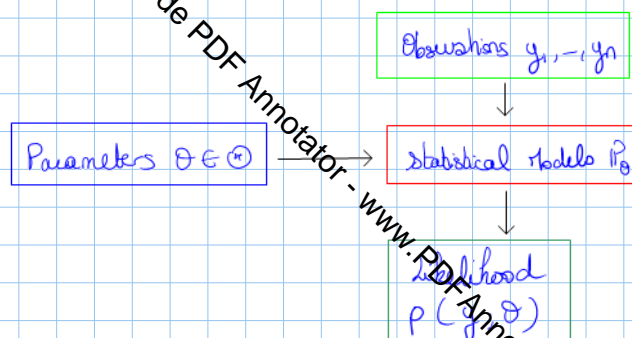


Statistical inference

To start this chapter, we will give a schematic overview of the landscape:



The questions that arise are:

- (i) Can we estimate the unknown parameter $\theta \in \Theta$ so that the model fits the data as good as possible \rightarrow Point estimation
- (ii) How can we take into account the uncertainty of this estimation to know for example if $\theta < \theta^*$ or $\theta \geq \theta^*$ where θ^* is a reference value that we know \rightarrow hypothesis testing
- (iii) Provide the precision of the estimation \rightarrow Confidence intervals.

these will be the 3 topics we will see in the following

The first notion to be exhibited in the inference issue is what is called identifiability

def: the model is said to be identifiable if $\theta_1 \neq \theta_2$ implies $P_{\theta_1} \neq P_{\theta_2}$

ex: the Exponential densities forms a family of identifiable models. So does the normals $N(\mu, \sigma^2)$ with $\theta = (\mu, \sigma^2)$

Goal of parametric inference = identify the unknown P_0 from the n sample of observations (Y_1, \dots, Y_n)

More generally, one may want to estimate a function $g(\theta) \in \mathbb{R}^k$ of $\theta \in \Theta$

Def: An estimator \hat{g} of $g(\theta)$ is a function of the n sample (Y_1, \dots, Y_n) taking its values in $g(\Theta)$. A sequence of estimators $(\hat{g}_n)_{n \geq 0}$ of $g(\theta)$ where \hat{g}_n is a function of (Y_1, \dots, Y_n) of size n is said consistent if for all $\theta \in \Theta$ we have

$$\lim_{n \rightarrow \infty} \hat{g}_n = g(\theta) \quad P_{\theta}\text{-a.e.}$$

The question now is how do we construct consistent estimators. We will see several of them in the next sections.

Remark: Suppose we have a convergent estimator $\hat{\theta}_n$ of θ . If g is continuous then $g(\hat{\theta}_n)$ is a convergent estimator of $g(\theta)$

I Method of moments

We start this section with an illustrative example.

Example: let us consider the beta distribution family $\mathcal{B} = \{ \beta(a, b), a > 0, b > 0 \}$ and $\theta = (a, b)$
 We want to estimate θ .

let $Y_1 \sim \beta(a, b)$, therefore $E_0[Y_1] = \frac{a}{a+b}$ and $E_0[Y_1(1-Y_1)] = \frac{ab}{(a+b)(a+b+1)}$

thanks to the strong law of large numbers, we can propose the following convergent estimators of a and b :

$$\hat{a}_n = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{and} \quad \hat{b}_n = \frac{1}{n} \sum_{i=1}^n Y_i(1-Y_i) \quad \text{where } Y_i \text{ i.i.d } \beta(a, b)$$

thanks to the previous remark, $\hat{a}_n = \frac{\hat{c}_n \hat{d}_n}{\hat{c}_n - \hat{d}_n - \hat{c}_n^2}$ and $\hat{b}_n = \frac{(1 - \hat{c}_n) \hat{d}_n}{\hat{c}_n - \hat{d}_n - \hat{c}_n^2}$ are convergent estimators of a and b .

the method of moments consists in finding a function m , invertible and continuous, a measurable function g s.t $E_0[g(Y)]$ and $m(\theta) = E_0[g(Y)] \quad \forall \theta \in \Theta$

def: An estimator based on the method of moments is

$$\hat{\theta}_n = m^{-1} \left(\frac{1}{n} \sum_{i=1}^n g(Y_i) \right)$$

Remark: In most of the cases, the function g is polynomial

Proposition: If $E_0[|g(Y)|] < \infty$ and m^{-1} is continuous, $\hat{\theta}_n \rightarrow \theta$ P.s.a.e

