

# Introduction to Sampling Distributions and Statistical Estimation

James H. Steiger

November 12, 2003

## 1 Topics for this Module

1. Parameters and Statistics
  - (a) Reversing the Information Flow in Statistical Inference
2. Sampling Distributions
3. Sampling Error
4. Principles of “Good Estimation”
  - (a) Unbiasedness
  - (b) Consistency
  - (c) Efficiency
  - (d) Sufficiency
  - (e) Maximum Likelihood
5. Practical vs. Theoretical Considerations

## 2 Parameters and Statistics

In statistical estimation we use a *statistic* (a function of a sample) to estimate a *parameter*, a numerical characteristic of a statistical population. In the preceding discussion of the binomial distribution, we discussed a well-known statistic, the sample proportion  $\hat{p}$ , and how its long-run distribution over repeated samples can be described, using the binomial process and the binomial distribution as models. We found that, if the binomial model is correct, we can describe the exact distribution of  $\hat{p}$  (over repeated samples) if we know  $N$  and  $p$ , the parameters of the binomial distribution. The situation can be diagrammed as in the top of Figure 1 below.

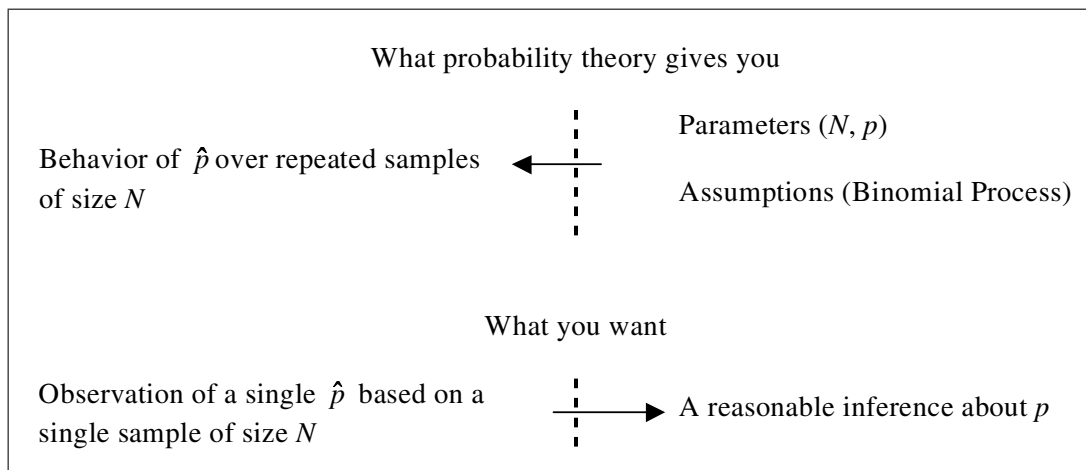


Figure 1: *Reversing the information flow in statistical inference*

### 2.1 Reversing the Information Flow in Statistical Inference

Probability theory tells us about the *long run* behavior of  $\hat{p}$ , but this requires specification of precisely what we do not know, i.e.,  $p$ , the proportion in the population. However, what we would like to have is something different from what probability theory provides directly. Generally what we have is *one* sample, of size  $N$ , and what we would *like* probability theory to provide for us is knowledge about  $p$  on the basis of the information in our data. But

it doesn't, and so probability theory, as important as it is, provides just the beginning point for statistical inference.

Looking back at Figure 1, what we would like is to turn around the direction of information flow. In probability theory, knowledge of  $p$  and  $N$  leads to knowledge about the long run behavior of  $\hat{p}$ . In statistical inference, we would like something else — a method to use knowledge of  $\hat{p}$  and  $N$  to lead to knowledge of  $p$ .

### 3 Sampling Distributions

A *statistic* is any function of the sample. Over repeated samples, statistics will almost always vary in value. So, over repeated samples, a statistic will have a *sampling distribution*. Sampling distributions have several characteristics:

1. Exact sampling distributions are difficult to derive
2. They are often different in shape from the distribution of the population from which they are sampled
3. They often vary in shape (and in other characteristics) as a function of  $N$ .

### 4 Sampling Error

Consider any statistic  $\hat{\theta}$  used to estimate a parameter  $\theta$ . For any given sample of size  $N$ , it is virtually certain that  $\hat{\theta}$  will not be equal to  $\theta$ . We can describe the situation with the following equation in random variables

$$\hat{\theta} = \theta + \varepsilon$$

where  $\varepsilon$  is called “sampling error,” and is defined tautologically as

$$\varepsilon = \hat{\theta} - \theta \tag{1}$$

i.e., the amount by which  $\hat{\theta}$  is wrong. In most situations,  $\varepsilon$  can be either positive or negative.

## 5 Principles of “Good Estimation”

A statistic that is used to estimate a particular parameter is called an *estimator* of that parameter. A realized value of the estimator is called an *estimate* of the parameter. Although some of the estimators that we use (like the sample mean  $\bar{X}_\bullet$  and the sample proportion  $\hat{p}$ ) are the sample analogs of the population quantities they estimate, many other estimators (for example,  $s^2$ , the sample variance) are not.

For any parameter, there are many possible estimators. Generally, an estimator in wide use has achieved popularity because it satisfies one or more *optimality criteria*, i.e. qualities that a good estimator is supposed to have. Below, we discuss a number of commonly used criteria for a good estimator.

