood of obtaining the sample*  $y_1, \dots, y_n$  when the parameter is  $\theta$ ,  $L(\theta) = L(y_1, \dots, y_n | \theta)$ , is defined to be

*i)*  $p(y_1, \dots, y_n | \theta)$  for discrete  $Y_1, \dots, Y_n$

*ii)*  $f(y_1, \dots, y_n | \theta)$  for continuous  $Y_1, \dots, Y_n$

- what  $L(\theta_1) > L(\theta_2)$  means for given  $Y_1 = y_1, \dots, Y_n = y_n$ ?

(Def 9.4 *continued*) If the set of random variables  $Y_1, \dots, Y_n$  is a random sample from  $p(y | \theta)$  or  $f(y | \theta)$  (i.e.,  $Y_1, \dots, Y_n \sim^{iid} p(y | \theta)$  or  $f(y | \theta)$ ), then the *likelihood of the sample*  $L(\theta) = L(y_1, \dots, y_n | \theta)$  is

i) for discrete  $Y_1, \dots, Y_n$ ,

$$p(y_1, \dots, y_n | \theta) = \prod_{i=1}^n p(y_i | \theta) = p(y_1 | \theta) \times \dots \times p(y_n | \theta)$$

ii) for continuous  $Y_1, \dots, Y_n$ ,

$$f(y_1, \dots, y_n | \theta) = \prod_{i=1}^n f(y_i | \theta) = f(y_1 | \theta) \times \dots \times f(y_n | \theta)$$

(Thm 9.4: factorization criterion) Let  $U$  be a statistic based on the random sample  $Y_1, \dots, Y_n$ . Then  $U$  is a *sufficient statistic* for the estimation of  $\theta$  iff  $L(\theta) = L(y_1, \dots, y_n | \theta)$  can be factored into two nonnegative functions :

$$L(y_1, \dots, y_n | \theta) = g(u, \theta) \times h(y_1, \dots, y_n)$$

where  $h(y_1, \dots, y_n)$  is either a constant or a function of  $y_1, \dots, y_n$  (not a function of  $\theta$ ).

(Example 9.5)

(Exercises) Let  $Y_1, \dots, Y_n$  be IID samples from  $N(\mu, \sigma^2)$ . Then,

1) If  $\mu$  is unknown and  $\sigma^2$  is known, show  $\bar{Y}$  is sufficient for  $\mu$ .

2) If  $\mu$  is known and  $\sigma^2$  is unknown, show  $\sum_{i=1}^n (Y_i - \mu)^2$  is sufficient for  $\sigma^2$ .

3) If  $\mu$  and  $\sigma^2$  are unknown, show  $\sum_{i=1}^n Y_i$  and  $\sum_{i=1}^n Y_i^2$  is jointly sufficient for  $\mu$  and  $\sigma^2$

- Note for sufficient statistics

i) there are many possible sufficient statistics for one parameter.

- the random sample itself,  $Y_1, \dots, Y_n$

- $Y_{(1)} \leq \dots \leq Y_{(n)}$  : the set of order statistics from  $Y_1, \dots, Y_n$

- Example 9.5: a one-to one function of  $\bar{Y}$

ii) we want to find a sufficient statistic that reduces the data in the sample as much as possible by using (Thm 9.4: factorization criterion)

- sufficient statistics are useful to develop unbiased estimators with minimum variance (See 9.5).

## 9.5 Minimum-Variance Unbiased Estimation

### Motivation

Search an estimator with good properties: unbiasedness, sufficiency and minimum variance.

### Procedure

S1 Find a sufficient statistic  $U$  for  $\theta$

S2 Check if  $E(U) = \theta$

S3 If yes,  $U$  is a *minimum-variance unbiased estimator* (MVUE) for  $\theta$ .

If no, find a function of  $U$ , say  $h(U)$  such that  $E(h(U)) = \theta$  from  $E(U)$ . Then  $h(U)$  is a *MVUE* for  $\theta$ .