Moreover if $\forall \theta \in \Theta, E[|g(Y)|^2] < \infty$ and if m is differentiable then

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow{\mathcal{L}} \mathcal{P} \left(0, \frac{1}{m'(\theta)^2} \text{Var}_0(g(Y)) \right)$$

Proposer en exercice les cas lois exponentielle et Cauchy

This definition can be generalized to multidimensional parameters: $\theta \in \Theta \subset \mathbb{R}^d, d \geq 1$

In this case, it is most of the time impossible to identify θ with a single function g .

let $g_l: \mathbb{R} \rightarrow \mathbb{R}, l \in \{1, \dots, d\}$ and consider $y \mapsto (g_1(y), \dots, g_d(y)), y \in \mathbb{R}$

such that $m_l(\theta) = E_0(g_l(Y)) = \int_{\mathbb{R}} g_l(y) f_0(y) dy; l=1, \dots, d$ has a unique solution

def: An estimator by the method of moments associated to the $(g_l)_l$ is the solution of

$$m_l(\hat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n g_l(Y_i), \quad l=1, \dots, d$$

Notation: $\vec{m}(\theta) = E_0[\vec{g}(Y)] = (E_0[g_1(Y)], \dots, E_0[g_d(Y)])$

and $\theta = \vec{m}^{-1}(m_1(\theta), \dots, m_d(\theta))$

To estimate θ we set $\hat{\theta}_n = \vec{m}^{-1} \left(\frac{1}{n} \sum_{i=1}^n g_1(Y_i); \dots; \frac{1}{n} \sum_{i=1}^n g_d(Y_i) \right)$

Proposition: If \vec{m} is continuous, invertible with a continuous inverse, then the estimator by the method of moments is well defined and

$$\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{\text{P.s.a.e}} \theta$$

Moreover, if \vec{m}^{-1} is differentiable and $E_0[g_l(Y)^2] < \infty$

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{P}(0, V(\theta))$$

II Maximum likelihood estimator

Example: For an investigation, one is looking for a person of 1,80 meter high. The first question to ask is: Is this person a male or a female?

this looks like the following reaction



this writes: $p(1.80; \theta_H) \geq p(1.60; \theta_F)$
where $\theta_H = (\mu_H, \sigma_H^2)$ and
 $\theta_F = (\mu_F, \sigma_F^2)$

We have chosen the set of parameters which maximises the quantity $\theta \mapsto p(1,80;\theta)$

1) Definition of the maximum likelihood estimator (MLE)

def: we likelihood (function) associated to a statistical experiment \mathcal{E}^\wedge the following function:

$\theta \in \Theta \mapsto L_n(\theta; y_1, y_2, \dots, y_n) = \frac{1}{n} \sum_{i=1}^n \log f(\theta; y_i)$ (noted $L_n(\theta)$)
where $f(\theta; y)$ is the probability density function of Y (w.r.t. the Lebesgue measure)

Note that the likelihood function is a random function, depending on the observations

Example: For the normal distribution above: $\mathcal{P}_\theta = \{ \mathcal{P}(\mu, \sigma^2) ; \theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}_+^{\times} \}$

$$P_{\theta}(dy) = f(\theta, y) dy = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^2\right) dy$$

$$\text{and } \mathcal{L}_n(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right)$$

Faire en exercice le cas cauchy, Poisson (loi discrète) et mélange (cf Slac p 92-93)

def we call a maximum likelihood estimator, any estimator $\hat{\theta}_n^{MLE}$ such that

$$L_n(\hat{\theta}_n^{\text{MLE}}; y_1, \dots, y_n) = \max_{\theta \in \Theta} L_n(\theta; y_1, \dots, y_n)$$

That is to say $\hat{\theta}_n^{MLE} \in \underset{\theta \in \Theta}{\operatorname{argmax}} L_n(\theta; y_1, \dots, y_n)$

Crée avec une version d'essai de PDF Annotator - www.pdfannotator.com

Note that the MLE may not exist and may not be unique

def: application: $\theta \in \Theta \mapsto \ln(\theta; y_1, \dots, y_n) = \frac{1}{n} \log L_n(\theta; y_1, \dots, y_n)$

$$= \frac{1}{n} \sum_{i=1}^n \log f(\theta; y_i)$$

is all defined if $f(\theta, y) > 0$ and is called the log likelihood function. Imposing

$\log 0 = -\infty$, one can generalise this definition.

