# **COURSE**

3

## Sampling and Estimation

Sampling and estimation are fundamental in inferential statistics. They allow us to draw conclusions about a large population based on the analysis of a smaller representative subset (sample). In this chapter, we introduce the essential notions of sampling, random sampling distributions, and estimation of unknown parameters.



## Sampling

## 3.1.1 Concept of Sampling

## Definition 1.1

Consider a population  $\Omega$  of size N. A **sample** is a subset of this population. A sample of size n is thus a list of n individuals  $(\omega_1, \omega_2, ..., \omega_n)$  drawn from the parent population.

## Example 1.1

Consider a population composed of 5 students. We are interested in the weekly time devoted by each student to studying statistics.

$$\Omega = \{A, B, C, D, E\}, \quad N = 5$$

Student	Study Time (h)
A	7
В	3
С	6
D	10
E	4

## Definition 1.2

Sampling is the process of selecting samples. The ratio t of the sample size n to the population size N from which it is drawn is called the **sampling rate** or **sampling fraction**, i.e.

$$t = \frac{n}{N}$$

## Example 1.2

If we draw samples of size 2, then  $t = \frac{2}{5}$  (see Example 1.1).

## Definition 1.3

A random sample is a selection of n individuals from a parent population such that all possible combinations of n individuals have the same probability of being selected. Other types of sampling exist, but we will focus exclusively on random sampling.

## Remark

We aim to describe a qualitative or quantitative characteristic C of a population  $\Omega$  by studying the results obtained from a sample of size n.

## Example 1.3

- 1. For a given population, we may study quantitative characteristics such as weight or height.
- 2. For a given population, we may study qualitative characteristics such as eye color or hair color.
- 3. In the initial example, the characteristic studied is the weekly time devoted to studying statistics.

## Definition 1.4

Let C be a quantitative characteristic defined on a parent population  $\Omega$ . C is the realization of a random variable X defined on  $\Omega$ :

$$X: \Omega \to \mathbb{R}, \quad \omega_i \mapsto X(\omega_i) = x_i$$

A sample of values of X is the list of observed values  $(x_1, x_2, ..., x_n)$  taken by X on a sample  $(\omega_1, ..., \omega_n)$  of the population  $\Omega$ . The coordinates can be regarded as realizations of a random vector  $(X_1, ..., X_n)$  called an n-sample of X, where the  $X_i$  are independent and identically distributed (i.i.d.) random variables with the same distribution as X.

## Definition 1.5

Any random variable that can be expressed in terms of the random variables  $X_1, ..., X_n$  is called a **statistic**.

## Example 1.4

$$X_i$$
 and  $\overline{X} = \frac{1}{n} \sum_{i=1}^n X_i$  are examples of statistics.

#### Remark

If we extract several samples of the same size n, the results we obtain will vary depending on the sample considered. We call this variability **sampling fluctuations**. To make reliable inferences about the parent population, we must study the probability laws governing these fluctuations.

3.1. SAMPLING

## 3.1.2 Sampling Distributions

## Sample Mean and Sample Variance

## Definition 1.6

Consider a population  $\Omega$  whose elements possess a quantitative characteristic C that is the realization of a random variable X with expectation  $\mu$  and standard deviation  $\sigma$ . Assume the population is infinite or that sampling is done with replacement.

We draw a sample  $(X_1, ..., X_n)$  from X, giving observed values  $(x_1, ..., x_n)$ . The **sample mean** is given by:

$$\overline{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

The corresponding random variable is:

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

Similarly, the **sample variance** is:

$$v = \frac{1}{n} \sum_{i=1}^{n} (x_i - \overline{x})^2$$

and the associated random variable:

$$V = \frac{1}{n} \sum_{i=1}^{n} (X_i - \overline{X})^2$$

We define the random variable  $S^2$ , called the **unbiased sample variance**, as:

$$S^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{i} - \overline{X})^{2} = \frac{n}{n-1} V$$

## 3.1.3 Sample Proportion

## Definition 1.7

Sometimes, the characteristic to be estimated is not quantitative but qualitative. In this case, we seek the proportion p of individuals in the population possessing that characteristic. The proportion p is estimated from the results obtained in a sample of size n.

