

STATISTICS PRE-COURSE
PART 3
BASICS OF INFERENCE STATISTICS

Alfonso Russo

Department of Economics and Finance

Tor Vergata University of Rome

September 2024

PART III SYLLABUS

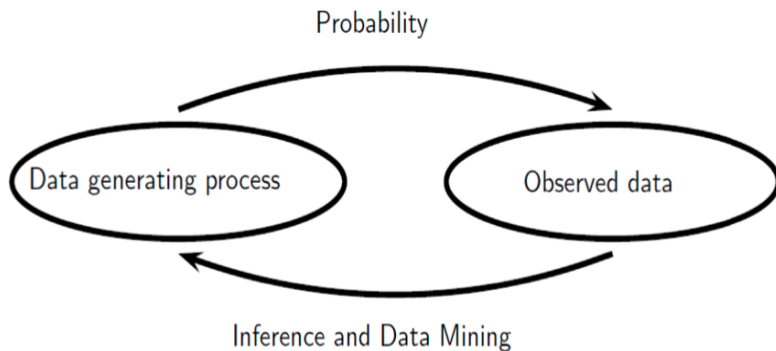
- 1 Basics of Inference
- 2 Confidence Intervals
- 3 Central Limit Theorem
- 4 Point Estimation
- 5 Statistical Models
- 6 Main Estimators

PROBABILITY AND INFERENCE

FROM POPULATION TO DATA AND BACK

- **Probability** starts from population, which is described by means of probability statements and functions, and predicts what happens in a sample extracted from it.
- **Inference** starts from a sample and describes the observed data in order to recover the data generating process.

INFERENCE IN A FIGURE



- **Estimation:** recover some *parameters* explaining the phenomena that generates the data
 - **point estimate:** a single number which is our best guess for a parameter
 - **interval estimate:** an interval of values which is believed to contain the true value of a parameter
- **Hypothesis testing:** using data to validate certain statements or predictions

- Suppose that you are interested in the (observable) random variable X having support \mathcal{X} . We denote its probability law by $f_X^*(x)$. Assume we do not know the exact probability law of X but that we know it actually belongs to a family \mathcal{F} of probability laws.
- Typically, \mathcal{F} is a parametric family:

$$\mathcal{F} = \{f_X(x | \theta), \theta \in \Theta\} \quad (1)$$

- Each member of the family depends on an unknown quantity θ that can take values in the *parameter space* Θ
- In Inference we want to determine, through an observed sample, what member of \mathcal{F} (or what value θ^* of θ) identifies the probability law of X .
- In other words, we assume that we know the functional form of the probability law (Bernoulli, Gamma, Normal, etc) but we do not know the value of the unknown parameter(s) to characterise it.

THE ELEMENTS OF A STATISTICAL MODEL

To characterise a random variable X , we need three basic elements:

- \mathcal{X}
- $f_X(x | \theta)$
- Θ

A statistical model is then the triple

$$\{\mathcal{X} ; f_X(x | \theta) ; \Theta\} \quad (2)$$

Note that

- X might be multidimensional
- $\Theta \in \mathbb{R}; \Theta \in \mathbb{R} \times \mathbb{R}$

Formulate the following statistical models:

- Bernoulli
- Geometric
- Normal
- Uniform

A **Random sample** is a collection of random variables $X_1, \dots, X_n \sim f_{X_1, \dots, X_n}$ that are

- independent

$$f_{X_1, \dots, X_n} = \prod_{i=1}^n f_{X_i}(x_i) \quad (3)$$

- identically distributed

$$f_{X_i}(x_i) = f_X(x_i) \quad \forall i \quad (4)$$

As a consequence

$$f_{X_1, \dots, X_n} = \prod_{i=1}^n f_X(x_i) \quad (5)$$

An **observed sample** (x_1, \dots, x_n) is a specific realisation of the random sample

HOW POPULATED IS THE POPULATION?

Find the difference!