### Why?

(Thm 9.5) **The Rao-Blackwell Theorem**

Let  $\hat{\theta}$  be an unbiased estimator for  $\theta$  such that  $V(\hat{\theta}) < \infty$ . If  $U$  is a sufficient statistic for  $\theta$ , define  $\hat{\theta}^* = E(\hat{\theta} | U)$ . Then, for all  $\theta$ ,

$$E(\hat{\theta}^*) = \theta \quad \text{and} \quad V(\hat{\theta}^*) \leq V(\hat{\theta}).$$



## Why? (continued)

- Given an unbiased estimator  $\hat{\theta}$  for  $\theta$  and a sufficient statistic  $U$  for  $\theta$ , there is a function of  $U$  that is also an unbiased estimator for  $\theta$  and has no larger variance than  $\hat{\theta}$ .
- Which sufficient statistic should we use in this theorem? a sufficient statistic identified from *the factorization criterion* (Thm 9.4)

(Example 9.6)

(Example 9.9)

(Example 9.8)

## 9.6 The Method of Moments

### Motivation

The  $k$ -th sample moment  $m'_k = \frac{1}{n} \sum_{i=1}^n Y_i^k$  should provide good estimates of the corresponding  $k$ -th population moment  $E(Y^k)$  where  $E(Y^k)$  are functions of the population parameters.

### Procedure

Suppose there are  $r$  unknown parameters. Then solve the  $r$  equations  $E(Y^k) = m'_k = \frac{1}{n} \sum_{i=1}^n Y_i^k$  for  $k = 1, 2, \dots, r$  with respect to  $r$  unknown parameters. We call the solutions to the  $r$  equations the *Moment Estimators* for the parameters.

[Note]

1.  $m'_k = \frac{1}{n} \sum_{i=1}^n Y_i^k$  is a consistent estimators of the corresponding  $E(Y^k)$ .
2. This method is easy to employ and usually provides consistent estimators for respective parameters
3. The estimators derived from this method are often not functions of sufficient statistics, so that their variances are sometimes larger than others(i.e., are sometimes not very efficient)

(Example 9.11, 9.12)

(Example 9.13))

(Exercise) Suppose  $Y_1, \dots, Y_n$  are IID samples from  $N(\mu, \sigma^2)$ , and  $\mu$  and  $\sigma^2$  are unknown. Find the method-of-moments estimators of  $\mu$  and  $\sigma^2$ .

## 9.7 The Method of Maximum Likelihood

### Motivation

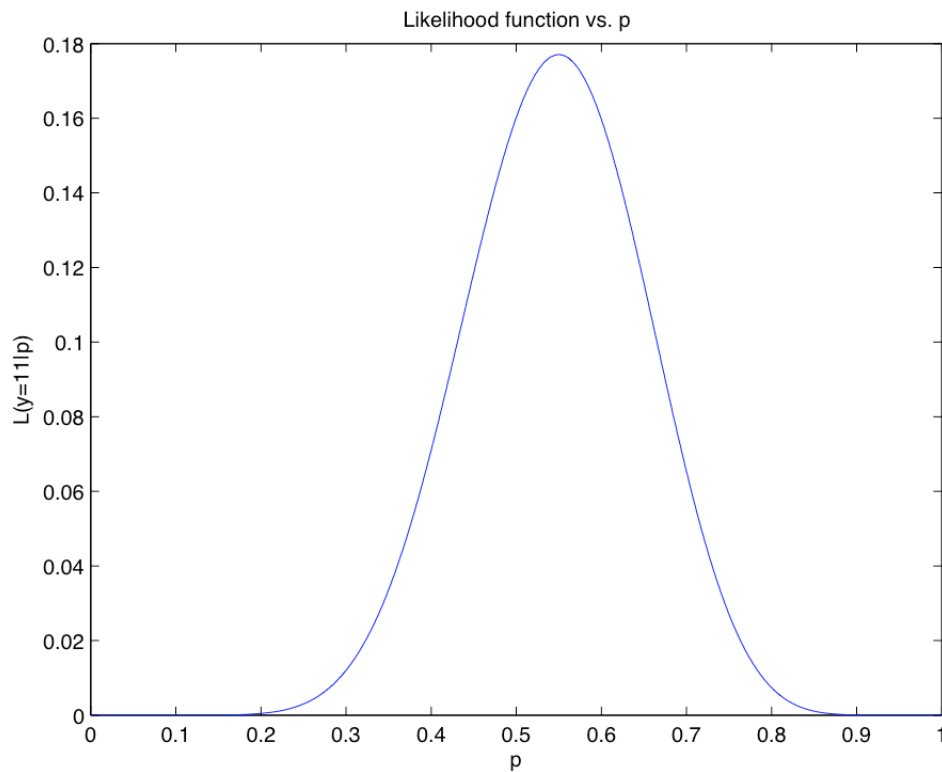
Let  $y_1, \dots, y_n$  be sample observation taken on corresponding random variables  $Y_1, \dots, Y_n$  whose distribution depends on  $\theta$ . Then, we would like to get the value of  $\theta$  (i.e., an estimate) which makes the observed data,  $y_1, \dots, y_n$  most likely (i.e., maximize the likelihood of obtaining the observed sample  $y_1, \dots, y_n$ ).

