

Slides for the course

Statistics and econometrics

Part 2: Estimation

European University Institute

Andrea Ichino

September 18, 2014

Outline

The problem of estimation

The method of Maximum Likelihood

The ML estimator for the normal distribution

The Score and the Fisher Information

The method of Moments

The MM estimator for the normal distribution

A case in which the MM and ML estimators do not coincide

Small sample distribution of an estimator

Section 1

The problem of estimation

The problem of estimation

In the pre-course you have seen how to model data and phenomena with statistical distributions that depend on unknown parameters:

Example: Bernoulli distribution.

$$X = \begin{cases} 1 & \text{with probability } q \\ 0 & \text{with probability } 1 - q \end{cases} \quad (1)$$

$$X \sim p_X(X = x|q) = q^x(1 - q)^{1-x} = \begin{cases} q & \text{for } x = 1 \\ 1 - q & \text{for } x = 0 \end{cases} \quad (2)$$

How can we use a random sample of data to get information on the parameter q in a given population to study the random variable X ?

(Note: pay attention to the meanings of X , x , 1, 0, q and $p_X(\cdot)$...)

Some notational conventions

- ▶ $f_X(X = x|\theta) = f_X(x|\theta)$ is the probability density function (pdf) of the random variable X evaluated at the realization x , given the parameter θ .
- ▶ $F_X(X = x|\theta) = F_X(x|\theta)$ is the correspondent cumulative distribution.
- ▶ $p_X(\cdot)$ and $P_X(\cdot)$ are used in case it is necessary to make explicit reference to a discrete pdf.

Estimators and estimates

An *estimator* is any function of a random sample that gives information on the parameters of the distribution of a random variable.

- ▶ Analogy: the recipe for a cake ...

A random sample is a set of observed realizations of a random variable, satisfying certain properties (to be discussed below)

- ▶ Analogy: the ingredients required by the recipe ...

An *estimate* is the specific value that the estimator takes when it is evaluated using the observed realizations of a specific sample

- ▶ Analogy: the cake that you actually get when you mix the ingredients according to the recipe (hopefully a good cake ...)

How do we choose between estimators

Estimators (like recipes for cakes) have properties that:

- ▶ translate into a quality of the estimate (how good is the cake)
- ▶ for given quality of the data (how good are the ingredients).

The estimator is a random variable (with a distribution) because it is a function of the sample observations which are random variables.

The estimate is a number: a function of the sample realizations.

To choose between estimators we will have to study their properties:

- ▶ Unbiasedness
- ▶ Efficiency
- ▶ Consistency
- ▶ Other asymptotic properties

Understanding “randomness” (and independence)

If $\{X_1 \dots X_i \dots X_n\}$ are independent draws from a population with density function $f(X|\theta)$, then

- ▶ $\{X_1 = x_1 \dots X_i = x_i \dots X_n = x_n\}$ is a random sample from the population defined by $f(X|\theta)$.
- ▶ each draw X_i is a random variable and x_i is its sample realization.

Exercise questions:

- ▶ Is your class a random sample of PhD students in economics?
- ▶ Subjects with A-L names, are a random sample of Italians?
- ▶ If I toss a fair coin for each of you, those who get a tail are a random sample of the class?

Knowing that $X_i = x_i$ should not have info on the realization of X_j .

Section 2

The method of Maximum Likelihood

The likelihood function

Let $x_1, \dots, x_j, \dots, x_n$ be a random sample of size n from the random variable $X \sim f_X(x|\theta)$, where θ is an unknown (set of) parameters. The *likelihood function* is the product of the pdf evaluated at the n realizations x_j :

$$L(X|\theta) = \prod_{i=1}^n f_X(x_i|\theta) \quad (3)$$

Note that the likelihood is:

- ▶ a function of the parameter θ , given the realizations x_j ;
- ▶ the pdf of observing the sample realizations as a function of θ ;

The method of Maximum Likelihood estimates θ as the value of the parameter θ that maximizes the likelihood *given the realized sample*.

Regularity conditions

Under some regularity conditions, the ML estimator is obtained by standard maximisation of the likelihood:

- ▶ first derivative equal to zero (first order condition);
- ▶ negative second derivative (second order condition).

These conditions are:

- ▶ $f(X|\theta)$ must be
 - ▶ continuous
 - ▶ with continuous first order and second order derivatives
- ▶ the set of values X for which $f(X|\theta) \neq 0$ must not depend on θ , i.e. the support of the underlying distribution cannot depend on the parameter to be estimated.