- **Performance analysis of Statistics students this year:** Some of your tests are collected to assess your statistics knowledge.
- **Apple New Product:** Apple has developed a new device and some units are tested to estimate average duration.
- **How rich could Temple lawyers get?:** Lawyers in Temple are surveyed to estimate how remunerative are careers in Law.
- **Humidity in desertic areas:** Deserts worldwide are sampled to study humidity levels.

$$\left\{ \mathcal{X}^n ; f_n(x_n | \theta) = \prod_{i=1}^n f_X(x_i | \theta) ; \theta \in \Theta \right\} \quad (6)$$

EXAMPLE

Let X_1, \dots, X_n be i.i.d (independent and identically distributed) from a $Poisson(\lambda)$

The joint distribution takes form

$$\begin{aligned} f_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \prod_{i=1}^n f_X(x_i) \\ &= \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \\ &= \frac{1}{\prod_{i=1}^n x_i!} e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \end{aligned} \tag{7}$$

- **Parameter:** numerical characteristics of the population that we are able to recover (typically unknown)
 - Example: λ in a Poisson
- **Statistic:** a function of the sample that does not directly depend on any unknown parameter
 - Example: $S(X_1, \dots, X_n) = X_{(n)} - X_{(1)}$
- **Estimator:** a statistic used to estimate a population parameter
 - Example: $T(X_1, \dots, X_n) = \bar{X}$ is an estimator for μ
- **Estimate:** the value of an estimator corresponding to an observed sample
 - Example: $T(x_1, \dots, x_n) = \bar{x}$ is the value of \bar{X} coming from (x_1, \dots, x_n)

In order to assess the IQ of Tor Vergata students, we interview 10 people and use the sample mean \bar{X} as an estimator for the population mean μ .

- observed sample: $x = (x_1 = 95, x_2 = 104, x_3 = 104, x_4 = 95, x_5 = 88, x_6 = 126, x_7 = 77, x_8 = 112, x_9 = 111, x_{10} = 105)$
- estimate: $T(x_1, \dots, x_{10}) = \bar{x} = 101.7$

VARIABILITY OF ESTIMATORS

If we collected another sample, we would obtain different results

- 2nd observed sample: $x' = (123, 119, 94, 116, 106, 91, 88, 107, 91, 103)$
- estimate: $T(x'_1, \dots, x'_n) = \bar{x}' = 103.8$

Since it is a function of random objects, an estimator is a *random variable*, and the estimates are its *realisations*

There is no "universal estimator". We must choose it according to

- the distribution of the data

- we would not try to estimate the maximum of a discrete variable with a continuous value

- the parameter of interest

- we would not try to estimate the mean and the variance of a Normal distribution using the same estimator

If the parameter of interest is the expected value of the population $\mathbb{E}[X]$, then the obvious candidate is the **sample mean**

$$\bar{X} = n^{-1} \sum_{i=1}^n X_i \quad (8)$$

Good Properties:

- from the *Law of Large Numbers* we know that $\bar{X} \rightarrow \mathbb{E}[X]$ as $n \rightarrow \infty$
- the *Central Limit Theorem* provides us with an approximate distribution of \bar{X}

HOW DO WE DEFINE AN ESTIMATOR

The aim of an estimator is to recover and characterise the unknown distribution that generated the data

There are several *automatic* ways to derive an estimator, depending on how we use the data to estimate the generating distribution

■ Method of Moments:

- Find a distribution that has some features in common with the observed sample

■ Maximum Likelihood:

- find a distribution that maximises the probability of observing the sample at hand

The core idea behind the MM is to equate population and sample moments

$$\left\{ \begin{array}{l} \mathbb{E}[X] = n^{-1} \sum_{i=1}^n X_i \\ \mathbb{E}[X^2] = n^{-1} \sum_{i=1}^n X_i^2 \\ \mathbb{E}[X^3] = n^{-1} \sum_{i=1}^n X_i^3 \\ \dots \\ \dots \end{array} \right. \quad (9)$$

The MOM gives *consistent* estimators for population parameters. Weak assumptions needed. However, these will be biased!

METHOD OF MOMENTS

Suppose you want to estimate k unknown parameters $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$ that characterise the distribution $f_Y(y; \boldsymbol{\theta})$ of the random variable Y . Suppose further that *population moments* (the moments of the "true" distribution) can be expressed as functions of the unknown parameters

$$\begin{aligned}\phi_1 &:= \mathbb{E}[Y] = h_1(\theta_1, \dots, \theta_k) \\ \phi_2 &:= \mathbb{E}[Y^2] = h_2(\theta_1, \dots, \theta_k) \\ &\quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \\ \phi_k &:= \mathbb{E}[Y^k] = h_k(\theta_1, \dots, \theta_k)\end{aligned}\tag{10}$$

METHOD OF MOMENTS

Working with a sample (y_1, \dots, y_n) we could use the j th sample moment as an estimator $\hat{\phi}_j$ for ϕ_j .