Note that $\hat{\theta}_n^{MLE} \in \arg\max_{\theta \in \Theta} \ln(\theta, y_1, \dots, y_n)$

2) Likelihood equation

If the maximum of $\theta \mapsto \ln(\theta)$ (or $L_n(\theta)$) is not reached on the boundary of Θ and if L_n is continuously differentiable, then a necessary condition for $\hat{\theta}_n^{MLE}$ to be the MLE is

$$\nabla_{\theta} \ln(\theta, y_1, \dots, y_n) \Big|_{\theta = \hat{\theta}_n^{MLE}} = 0 \quad \text{or equivalently} \quad \nabla_{\theta} \ln(\theta, y_1, \dots, y_n) \Big|_{\theta = \hat{\theta}_n^{MLE}} = 0$$

this is the likelihood equation.

3) Convergence:

We define the entropy of P_{θ} as $H_{\theta} = -\mathbb{E}_{\theta} [\log p(Y, \theta)]$ which write for a r.v with pdf

$$f(\theta, x) : H_{\theta} = - \int f(\theta, x) \log f(\theta, x) dx$$

This quantity measure the "disorder" of a probability measure.

Theorem: we assume that:

- (i) the model is identifiable, the set Θ is compact and $\Delta = \sup_{\theta} f(\theta, y)$ does not depend on θ
- (ii) $(\theta, y) \mapsto f(\theta, y)$ is bounded
- (iii) $\theta \mapsto f(\theta, y)$ is continuous
- (iv) H_{θ} is well defined for all $\theta \in \Theta$

then the MLE $\hat{\theta}_n^{MLE}$ is a convergent estimator of θ

Comments: this theorem is really restrictive and of course there exist other theorems that generalises this one. the thing to know is that the smoother the likelihood, the better the convergence.

4) Exemples:

a) The gaussian model

We recall that $\theta = (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+$

$$\ln(\theta, y_1, \dots, y_n) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2$$

the likelihood equation writes: $\nabla_{\mu} \ln(\theta) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu)$

$$\nabla_{\sigma^2} \ln(\theta) = -\frac{n}{2\sigma^2} + \frac{1}{4\sigma^4} \sum_{i=1}^n (y_i - \mu)^2$$

this yields : $\hat{\theta}_{MLE} = \left(\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i ; \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2 \right)$

We check that this critical point is the unique maximum

En examinant la loi log normale, des bernoulli, des uniformes pour introduire le cas des indicatrices et l'exemple où il n'y a pas de MLE. → Page p 100-101

the maximum likelihood is given by one of the maximum of a specific function - this scheme is however more general and one can define other estimators of the same form. this is what we will see in the sequel

In the same spirit, an estimator can be one zero of a given function. this is another class of estimators.

the following section introduces two categories: Π and Σ estimators.

III - 2- and Π - estimators

Problem of the method of moments: Requires an integrable function which may not be easy to find. the following classes of estimators do not depend on such a restrictive function. Moreover there exist methods to automatically select the best one (in an asymptotic way)

1) Σ - estimators:

Let first assume that $\Theta \subset \mathbb{R}$, let us denote $m_g(\theta) = \int_{\mathbb{R}} g(x) P_\theta(x) dx$, then the method of moments can be written as:

$$\int_{\mathbb{R}} (m_g(\theta) - g(x)) P_\theta(x) dx = 0 \quad \text{where } g \text{ has to be chosen.}$$

One can generalize this definition.

Let now $\Theta \subset \mathbb{R}^d$, $d \geq 1$ and $\phi: \Theta \times \mathbb{R} \rightarrow \mathbb{R}^d$ such that $\int_{\mathbb{R}} \phi_l(\theta, x) P_\theta(x) dx = 0, \forall 1 \leq l \leq d$

def: Given $\phi: \Theta \times \mathbb{R} \rightarrow \mathbb{R}^d$ as above, we call Σ -estimator associated to ϕ any estimator which solve the empirical form of the equality: that is to say:

$$\frac{1}{n} \sum_{i=1}^n \phi_l(\hat{\theta}_n, y_i) = 0 \quad \forall 1 \leq l \leq d$$

2) Π - estimators:

Let now $\Psi: \Theta \times \mathbb{R} \rightarrow \mathbb{R}$ such that $\forall \theta \in \Theta \subset \mathbb{R}^d$, $d \geq 1$ the function

$$a \mapsto \int_{\mathbb{R}} \Psi(a, y) P_\theta(y) dy \quad \text{has a maximum in } \theta \in \Theta$$

def: An Π -estimator associated to Ψ (called contrast function) is given by

$$\hat{\theta}_n \in \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \Psi(\theta, y_i)$$

Note that an Π -estimator (up to regularity conditions) is also a Σ -estimator by differentiation.

