

## Chapter 8

# Point Estimation

Inferential statistics consists of two main branches — estimation and hypothesis testing. Point estimation selects a value of a sample statistic to estimate a population parameter. We call the value this statistic a *point estimate*. Since statistics are random variables, the point estimate is simply a realization of the random variable. For example,  $\bar{X}$  and  $S^2$  are used to find point estimates  $\bar{x}$  and  $s^2$  for the mean and variance of a population. Point estimates are numbers even though the statistics or *point estimators* themselves are random variables. We will use lower case to denote the point estimates and upper case to denote the statistic. There are also *interval estimators* and *interval estimates*, which we will discuss after this chapter. Since our point estimators or statistics are random variables, we look at their sampling distribution. However, the sampling distribution comes after selecting which point estimator to use. Point estimators have various properties we use to decide which estimator to use, including *unbiasedness*, *efficiency*, and *consistency*.

### 8.1 Unbiased Estimators

**Definition** A statistic  $\hat{\theta}$  is an *unbiased estimator* of the population parameter  $\theta$  if and only if

$$\mathbb{E}[\hat{\theta}] = \theta.$$

Unbiasedness means that across samples the estimator will be the parameter value on average. On average means over repeated samples. For instance, suppose we take  $M$  different samples of size  $n$  and estimate  $\hat{\theta}$   $M$  times where  $\theta = 1.11$ .

$m$	$\hat{\theta}$
1	1.10
2	1.05
3	1.12
.	.
.	.
.	.
$M$	1.09

If  $\hat{\theta}$  is an unbiased estimator of  $\theta$ , then the average of the  $M$  estimates should be close to  $\theta$ .

**Example** We showed that  $\mathbb{E}[\bar{X}] = \mu$  when  $X_i$  is *i.i.d.*, so the sample mean is an unbiased estimator for the population mean when the random sample is independently and identically distributed.

**Example** Suppose we have a random sample  $X_1, X_2, \dots, X_n$  from the population represented by the random variable  $x$  with a PDF

$$f(x) = e^{-(x-\delta)}, \quad x > \delta.$$

$$\begin{aligned}
\mathbb{E}[X] &= \int_{\delta}^{\infty} x \cdot e^{-(x-\delta)} dx \\
&= xe^{-(x-\delta)} \Big|_{\delta}^{\infty} + \int_{\delta}^{\infty} e^{-(x-\delta)} dx \\
&= \delta - e^{-(x-\delta)} \Big|_{\delta}^{\infty} \\
&= \delta - (0 - 1) \\
&= 1 + \delta
\end{aligned}$$

We know  $\mathbb{E}[\bar{X}] = 1 + \delta$ . So while the sample mean is an unbiased estimator for the mean, it is a *biased* estimator for the unknown parameter  $\delta$ . An unbiased estimator for  $\delta$  will simply be  $\bar{X} - 1$  where  $-1$  is the *bias correction factor*.

**Definition** The *finite sample bias*  $b_n(\theta)$  is

$$b_n(\theta) = \mathbb{E}[\hat{\theta}] - \theta.$$

The bias we defined is a finite sample bias, because the sample size is finite. We can also define asymptotic bias, which is when the sample size  $n \rightarrow \infty$ .

**Definition**  $\hat{\theta}$  is an *asymptotically unbiased estimator* of  $\theta$  if and only if

$$\lim_{n \rightarrow \infty} b_n(\theta) = 0.$$

Similarly, the *asymptotic bias* for a estimator  $\hat{\theta}$  is

$$\lim_{n \rightarrow \infty} b_n(\theta) = \lim_{n \rightarrow \infty} (\mathbb{E}(\hat{\theta}) - \theta).$$

**Proposition 8.1.1**  $S^2$  is an unbiased estimator of  $\sigma^2$  and  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  is an asymptotically unbiased estimator of  $\sigma^2$ .