### 5.1 Unbiasedness

**Definition 5.1 (*Unbiased Estimator*)** An estimator  $\hat{\theta}$  of a parameter  $\theta$  is unbiased if  $E(\hat{\theta}) = \theta$ , or, equivalently, if  $E(\varepsilon) = 0$ , where  $\varepsilon$  is sampling error as defined in Equation 1.

Ideally, we would like the positive and negative errors of an estimator to balance out in the long run, so that, on average, the estimator is neither high (an overestimate) nor low (an underestimate).

### 5.2 Consistency

We would like an estimator to get better and better as  $N$  gets larger and larger, otherwise we are wasting our effort gathering a larger sample. If we define some error tolerance  $\epsilon$ , we would like to be sure that sampling error  $\varepsilon$  is almost certainly less than  $\epsilon$  if we let  $N$  get large enough. Formally we say the following.

**Definition 5.2 (*Consistency*)**. An estimator  $\hat{\theta}$  of a parameter  $\theta$  is consistent if for any error tolerance  $\epsilon > 0$ , no matter how small, a sequence of statistics  $\hat{\theta}_N$  based on a sample of size  $N$  will satisfy the following

$$\lim_{N \rightarrow \infty} \Pr \left( \left| \hat{\theta}_N - \theta \right| < \epsilon \right) = 1 \quad (2)$$

**Example 5.1 (An Unbiased, Inconsistent Estimator)** Consider the statistic  $D = (X_1 + X_2)/2$  as an estimator for the population mean. No matter how large  $N$  is,  $D$  always takes the average of just the first two observations. This statistic has an expected value of  $\mu$ , the population mean, since

$$E\left(\left[\frac{1}{2}X_1 + \frac{1}{2}X_2\right]\right) = \frac{1}{2}E(X_1) + \frac{1}{2}E(X_2) = \frac{1}{2}\mu + \frac{1}{2}\mu = \mu$$

but it does not keep improving in accuracy as  $N$  gets larger and larger.

### 5.3 Efficiency

All other things being equal, we prefer estimators with a smaller sampling errors. Several reasonable measures of “smallness” suggest themselves: (a) the average absolute error, and (b) the average squared error. Consider the latter. The variance of an estimator can be written

$$\sigma_{\hat{\theta}}^2 = E\left(\hat{\theta} - E(\hat{\theta})\right)^2$$

and when the estimator is unbiased,  $E(\hat{\theta}) = \theta$ , so the variance becomes

$$\sigma_{\hat{\theta}}^2 = E\left(\hat{\theta} - \theta\right)^2 = E(\varepsilon^2)$$

since  $\hat{\theta} - \theta = \varepsilon$ .

For an unbiased estimator, the sampling variance is also the average squared error, and is a direct measure of how inaccurate the estimator is, on average. More generally, though, one can think of sampling variance as the randomness, or noise, inherent in a statistic. (The parameter is the “signal.”) Such noise is generally to be avoided.

Consequently, the *efficiency* of a statistic is inversely related to its sampling variance, i.e.

$$Efficiency(\hat{\theta}) = \frac{1}{\sigma_{\hat{\theta}}^2}$$

The *relative efficiency* of two statistics is the ratio of their efficiencies, which is the inverse of the ratio of their sampling variances.

**Example 5.2 (Relative Efficiency)** Suppose statistic  $A$  has a sampling variance of 5, and statistic  $B$  has a sampling variance of 10. The relative efficiency of  $A$  relative to  $B$  is 2.

## 5.4 Sufficiency

An estimator  $\hat{\theta}$  is *sufficient* for estimating  $\theta$  if it uses all the information about  $\theta$  available in a sample. The formal definition is as follows:

**Definition 5.3 (*Sufficient Statistic*)** *Recalling that any statistic is a function of the sample, define  $\hat{\theta}(S)$  to be a particular value of an estimator  $\hat{\theta}$  based on a specific sample  $S$ . An estimator  $\hat{\theta}$  is a sufficient statistic for estimating  $\theta$  if the conditional distribution of the sample  $S$  given  $\hat{\theta}(S)$  does not depend on  $\theta$ .*

The fact that once the distribution is conditionalized on  $\hat{\theta}$  it no longer depends on  $\theta$ , shows that all the information that  $\theta$  might “reveal in the sample” is captured by  $\hat{\theta}$ .

## 5.5 Maximum Likelihood

The likelihood of a sample of  $N$  independent observations is simply the product of the probability densities of the individual observations. Of course, if you don’t know the parameters of the population distribution, you cannot compute the probability density of an observation. The *principle of maximum likelihood* says that the best *estimate* of a population parameter is the one that makes the sample most likely. Deriving estimators by the principle of maximum likelihood often requires calculus to solve the maximization problem, and so we will not pursue the topic here.

## 6 Practical vs. Theoretical Considerations

In any particular situation, depending on circumstances, you may have an overriding consideration that causes you to ignore one or more of the above considerations — for example the need to make as small an error as possible when using your own data. In some situations, any additional error of estimation can be extremely costly, and practical considerations may dictate a biased estimator if it can be guaranteed that a bias can reduce  $\varepsilon$  for that sample.