We will see cases in which these conditions are violated and still we can derive the ML estimator of the parameter of interest.

The log-likelihood and the maximization problem

Under the regularity conditions it is simpler to maximize the log likelihood (given monotonicity of log)

$$\max_{\theta} l(X|\theta) = \ln(L(X|\theta)) = \sum_{i=1}^n \ln(f_X(x_i|\theta)) \quad (4)$$

so that the ML estimator of θ is:

$$\hat{\theta} = \mathit{arg} \max_{\theta} \left\{ \ln L(X|\theta) = \sum_{i=1}^n \ln(f_X(x_i|\theta)) \right\} \quad (5)$$

which has to satisfy the first and second order conditions

$$\frac{d \ln(L(X|\theta))}{d\theta} = \sum_{i=1}^n \frac{d \ln(f_X(x_i|\theta))}{d\theta} = \sum_{i=1}^n \frac{1}{f_X(x_i|\theta)} \frac{df_X(x_i, \theta)}{d\theta} = 0 \quad (6)$$

$$\frac{d^2 \ln(L(X|\theta))}{(d\theta)^2} < 0 \quad (7)$$

Subsection 1

The ML estimator for the normal distribution

Example: the Normal distribution

$$X \sim f_X(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (8)$$

The maximization problem is:

$$\max_{\mu, \sigma^2} L(X|\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \quad (9)$$

$$\max_{\mu, \sigma^2} l(X|\mu, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad (10)$$

The two first order condition are:

$$\frac{d \ln(L(X|\mu, \sigma^2))}{d\mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \quad (11)$$

$$\frac{d \ln(L(X|\mu, \sigma^2))}{d\sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2} \left(\frac{1}{\sigma^2} \right)^2 \sum_{i=1}^n (x_i - \mu)^2 = 0 \quad (12)$$

ML estimator for the normal distribution

Using (11), the sample mean is the ML estimator of the mean:

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n X_i \quad (13)$$

The corresponding estimate is

$$\mu_e = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \quad (14)$$

Using (12), the sample variance is the ML estimator of the variance:

$$\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})_{ML}^2 \quad (15)$$

The corresponding estimate is

$$\sigma_e^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (16)$$

Moments of the ML estimator of the mean of a normal

The ML estimator (as any estimator) is a random variable with moments:

$$E(\hat{\mu}) = \mu \quad \text{and} \quad \text{Var}(\hat{\mu}) = \frac{\sigma^2}{n} \quad (17)$$

$$E(\hat{\mu}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} n\mu = \mu \quad (18)$$

$$\begin{aligned} \text{Var}(\hat{\mu}) &= E(\hat{\mu} - \mu)^2 = E\left(\frac{1}{n} \sum_{i=1}^n X_i - \mu\right)^2 = E\left(\frac{1}{n^2} \left(\sum_{i=1}^n (X_i - \mu)\right)^2\right) \\ &= \frac{1}{n^2} E\left(\sum_{i=1}^n \sum_{j=1}^n (X_i - \mu)(X_j - \mu)\right) = \frac{n}{n^2} \sigma^2 = \frac{\sigma^2}{n} \end{aligned} \quad (19)$$

where $E(.) = n\sigma^2$ when $i = j$ and $E(.) = 0$ when $i \neq j$ because of X_i and X_j are independent.

Subsection 2

The Score and the Fisher Information

The Score of the likelihood function

The Score of the (log) likelihood $l(X|\theta)$ is the gradient (i.e. the vector of partial derivatives), with respect to the parameters θ

$$S(\theta, X) = \frac{\partial \ln(L(X|\theta))}{\partial \theta} = \frac{1}{L(X|\theta)} \frac{\partial L(X|\theta)}{\partial \theta} \quad (20)$$

The score

- ▶ measures the sensitivity of the likelihood to changes of the parameters for given X ;
- ▶ plays an important role in many applications of ML estimation.