**Proof**

$$\begin{aligned}
\mathbb{E}[S^2] &= \mathbb{E} \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] \\
&= \frac{1}{n-1} \mathbb{E} \left[ \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2 \right] \\
&= \frac{1}{n-1} \mathbb{E} \left[ \sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2 \right] \\
&= \frac{1}{n-1} \left[ \sum_{i=1}^n \mathbb{E}(X_i - \mu)^2 - n\mathbb{E}(\bar{X} - \mu)^2 \right] \\
&= \frac{1}{n-1} \left[ \sum_{i=1}^n \sigma^2 - n \frac{\sigma^2}{n} \right] \\
&= \frac{1}{n-1} [\sigma^2(n-1)] \\
&= \sigma^2
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right] &= \frac{1}{n} \mathbb{E} \left[ \sum_{i=1}^n ((X_i - \mu) - (\bar{X} - \mu))^2 \right] \\
&\quad \cdot \\
&\quad \cdot \\
&\quad \cdot \\
&= \frac{(n-1)\sigma^2}{n} \\
&\rightarrow \sigma^2 \text{ as } n \rightarrow \infty \quad \blacksquare
\end{aligned}$$

The naive sample variance estimator has a finite sample bias of  $-\frac{\sigma^2}{n}$ , but this bias disappears as the sample size grows. We have to multiply by the bias correction factor of  $\frac{n-1}{n}$  to obtain an unbiased estimator from the naive sample variance. Note that asymptotic unbiasedness depends on the sample size but unbiasedness does not. In this sense, unbiasedness is a stronger property. Typically, obtaining unbiased estimators is rather difficult but obtaining large samples is easy, so we often opt for biased estimators even though we desire this property.

## 8.2 Efficiency

Efficiency is a relative property that trades off bias of an estimator with its precision. Precision is the inverse of the variance. More precise estimators have smaller variance. Many estimators yield precise but biased estimates, so efficiency provides a criterion to guide us in trading off precision and bias. The criterion often used is called the *mean-squared error*.

**Definition** The mean-squared error (*MSE*) for an estimator  $\hat{\theta}$  is

$$MSE(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta]$$

where  $\mathbb{E}[\hat{\theta}]$  is not necessarily  $\theta$ . We can express the *MSE* as

$$MSE(\hat{\theta}) = bias^2 + Var(\hat{\theta}),$$

which is the squared bias plus the variance.

Now we can define efficiency.

**Definition** A estimator  $\hat{\theta}$  is *more efficient* than an estimator  $\hat{\phi}$  if and only if

$$MSE(\hat{\theta}) < MSE(\hat{\phi}).$$

Thus for unbiased estimators  $\tilde{\theta}$  and  $\tilde{\phi}$ ,  $\tilde{\theta}$  is more efficient than  $\tilde{\phi}$  if and only if

$$Var(\tilde{\theta}) < Var(\tilde{\phi}).$$

Efficiency is a comparative property. An estimator can only be more or less efficiency – not efficient. Efficiency is a notion of reliability, because it weights both bias and precision. We can also define *asymptotic efficiency* which compares the mean-squared error as the sample size grows.

**Example** Consider a new estimator  $Y = \frac{\bar{X}}{2}$ . Which estimator for the mean is more efficient?  $Y$  or  $\bar{X}$ ?

$$\begin{aligned}
Var(Y) &= \frac{\sigma^2}{4n} < \frac{\sigma^2}{n} = Var(\bar{X}) \\
b_n(Y) &= -\frac{\mu}{2} > b_n(\bar{X}) = 0 \\
MSE(Y) &= \frac{\mu^2}{4} + \frac{\sigma^2}{4n}, \quad MSE(\bar{X}) = \frac{\sigma^2}{n}
\end{aligned}$$

Then  $Y$  is more efficient than  $\bar{X}$  if and only if  $n\mu^2 < 3\sigma^2$ . So when the mean is small compared to the variance, we can shrink the sample mean to obtain biased estimates of the population mean with a high precision gain compared to the bias introduced over the sample mean.

### 8.3 Consistency

We use the mean-squared error to measure how close our estimate is to the true population parameter. We have weaker notions of closeness as the sample size grows. It is often difficult to find unbiased estimators, but we can be content with estimators that get close to the truth in large samples. We call this notion *consistency*. There are several types of consistency. We will focus on *mean-squared error consistency*.