The observed proportion f in a sample is the realization of a random variable F, representing the frequency of appearance of this characteristic in the sample. F is called the **sample** proportion or statistical frequency:

$$F = \frac{K}{n}$$

where K is the random variable counting the number of occurrences of the characteristic in the sample of size n. By definition,  $K \sim B(n, p)$ , so that:

$$E(K) = np$$
,  $Var(K) = npq$  with  $q = 1 - p$ .

Therefore,

$$E(F) = p, \quad Var(F) = \frac{pq}{n}.$$

#### Remark

For  $n \geq 30$ , with  $np \geq 15$  and  $nq \geq 15$ , F can be approximated by a normal distribution:

$$F \sim N\left(p, \sqrt{\frac{pq}{n}}\right)$$
.

# 2 Estimation

#### 3.2.1 Estimators

## Definition 2.1

To **estimate a parameter** means to find an approximate value based on the results obtained from a sample.

An **estimator**  $\hat{\theta}$  of the unknown parameter  $\theta$  is a function that assigns to each set of observations  $(x_1, ..., x_n)$  an estimated value  $\hat{\theta}$ :

$$\hat{\theta}:(x_1,...,x_n)\longmapsto\hat{\theta}=f(x_1,...,x_n)$$

Hence,  $\hat{\theta}$  is a random variable. We can compute its expectation  $E(\hat{\theta})$  and variance  $Var(\hat{\theta})$ . These quantities measure the quality of the estimator for the parameter  $\theta$ .

## Example 2.1

Estimating the average height of a population from the empirical mean of a sample taken from that population.

3.2. ESTIMATION 5

## Definition 2.2

An estimator  $\hat{\theta}$  is said to be **unbiased** if the mean of its sampling distribution equals the true value of the parameter  $\theta$ :

$$E(\hat{\theta}) = \theta.$$

Otherwise, it is said to be **biased**.

The **bias** of an estimator is defined as:

$$Bias(\hat{\theta}) = E(\hat{\theta}) - \theta.$$

## Remark

The absence of bias does not necessarily imply that an estimator is efficient. A parameter can have multiple unbiased estimators. In such cases, efficiency is compared using their variances: an estimator with smaller variance provides estimates closer to the true value of  $\theta$ .

## Definition 2.3

An unbiased estimator  $\hat{\theta}_1$  is said to be **efficient** if, for any other unbiased estimator  $\hat{\theta}_2$ :

$$E(\hat{\theta}_1) = E(\hat{\theta}_2) = \theta$$
 and  $Var(\hat{\theta}_1) < Var(\hat{\theta}_2)$ .

## Definition 2.4

An estimator  $\hat{\theta}$  is said to be **consistent** (or convergent) if its distribution tends to concentrate around the true value  $\theta$  as the sample size increases, i.e.:

$$\lim_{n \to +\infty} Var(\hat{\theta}) = 0.$$

#### **Common Estimators**

(A) Quantitative Characteristic Let X be a random variable with mean  $\mu$  and standard deviation  $\sigma$  defined on a parent population  $\Omega$ , and let  $(X_1, ..., X_n)$  be a random sample.

## Properties

- 1.  $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$  is an unbiased and consistent estimator of  $\mu$   $(E(\overline{X}) = \mu)$ .
- 2.  $V = \frac{1}{n} \sum_{i=1}^{n} (X_i \overline{X})^2$  is a biased estimator of  $\sigma^2$ .
- 3.  $S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i \overline{X})^2 = \frac{n}{n-1} V$  is an unbiased and consistent estimator of  $\sigma^2$ .