Note that under the regularity conditions stated above, the first order condition to obtain the ML estimator can be written as:

$$S(\hat{\theta}_{ML}, X) = 0 \quad (21)$$

The mean of the score

A useful property of the score is that its mean (integrating over X at the true θ) is zero

$$\begin{aligned} E_X(S(\theta, X)) &= \int_{-\infty}^{+\infty} \frac{\partial \ln(L(X|\theta))}{\partial \theta} L(X|\theta) dX \\ &= \int_{-\infty}^{+\infty} \frac{1}{L(X|\theta)} \frac{\partial L(X|\theta)}{\partial \theta} L(X|\theta) dX \\ &= \int_{-\infty}^{+\infty} \frac{\partial L(X|\theta)}{\partial \theta} dX \\ &= \frac{\partial \left(\int_{-\infty}^{+\infty} L(X|\theta) dX \right)}{\partial \theta} = \frac{\partial(1)}{\partial \theta} = 0 \end{aligned}$$

Thus the ML estimator of θ is the value $\hat{\theta}_{ML}$ that makes the realization of the score equal to its expected value at the true θ (under regularity conditions).

An example: Binomial likelihood

Consider the likelihood of a binomial sample realisation ($n = 1$):

$$X \sim q^x(1 - q)^{1-x} = L(X|q) \quad \text{where } x = \{0, 1\}. \quad (22)$$

The log likelihood is

$$l(X|q) = x \ln(q) + (1 - x) \ln(1 - q) \quad (23)$$

and the score is

$$\frac{\partial l(X|q)}{\partial q} = \frac{x}{q} - \frac{1 - x}{1 - q} \quad (24)$$

The expected value of the score is

$$\begin{aligned} E\left(\frac{\partial l(X|q)}{\partial q}\right) &= E\left(\frac{x}{q} - \frac{1 - x}{1 - q} \mid x = 1\right)q + E\left(\frac{x}{q} - \frac{1 - x}{1 - q} \mid x = 0\right)(1 - q) \\ &= \frac{1}{q}q - \frac{1}{1 - q}(1 - q) = 0 \end{aligned} \quad (25)$$

The variance of the score: Fisher Information

Since the mean of the score is zero, its variance can be written as

$$\begin{aligned}\mathcal{I}_n(\theta) &= E_X \left((S(\theta, X))^2 \right) = E_X \left(\left(\frac{\partial \ln(L(X|\theta))}{\partial \theta} \right)^2 \right) \\ &= \int_{-\infty}^{+\infty} \left(\frac{\partial \ln(L(x|\theta))}{\partial \theta} \right)^2 L(X|\theta) dX \\ &= - \int_{-\infty}^{+\infty} \left(\frac{\partial^2 \ln(L(x|\theta))}{\partial \theta^2} \right) L(X|\theta) dX \\ &= -E_X \left(\frac{\partial^2 \ln(L(X|\theta))}{\partial \theta^2} \right)\end{aligned}\tag{26}$$

- ▶ it increases with the (absolute value of the) second derivative;
- ▶ it measures concavity, and thus how precise is ML estimation;
- ▶ the subscript n indicates that $\mathcal{I}_n(\theta)$ is the Fisher Information for the likelihood of a sample of n random variables.
- ▶ it is a $k \times k$ symmetric matrix for k parameters θ .

Equivalence of the expressions for $\mathcal{I}_n(\theta)$ in (26)

Using the simplified notation:

$$S(\theta, X) = l_\theta(\theta, X) \quad (27)$$

we can demonstrate the second to third line step in previous slide:

$$E_X \left((l_\theta(\theta, X))^2 \right) = -E_X (l_{\theta\theta}(\theta, X)) \quad (28)$$

Given

$$\int_{-\infty}^{+\infty} e^{l(\theta, X)} dX = 1 \quad (29)$$

take the derivative with respect to θ on both sides

$$\int_{-\infty}^{+\infty} l_\theta(\theta, X) e^{l(\theta, X)} dX = 0 \quad (30)$$

...

Expressions for $\mathcal{I}_n(\theta)$ in (26) (cont.)

and do it again on both sides

$$\int_{-\infty}^{+\infty} [l_{\theta\theta}(\theta, X) + (l_{\theta}(\theta, X))^2] e^{l(\theta, X)} dX = 0 \quad (31)$$

rearranging

$$-\int_{-\infty}^{+\infty} l_{\theta\theta}(\theta, X) e^{l(\theta, X)} dX = \int_{-\infty}^{+\infty} (l_{\theta}(\theta, X))^2 e^{l(\theta, X)} dX \quad (32)$$

which proves the result

$$-E_X(l_{\theta\theta}(\theta, X)) = E_X\left((l_{\theta}(\theta, X))^2\right)$$

Fisher Information for samples of size n and 1

The Fisher Information increases with the size of the sample.