3) Convergence

We will give basic arguments on $\{P_\theta, \theta \in \Theta\}$ to ensure convergence of Σ and Π estimators. You can find more complex theorem with lighter conditions in the litterature.

See in particular the book of Van der Vaart on asymptotic Statistics.

$$n(a) = \frac{1}{n} \sum_{i=1}^n \Psi(a, y_i), \quad a \in \mathcal{A} \quad \text{and} \quad H(a, \vartheta) = \mathbb{E}_{\vartheta} [\Psi(a, Y)]$$

Proposition: Assume $\odot \subset \mathbb{R}^d$, $d \geq 1$ and $\hat{\Theta}_n$ the n estimator is well defined. Assume in addition

l'essai de (ii) $\sup_{a \in \mathcal{O}} |\eta_n(a) - \eta(a, \theta)| \xrightarrow{P_0} 0$

$$\forall \epsilon > 0 \quad \sup_{|a-\theta| > \epsilon} V(a, \theta) < V(\theta, \theta) \quad (\text{condition of maximum})$$

(iii) $\mathbb{P}_n(\hat{\theta}_n) \geq \Gamma_n(\theta) - \varepsilon_n$ with $\varepsilon_n \xrightarrow{\mathbb{P}_0} 0$

then the n -coherent $\hat{\Theta}_n$ is consistent (or convergent): $\hat{\Theta}_n \xrightarrow{P_\theta} \Theta$

ii) 2-estimators

Let $Z_n(a) = \frac{1}{n} \sum_{i=1}^n \phi(a, y_i)$, $a \in \Theta$ and $Z(a, \mathcal{D}) = \mathbb{E}_{\mathcal{D}}(\phi(a, y))$, $a \in \Theta$

Proposition: We assume that the 2-estimate is well defined and:

(i) $\sup_{a \in \mathcal{Q}} \|Z_n(a) - Z(a, \theta)\| \xrightarrow{P_0} 0$ ($\|\cdot\|$ is the Euclidean norm)

$$(ii) \quad \forall \varepsilon > 0 \quad \inf_{|a-a| \geq \varepsilon} \|Z(a, \theta)\| > 0 = \|Z(0, \theta)\|$$
$$(iii) \quad Z_n(\hat{\Theta}_n) \xrightarrow{P_\theta} 0$$

then $\hat{\Theta}_n$ is constant (or convergent) i.e. $\hat{\Theta}_n \xrightarrow{P_0} \Theta$

One can also prove Central limit theorems to have the speed of convergence. \rightarrow analyse asymptotique?

IV which estimator is the best?

The standard approach to making such judgments is called decision theory.

For estimation, this approach begins with the choice of a loss function L such that $L(\theta, d)$ quantifies the loss associated with estimating $g(\theta)$ by the value d . It is natural to assume that $L(\theta, g(\theta)) = 0$ so that there is no loss with the correct answer. Moreover we assume also $L(\theta, d) \geq 0$ if θ and d .

Because Y is random, $L(0, \hat{\sigma}(Y))$ is also random and can be large even if the estimator is excellent. Therefore, to judge the optimality of an estimator, one has to think in average.

def: Let $R(\theta, \delta) = \mathbb{E}_\theta [L(\theta, \delta(Y))]$. R is called the risk function

Example: $Y \sim \text{Bin}(100, \theta)$, $\theta \in [0, 1]$

this is equivalent of n coin tosses with a coin with $\theta = P(\text{head})$
therefore an estimator of θ is $\hat{\theta}(Y) = \frac{Y}{100}$

If we consider the loss function given by $L(\theta, d) = (\theta - d)^2$ which is the quadratic loss, then $R(\theta, \delta) = \mathbb{E}_\theta \left[\left(\theta - \frac{Y}{100} \right)^2 \right] = \frac{\sigma^2(1-\theta)}{100}$

Cr    avec une version d'essai de PDF Annotator - www.pdfannotator.com

The risk is of the following form



$$R(\theta, \delta_1)$$

$$R(\theta, \delta_0)$$

$$R(\theta, \delta_2)$$

Now the question is: can we compare 2 estimators?