More formally,

$$\hat{\phi}_j = n^{-1} \sum_{i=1}^n y_i^j \quad \text{for } j = 1, \dots, k \quad (11)$$

The Method of Moments estimator of θ , denoted by $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ is the solution to the following system

$$\begin{aligned} \hat{\phi}_1 &= h_1(\hat{\theta}_1, \dots, \hat{\theta}_k) \\ \hat{\phi}_2 &= h_2(\hat{\theta}_1, \dots, \hat{\theta}_k) \\ &\vdots \quad \vdots \quad \vdots \quad \vdots \\ \hat{\phi}_k &= h_k(\hat{\theta}_1, \dots, \hat{\theta}_k) \end{aligned} \quad (12)$$

- Let X_1, \dots, X_n be a random sample from a population with probability law

$$f_X(x | \theta) = \frac{1}{\theta} x^{\frac{1-\theta}{\theta}} \quad 0 < x < 1. \quad (13)$$

Find the MoM estimator for θ .

- Consider a sample $X_1, \dots, X_n \sim \mathcal{N}(1, \sigma^2)$. Find the MOM estimator for σ^2 .

THE LIKELIHOOD FUNCTION

BASIC INTUITION

Let $X \sim \text{Binomial}(n, p)$. The probability mass function

$$\mathbb{P}(X = x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad (14)$$

gives us the probability of observing a specific value x .

Assume that we know $n = 10$ and we suddenly observe $x = 8$.

■ **With $p = 0.5$:** $\mathbb{P}(X = 8) = \binom{10}{8} (0.5)^8 (0.5)^2 = 0.043$

■ **With $p = 0.7$:** $\mathbb{P}(X = 8) = \binom{10}{8} (0.7)^8 (0.3)^2 = 0.233$

For $x = 8$ a value of the parameter $p = 0.7$ seems more likely than $p = 0.5$.

THE LIKELIHOOD FUNCTION

BASIC INTUITION

When we fix a realisation x and we consider it to be a function of the unknown parameter p , the p.m.f. $\binom{n}{x} p^x (1-p)^{n-x}$ gives us a measure of how compatible is x is p .

This is called the **Likelihood** of p .

NB The Likelihood tells us how **plausible** a value of the parameter is, but it does not measure its **probability**

THE LIKELIHOOD FUNCTION

MORE FORMALLY

- Given a statistical model $\{\mathcal{X}^n ; f_n(\mathbf{x}_n | \theta) ; \theta \in \Theta\}$, the function $\mathcal{L} : \Theta \rightarrow \mathbb{R}^+$ defined as

$$\mathcal{L}(\theta | \mathbf{x}_n) \propto f_n(\mathbf{x}_n | \theta) , \quad \theta \in \Theta \quad (15)$$

is called the likelihood function associated to the observed sample \mathbf{x}_n .

- Since $\mathbf{x}_n = (x_1, \dots, x_n) \in \mathbb{R}$, $\mathcal{L}(\cdot)$ is a function of θ .
- Values of θ for which $\mathcal{L}(\theta | \mathbf{x}_n)$ is higher, are the most "compatible" with the observed sample.

The **Maximum Likelihood Estimator (MLE)** for a parameter θ , is the solution to the following maximisation problem

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} \mathcal{L}(\theta \mid x_1, \dots, x_n) \\ &= \log \left(\mathcal{L}(\theta \mid x_1, \dots, x_n) \right) \\ &= \ell(\theta \mid x_1, \dots, x_n)\end{aligned}\tag{16}$$

MAXIMUM LIKELIHOOD ESTIMATOR

Operationally, the steps to find the **MLE** are:

- 1 Compute the derivative of the log-likelihood and equate it to 0

$$\frac{d\ell(\theta \mid x_1, \dots, x_n)}{d\theta} = 0 \quad (17)$$

- 2 Isolate θ to find a candidate for the MLE
- 3 Check second derivative

- Consider a Bernoulli model for a generic sample of dimension n .

$$\mathcal{L}(\theta \mid \mathbf{x}_n) = f_n(\mathbf{x}_n \mid \theta) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i} \quad (18)$$