A useful relationship links $\mathcal{I}_n(\theta)$ to $\mathcal{I}_1(\theta)$.

$$\mathcal{I}_n(\theta) = n\mathcal{I}_1(\theta) \quad (33)$$

where $\mathcal{I}_1(\theta)$ is the Information computed for one generic observation in the sample.

This distinction is important because, as we will show:

- ▶ $\frac{1}{\mathcal{I}_n(\theta)}$ is the variance of an unbiased ML estimator which is also the Cramer-Rao lower bound, i.e. the lowest possible variance of an unbiased estimator when the sample has size n ;
- ▶ $\frac{1}{\mathcal{I}_1(\theta)}$ is the variance of the (normal) asymptotic distribution of the ML estimator.

Example: $\mathcal{I}_n(\theta)$ for the exponential distribution

$$X \sim f_X(X|\theta) = \theta e^{-\theta X} \quad (34)$$

The maximization problem is:

$$\max_{\theta} L(X|\theta) = \prod_{i=1}^n \theta e^{-\theta X_i} \quad (35)$$

$$\max_{\theta} l(X|\theta) = +n \ln \theta - \theta \sum_{i=1}^n X_i \quad (36)$$

which leads to the first order condition and ML estimator:

$$\frac{dl(X|\theta)}{d\theta} = S(\theta, X) = \frac{n}{\theta} - \sum_{i=1}^n X_i = 0 \quad (37)$$

$$\hat{\theta}_{ML} = \frac{n}{\sum_{i=1}^n X_i} \quad (38)$$

Note: another case of biased ML estimator. What if we had specified the exponential as $X \sim f_X(X|\lambda) = \frac{1}{\lambda} e^{-\frac{1}{\lambda} X}$? (See problem set 1).

Example: $\mathcal{I}_n(\theta)$ for the exponential distribution (cont.)

The Score is

$$S(\theta, \mathbf{X}) = l_{\theta}(\theta, \mathbf{X}) = \frac{n}{\theta} - \sum_{i=1}^n X_i \quad (39)$$

The Information for the sample of size n can be computed as

$$\mathcal{I}_n(\theta) = -E_X \left(\frac{\partial^2 \ln(L(\mathbf{X}|\theta))}{\partial \theta^2} \right) = \frac{n}{\theta^2} \quad (40)$$

while the Information for a generic sample observation X_i

$$\mathcal{I}_1(\theta) = -E_X \left(\frac{\partial^2 \ln(L(X_i|\theta))}{\partial \theta^2} \right) = \frac{1}{\theta^2} \quad (41)$$

In part 3 of the slides, we will see that because of the CLT

$$\text{AsyVar}(\hat{\theta}_{ML}) = \frac{1}{\mathcal{I}_1(\theta)} = \theta^2 \quad (42)$$

Section 3

The method of Moments

A convenient method for multiple parameters

The “Methods of Moments” constructs estimators using restrictions of the population’s moments that should be satisfied also in the sample (under random sampling).

Let x_1, \dots, x_n be a random sample of n draws from the random variable $X \sim f_X(x|\theta_1, \dots, \theta_k)$, where the parameters θ are unknown.

Suppose that k moments of X exist and let the j th moment be

$$E(X^j) = g_j(\theta_1, \dots, \theta_k) \quad (43)$$

The MM estimates the k parameters as solutions of the system:

$$E(X^j) = g(\theta_1, \dots, \theta_k) = \int_{-\infty}^{\infty} X^j f_X(X|\theta_1, \dots, \theta_k) dX = \frac{1}{n} \sum_{i=1}^n x_i^j$$

with one equation for each $j = 1, \dots, k$.

Subsection 1

The MM estimator for the normal distribution

The MM estimator for the normal distribution

$$f_X(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (44)$$

The system that defines the MM estimators of μ and σ^2 is:

$$\begin{aligned} E(X) &= g_1(\hat{\mu}, \hat{\sigma}^2) = \int_{-\infty}^{\infty} X f_X(x|\hat{\mu}, \hat{\sigma}^2) dX = \frac{1}{n} \sum_{i=1}^n X_i \\ E(X^2) &= g_2(\hat{\mu}, \hat{\sigma}^2) = \int_{-\infty}^{\infty} X^2 f_X(x|\hat{\mu}, \hat{\sigma}^2) dX = \frac{1}{n} \sum_{i=1}^n X_i^2 \\ \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n X_i \end{aligned} \quad (45)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2 \quad (46)$$

The MM and ML estimators coincide in this case.

