

1.017/1.010 Class 14

Estimation

Estimating Distributional Properties

We can use random samples to estimate the unknown properties of random variable x . These may be:

1. Distributional **parameters**, such as the lower and upper limits of a uniform distribution (a and b)
2. Other distributional properties such as the **mean, variance, 90th percentile** value, etc.

Parametric statistics assumes that the form of the x distribution of x is known. **Nonparametric statistics** makes no assumptions about distribution of x .

Estimators

An **Estimator** (or a **statistic**) $\hat{a}(x_1, x_2, \dots, x_N)$ is a function used to derive an **estimate** \hat{a} of the unknown distributional property a from a random sample x_1, x_2, \dots, x_N

$$\hat{a} = \hat{a}(x_1, x_2, \dots, x_N)$$

Example:

Suppose we want to estimate from the random sample x_1, x_2, \dots, x_N the mean of $F_x(x)$. A possible choice for the estimator is the sample mean m_x . Then:

$$a = E[x]$$
$$\hat{a} = \hat{a}(x_1, x_2, \dots, x_N) = m_x = \frac{1}{N} \sum_{i=1}^N x_i$$

How do we know if this is a “good” estimator?

Properties of Good Estimates

We would like the estimate \hat{a} to be as “close” as possible to the true value a .

Since \hat{a} is derived from a random sample (a set of random variables) it is also a random variable, defined by the probability distribution $F_{\hat{a}}(\hat{a})$. This distribution describes how much \hat{a} might deviate from a .

We can sometimes derive $F_{\hat{a}}(\hat{a})$ or we may use the Central Limit theorem to justify the assumption that $F_{\hat{a}}(\hat{a})$ is a **normal distribution** (if a is derived from a large random sample).

If $F_{\hat{a}}(\hat{a})$ is normal, it is completely defined by the mean $E[\hat{a}]$ and variance $Var[\hat{a}]$ of \hat{a} .

Even if $F_{\hat{a}}(\hat{a})$ is not normal we can use the mean and variance of \hat{a} to measure the quality of the estimate. "Good" estimates should be **unbiased** and **consistent**:

1. Unbiased

The expected value of the estimate is the true property value:

$$E[\hat{a}] = a$$

2. Consistent:

The variability of the estimate goes to zero as the sample size increases:

$$Var[\hat{a}] \rightarrow 0 \quad \text{as} \quad N \rightarrow \infty$$

If the estimate is unbiased and consistent it will converge to the true value as the sample size increases.

In general, the mean and variance of an estimate are most readily estimated with **stochastic simulation** (see exercise below). In some cases they can be derived directly from the estimator function \hat{a} .

Example – Properties of the Sample mean:

Reconsider the sample mean m_x used to estimate the mean $E[x]$ of $F_x(x)$:

$$\hat{a} = \hat{a}(x_1, x_2, \dots, x_N) = m_x = \frac{1}{N} \sum_{i=1}^N x_i$$

The mean and variance of this estimate are:

$$E[\hat{a}] = E[m_x] = E\left[\frac{1}{N} \sum_{i=1}^N x_i\right] = E[x] = a$$

$$\text{Var}[\hat{a}] = \text{Var}[\mathbf{m}_x] = \text{Var}\left[\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i\right] = \frac{1}{N} \text{Var}[\mathbf{x}]$$

So the sample mean is an **unbiased** and **consistent** estimator of the population mean.

For large samples the Central Limit Theorem implies that \mathbf{m}_x is normally distributed. The distribution of \mathbf{m}_x is centered on $E[\mathbf{x}]$ and gradually converges to $E[\mathbf{x}]$ as N increases.

Exercise: Comparing Alternative Estimates of Population Properties

Suppose that you have a sample of 10 observations of a random variable x that you believe to be exponentially distributed. Your objective is to estimate the 90 percentile value x_{90} of this variable. This value of the solution of the equation $F_x(x_{90}) = 0.9$ [i.e. $F_x^{-1}(0.9) = x_{90}$]

Propose at least two different methods for estimating x_{90} from the 10 observations.

Compare the performance of these alternative estimators with a stochastic simulation that performs the following steps:

1. Generate many (e.g. 1000) replicates, each consisting of 10 observations drawn from an exponential distribution. Select your own value of the distributional parameter $a = E[x]$.
2. For each replicate derive an estimate \hat{x}_{90} of x_{90} from each of your two proposed estimators.
3. For each estimator compute the sample mean and variance of \hat{x}_{90} over all replicates. Also, use the replicates to construct a histogram and a CDF plot (using MATLAB's `normplot` function) for each \hat{x}_{90} .
4. Determine whether your estimators are unbiased and consistent (check consistency by plotting the rerunning your simulation for a much larger number of observations).

Which of your estimators is better? Explain your reasoning.

Some relevant MATLAB functions: `exprnd`, `hist`, `normplot`, `mean`, `var`