**Definition** A statistic  $\hat{\theta}$  is a *consistent* estimator of the parameter  $\theta$  from some distribution if and only if for each  $\epsilon > 0$

$$\lim_{n \rightarrow \infty} \Pr(|\hat{\theta} - \theta| < \epsilon) = 1.$$

In other words, the statistic *converges in probability* to the true parameter value. We often denote this several ways:

$$\begin{aligned} \hat{\theta} &\rightarrow_p \theta \\ \text{plim}_{n \rightarrow \infty} \hat{\theta} &= \theta \end{aligned}$$

**Definition** A statistic  $\hat{\theta}$  is a *mean-squared error consistent* estimator of the parameter  $\theta$  from some distribution if and only if

$$\begin{aligned} \mathbb{E}[\hat{\theta}] &= \theta \\ \text{Var}(\hat{\theta}) &\rightarrow 0 \text{ as } n \rightarrow \infty \end{aligned}$$

This property corresponds to the *MSE* going to zero as the sample size grows large. Mean-squared error consistency means the sampling distribution of the statistic collapses onto the true value as the sample size grows, which implies convergence in probability.

**Proposition 8.3.1** *Mean-squared error consistency implies consistency.*

So far we have only looked at properties of estimators. Now we turn to looking at two examples of estimators. Each of the estimators aims to satisfy some objective.

### 8.4 Method of Moments

Method of moments estimation takes population parameters equates them to the sample analogue of the population moments. In many cases, we look at several moments and each moment forms an equation. We require as many moments as parameters obtain a unique solutions for the parameter estimators.

**Definition** The  $k^{th}$  order sample moment for set of observations  $x_1, x_2, \dots, x_n$  is

$$\frac{1}{N} \sum_{i=1}^N x_i^k$$

which corresponds to the  $k^{th}$  order moment  $\mathbb{E}[x_i^k]$  of the random variable  $x_i$ .

**Example** Suppose the random variable  $x_i \sim U(\alpha, 1)$  and we have  $n$  observations of  $x_i$ . What is the method of moments estimator of  $\alpha$  using the first moment? The first moment is

$$\mathbb{E}[x_i] = \frac{\alpha + 1}{2}$$

The sample analogue is

$$\frac{1}{N} \sum_{i=1}^N x_i$$

Hence, our method of moments estimator must satisfy

$$\frac{1}{N} \sum_{i=1}^N x_i = \frac{\hat{\alpha} + 1}{2} \Rightarrow \hat{\alpha} = 2\bar{x} - 1$$

Note that we could have used any moment to find an estimator of  $\alpha$ .

## 8.5 Maximum Likelihood

Maximum likelihood estimation (MLE) possesses many nice asymptotic properties like asymptotic normality and minimum variance. You will talk about this class of estimators extensively in econometrics. Here we will introduce them. In Baye's Rule,

$$P(C|I) = \frac{P(C)P(I|C)}{P(I)}$$

we call  $P(I|C)$  the likelihood. This probability represents the chances of observing the data (or realizations of the random variable  $I$ ) we do given some event ( $C$ ). The likelihood tells us the chances or likelihood of observing some data given some underlying parameters ( $\theta$ ). MLE aims to choose the underlying parameters to maximise the probability of observing the data we observe.

**Definition** Let  $x_1, \dots, x_n$  be the values of a random sample from a population with parameter  $\theta$ . The *likelihood function* of the sample is

$$L(\theta) = f(x_1, x_2, \dots, x_n; \theta),$$

which is the joint PDF for the data given parameter(s)  $\theta$ . The value of  $\theta$  that maximises  $L(\theta)$  is the maximum likelihood estimator of  $\theta$ .

We almost never maximise  $L(\theta)$ . A nice property of MLE is that we obtain the same estimator of  $\theta$  by maximising an increasing transformation of the likelihood function. We usually maximise  $\log L(\theta)$  instead.

**Example** Suppose  $X \sim \text{Binomial}(\theta)$ . Then the likelihood function is

$$L(\theta) = {}_nC_x \theta^x (1 - \theta)^{n-x} \Rightarrow \log L(\theta) = \log({}_nC_x) + x \log(\theta) + (n - x) \log(1 - \theta)$$

The loglikelihood function is concave in  $\theta$ , so the first order condition is necessary and sufficient for a maximum.

$$\frac{\partial \log L(\theta)}{\partial \theta} = \frac{x}{\hat{\theta}} - \frac{n - x}{\hat{\theta}} = 0$$

Solving for  $\hat{\theta}$ , we see  $\hat{\theta} = \frac{x}{n}$  is the maximum likelihood estimate of  $\theta$ . Hence  $\hat{\theta}_{MLE} = \frac{X}{n}$  is the maximum likelihood estimator of  $\theta$ .