Subsection 2

A case in which the MM and ML estimators do not coincide

The Pareto distribution

Let X denote individual income. The Pareto's Law claims that

$$P(X \geq x) = \left(\frac{\nu}{X}\right)^\theta \quad \Rightarrow \quad F_X(X|\theta, \nu) = 1 - \left(\frac{\nu}{X}\right)^\theta \quad (47)$$

where ν is the (known) minimum income in the population and $\theta > 1$.

Thus by differentiation the pdf is:

$$X \sim f_X(X|\theta, \nu) = \theta \nu^\theta \left(\frac{1}{X}\right)^{\theta+1}, \quad X > \nu; \theta > 1 \quad (48)$$

The MM estimator for θ in the Pareto distribution

$$E(X) = \int_{\nu}^{\infty} x \theta \nu^{\theta} \left(\frac{1}{x}\right)^{\theta+1} dx = \theta \nu^{\theta} \int_{\nu}^{\infty} x^{-\theta} dx = \frac{\theta \nu}{\theta - 1} \quad (49)$$

The Method of Moments estimates θ solving for $\hat{\theta}$:

$$E(X) = \frac{\hat{\theta} \nu}{\hat{\theta} - 1} = \frac{1}{n} \sum_{i=1}^n (X_i) = \bar{X} \quad (50)$$

which gives the estimator

$$\hat{\theta}^{MM} = \frac{\bar{X}}{\bar{X} - \nu} \quad (51)$$

And given a random sample x_1, \dots, x_n , an estimate

$$\theta_e^{MM} = \frac{\bar{x}}{\bar{x} - \nu} \quad (52)$$

The ML estimator for θ in the Pareto distribution

The likelihood of the random sample x_1, \dots, x_n is

$$L(X|\theta, \nu) = \prod_{i=1}^n \theta \nu^\theta \left(\frac{1}{X_i} \right)^{\theta+1} \quad (53)$$

$$l(X|\theta, \nu) = \ln L(X|\theta, \nu) = n \ln \theta + n \theta \ln \nu - (\theta + 1) \sum_{i=1}^n \ln X_i \quad (54)$$

The first order condition is

$$\frac{n}{\theta} + n \ln \nu - \sum_{i=1}^n \ln X_i = 0 \quad (55)$$

Hence the estimator and estimate are :

$$\hat{\theta}^{ML} = \frac{n}{-n \ln \nu + \sum_{i=1}^n \ln X_i}; \quad \theta_e^{ML} = \frac{n}{-n \ln \nu + \sum_{i=1}^n \ln x_i} \quad (56)$$

Section 4

Small sample distribution of an estimator

An estimator is a random variable with a distribution

Estimators are random variables because they are transformations of the sample draws which are random variables.

Distributions and related moments of an estimator can thus be derived using:

- ▶ the rules for random variable transformation;
- ▶ the Moment Generating Functions;
- ▶ the Characteristic Function.

There are however cases in which the small sample distribution of an estimator cannot be derived, while using asymptotic results one can obtain the asymptotic distribution of the estimator.

This possibility is crucially useful in empirical analysis, as we will see.

Example: distribution of the sample mean of a normal

Using the *Moment Generating Function* we can show that:

$$\hat{\mu} \sim \phi\left(\hat{\mu} \mid \mu, \frac{\sigma^2}{n}\right) \quad (57)$$

where $\phi(\cdot)$ is a normal distribution with mean μ and variance $\frac{\sigma^2}{n}$.

Consider (for $t \in \mathcal{R}$) the MGFs of X_i , $S = \sum_{i=1}^n X_i$ and $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$

$$M_{X_i}(t) = E_{X_i}(e^{tX_i}) = e^{\mu t + \frac{\sigma^2 t^2}{2}} \quad (58)$$

$$M_S(t) = \prod_{i=1}^n \left(e^{\mu t + \frac{\sigma^2 t^2}{2}} \right) = e^{n\mu t + n\frac{\sigma^2 t^2}{2}} \quad (59)$$

$$M_{\hat{\mu}}(t) = M_S\left(\frac{1}{n}t\right) = e^{\mu t + \frac{\sigma^2 t^2}{n^2}} \quad (60)$$

which is the MGF of a normal with mean μ and variance $\frac{\sigma^2}{n}$.