## Property

For a qualitative characteristic with population proportion p, the sample proportion F is an unbiased and consistent estimator of p.

#### 3.2.2 Confidence Intervals

#### Definition 2.5

Rather than determining a single approximate value of a parameter  $\theta$ , we may seek an **interval** that contains the true value of  $\theta$  with a specified probability.

Let X be a random variable whose distribution depends on the parameter  $\theta$ . A **confidence** interval of risk  $\alpha$  for  $\theta$  is defined by random variables  $A_n$  and  $B_n$  such that:

$$P(A_n \le \theta \le B_n) = 1 - \alpha.$$

The realized interval [a, b] is obtained from a sample  $(x_1, ..., x_n)$  as:

$$a = A_n(x_1, ..., x_n), b = B_n(x_1, ..., x_n).$$

## Remarks

- 1. The quantity  $1 \alpha$  is called the **confidence level** of the interval [a, b], i.e.  $P(a \le \theta \le b) = 1 \alpha$ .
- 2. In practice, we often have only one sample that provides a single confidence interval [a, b].
- 3. The parameter to be estimated may be a mean, a variance (for quantitative variables), or a proportion (for qualitative ones).

#### Confidence Interval for a Mean

We consider the case where X follows a normal distribution  $N(\mu, \sigma)$ , or when the sample size is large (n > 30) so that  $\overline{X}$  approximately follows the same law.

Given a sample  $(x_1,...,x_n)$ , we define:

$$m = \frac{x_1 + \dots + x_n}{n}, \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - m)^2.$$

(A) Case  $\sigma$  known

$$IC = \left[ m - t_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}; m + t_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

(B) Case  $\sigma$  unknown

$$IC = \left[ m - t_{1-\alpha/2, n-1} \frac{s}{\sqrt{n}}; m + t_{1-\alpha/2, n-1} \frac{s}{\sqrt{n}} \right]$$

where  $t_{1-\alpha/2, n-1}$  is the quantile of order  $1 - \frac{\alpha}{2}$  of Student's t distribution with n-1 degrees of freedom.

3.2. ESTIMATION 7

## Remark

If n > 30, then  $t_{1-\alpha/2, n-1} \approx t_{1-\alpha/2}$ .

#### Confidence Interval for a Variance

#### (A) Case $\mu$ known

$$IC = \left[ \frac{nv}{\chi_{1-\alpha/2}^2(n)} \, ; \, \frac{nv}{\chi_{\alpha/2}^2(n)} \right]$$

where  $v = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^2$ , and  $\chi^2_{1-\alpha/2}(n)$  and  $\chi^2_{\alpha/2}(n)$  are the chi-squared quantiles of orders  $1 - \frac{\alpha}{2}$  and  $\frac{\alpha}{2}$  respectively.

## (B) Case $\mu$ unknown

$$IC = \left[ \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2(n-1)}; \frac{(n-1)s^2}{\chi_{\alpha/2}^2(n-1)} \right]$$

#### Remark

If n > 30, we can approximate:

$$\chi_{\alpha}^2(n-1) \approx \frac{1}{2} \left(t_{\alpha} + \sqrt{2n-3}\right)^2.$$

Hence:

$$IC = \left[ \frac{2(n-1)s^2}{(t_{1-\alpha/2} + \sqrt{2n-3})^2}; \frac{2(n-1)s^2}{(t_{\alpha/2} + \sqrt{2n-3})^2} \right]$$

and the symmetry of the standard normal law ensures that  $t_{\alpha/2} = -t_{1-\alpha/2}$ .

#### Confidence Interval for a Proportion

From the approximation  $F \sim N(p, \sqrt{\frac{pq}{n}})$  with q = 1 - p, we deduce:

$$IC = \left[ f - t_{1-\alpha/2} \sqrt{\frac{f(1-f)}{n}} \, ; \, f + t_{1-\alpha/2} \sqrt{\frac{f(1-f)}{n}} \right]$$

where f is the sample proportion.