### Example

Suppose one tosses a coin 20 times, and count the number of heads observed. Let  $p$  the probability that one observes the head. Suppose we observed 11 heads and 9 tails. Then we would like to get the value of  $p$  which maximizes the likelihood of obtaining the observed data.

### Example(continued)

Let  $Y$  be the number of heads in 20 coin tosses. Then  $Y \sim b(20, p)$ . Since we observed 11 heads(i.e.,  $y = 11$ ), the likelihood of observing  $y = 11$  is  $L(y = 11 | p) = \binom{20}{11}p^{11}(1 - p)^9$ .



The value of  $p$  maximizing  $L(y = 11 | p)$  is 0.55.

## Procedure

- S1 Write down the likelihood  $L(\theta_1, \dots, \theta_k) = L(y_1, \dots, y_n \mid \theta_1, \dots, \theta_k)$
- S2 Find estimates for  $\theta_1, \dots, \theta_k$  that maximize  $L(\theta_1, \dots, \theta_k)$  or  $\ell(\theta_1, \dots, \theta_k) = \log L(\theta_1, \dots, \theta_k)$ . We call them the Maximum Likelihood Estimator(MLE)s for the parameters.
- take the derivative of  $\ell(\theta_1, \dots, \theta_k)$  with respect to  $\theta_1, \dots, \theta_k$ , set them to zero, and solve them  $\theta_1, \dots, \theta_k$ .
  - draw a plot of  $L(\theta)$  vs.  $\theta$  and find  $\hat{\theta}$  maximizing  $\ell(\theta)$  (when the range of  $\theta$  depends on the samples)

(Example 9.14)

(Example 9.15)

(Example 9.16)

[Important notes]

1. This estimation method often leads to MVUEs.

- The MLE is always some functions of any sufficient statistics  $U$  for  $\theta$ . Why?

$$: L(y_1, \dots, y_n | \theta) = g(u, \theta) \times h(y_1, \dots, y_n)$$

$$\Leftrightarrow \ell(y_1, \dots, y_n | \theta) = \ln g(u, \theta) + \ln h(y_1, \dots, y_n)$$

Then maximization of  $\ell(y_1, \dots, y_n | \theta)$  over  $\theta$  is equivalent to maximization of maximization of  $\ln g(u, \theta)$  over  $\theta$ , as  $\ln h(y_1, \dots, y_n)$  does not depend on  $\theta$ .

Moreover,  $\ln g(u, \theta)$  depends on the data only through the values of  $U$ .

- Therefore, if the MLE for  $\theta$  is adjusted to be unbiased, the resulting estimator often is an MVUE of  $\theta$ .

2. The MLE for  $\theta$  has the *invariance property*

- Suppose  $\hat{\theta}$  is the MLE for  $\theta$  and one is interested in estimating a function of  $\theta$ ,  $h(\theta)$ . Then the MLE  $\widehat{h(\theta)}$  for  $h(\theta)$  is  $\widehat{h(\theta)} = h(\hat{\theta})$ .

In (Example 9.16) the MLE of  $\sigma^2$  is  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$ . What is the MLE of  $\sigma$ ?