Following our example, one can construct other estimators: let $\delta_1(Y) = \frac{Y+3}{100}$ and $\delta_2(Y) = \frac{Y+3}{106}$

then $R(\theta, \delta_0) = \frac{\theta(1-\theta)}{100}$; $R(\theta, \delta_1) = \frac{9 + 100\theta(1-\theta)}{100^2}$ and $R(\theta, \delta_2) = \frac{(9-8\theta)(1+8\theta)}{106^2}$

If we plot these risks, one can see that δ_0 and δ_2 are better than δ_1 . However comparing δ_0 and δ_2 depends on the value of θ . (\rightarrow explain the plot)

However, as one does not know θ , one cannot choose the best estimator. This leads us to the following definition:

def: δ_1 is better to δ_2 if: for all $\theta \in \Theta$ $R(\theta, \delta_1) \leq R(\theta, \delta_2)$

If in addition $R(\theta, \delta_1) \neq R(\theta, \delta_2)$ for at least one $\theta \in \Theta$, we say that δ_2 is inadmissible.

As one can see in our example, there are not necessary better estimators in the whole family. However restricting to a subfamily, one may get a better estimator.

Remark: In 1-D: $R(\theta, \delta) = \text{Var}_\theta(\delta) + (\mathbb{E}_\theta[\delta] - g(\theta))^2$ for the quadratic case.

therefore, as we aim at minimizing $R(\theta, \delta)$ for all θ , it looks natural to choose an estimator δ such that $\mathbb{E}_\theta(\delta) = g(\theta)$ and with a minimal variance.

def: An estimator, δ , of $g(\theta)$ is integrable if $\mathbb{E}_\theta[\delta] < \infty$ for all $\theta \in \Theta$. It is unbiased if it is integrable and $\mathbb{E}_\theta[\delta] = g(\theta) \forall \theta \in \Theta$.

Note that if δ is unbiased and $\mathbb{E}_\theta[\delta^2] < \infty \forall \theta \in \Theta$, then $R(\theta, \delta) = \text{Var}_\theta(\delta)$

Example: If $X \in L^1(\mathbb{P}_\theta) \forall \theta \in \Theta$ then $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is an unbiased estimator of $\mathbb{E}_\theta[X]$ (thanks to the strong law of large numbers)

• $\mathcal{B} = \{U_{[0,\theta]}; \theta > 0\}$ then $\hat{\Theta}_n = \max_{1 \leq i \leq n} X_i$ is a consistent estimator of θ . However it is biased.

Indeed, since the X_i 's are independent: $\forall x \geq 0$

$$\mathbb{P}_\theta(\hat{\Theta}_n \leq x) = \mathbb{P}_\theta(\forall i \in \{1, \dots, n\} X_i \leq x) = \mathbb{P}(X_1 \leq x)^n$$

$$= \left(\frac{1}{\theta}\right)^n \mathbb{P}_\theta(X_1 \leq x) = \mathbb{P}(X_1 \leq x)^n$$

$$= \frac{1}{\theta^n} x^n \mathbb{1}_{[0,\theta]}(x)$$

this yields $\hat{\Theta}_n$ has a density function given by $\frac{n}{\theta^n} x^{n-1} \mathbb{1}_{[0,\theta]}(x)$

$$\text{and } \mathbb{E}_\theta[\hat{\Theta}_n] = \frac{n}{n+1} \theta$$

Demander à Eric pertinence de ce §

V Bayesian estimation:

- Preret: utilise données pour ajouter info
- Sebrubage la suite
- Introduire pdfs: $L^1 \rightarrow \text{MAP}$
 $L^2 \rightarrow \text{IE a posteriori}$

VI Unbiased Estimation

Jf time (ask Eric)

1) Minimum Variance Unbiased estimators

- UMVU
- Uniqueness of UMVU when T is a complete sufficient statistic.
- ex

2) why would not considering biased estimators?

- Explain thanks to an example the tradeoff bias-variance.