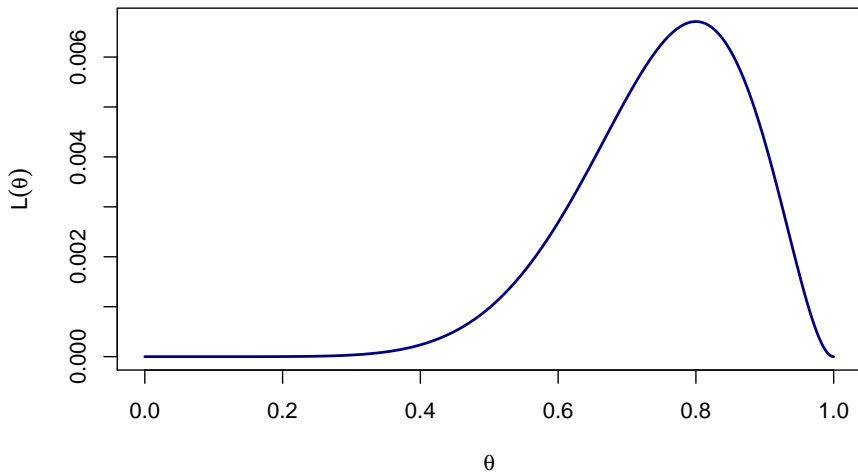
- I run the experiment 10 times and we observe $\sum_{i=1}^n x_i = 8$.

- The likelihood function is therefore

$$\mathcal{L}(\theta \mid \mathbf{x}_n) = \theta^8 (1 - \theta)^2, \quad \theta \in [0, 1] \quad (19)$$

BERNOULLI MAXIMUM LIKELIHOOD

- From visual inspection of the likelihood function



EXAMPLE

Let X_1, \dots, X_n be a random sample with $X_i \sim \text{Poisson}(\lambda)$. Then:

■ likelihood function:

$$\mathcal{L}(\lambda \mid x_1, \dots, x_n) = \frac{1}{\prod_{i=1}^n x_i!} e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \quad (20)$$

■ log-likelihood:

$$\ell(\lambda \mid x_1, \dots, x_n) = \log\left(\frac{1}{\prod_{i=1}^n x_i!}\right) - n\lambda + \log(\lambda) \sum_{i=1}^n x_i \quad (21)$$

EXAMPLE

- Compute the derivative of $\ell(\lambda | x_1, \dots, x_n)$ and equate it to 0

$$\frac{d\ell(\lambda | x_1, \dots, x_n)}{d\lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0 \quad (22)$$

- Isolate λ to get the MLE estimate

$$-n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0 \implies \hat{\lambda}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}_n \quad (23)$$

NB Even if $p_{X_1, \dots, X_n}(x_1, \dots, x_n | \lambda)$ denotes a discrete distribution, the Likelihood is still a continuous function in λ , hence it can be differentiated and maximised.

EXERCISES

- Consider a sample X_1, \dots, X_n of discrete random variables where $X_i \sim \text{Geometric}(\theta)$ with probability mass function given by

$$f_X(x|\theta) = \theta(1 - \theta)^{x-1} \quad \forall x \in \mathbb{N}_+ \text{ and } 0 < \theta < 1 \quad (24)$$

Find the MLE for θ .

- Let X_1, \dots, X_n be a random sample with pdf

$$f(x | x_0, \theta) = \theta x_0^\theta x^{-\theta-1} \quad (25)$$

Suppose $x_0 > 0$ is known and given. Find the MLE for θ .

EVALUATING POINT ESTIMATORS

Consider an estimator T for an unknown population parameter θ .

The estimator T is said to be **unbiased** if

$$\mathbb{E}[T] = \theta \quad (26)$$

The estimator T is **precise** if its variance $\mathbb{V}[T]$ is *small*

A "good" estimator is, on average, close to the real value of the parameter it is trying to recover and is always "on target".

The **Mean Squared Error** (MSE) evaluates the performance of an estimator T combining the concepts of Bias and Variance

$$MSE(T) = \mathbb{V}[T] + [Bias(T)]^2 \quad (27)$$

If $\mathbb{E}[T] = \theta$ the estimator T is unbiased and its MSE reduces to its variance.

Consistency

- the MSE can be alternatively defined as

$$MSE(T) = \mathbb{E}[(T - \theta)^2] \quad (28)$$

- when we get that

$$\lim_{n \rightarrow \infty} MSE[T] = 0 \quad (29)$$

the estimator T becomes closer and closer to the true parameter value θ as n grows. This important property is called **consistency** and it reassures us that increasing the sample size will improve the performance of our estimator