

Rappel: Estimation paramétrique: X_1, \dots, X_n une suite i.i.d. de variables/vecteurs aléatoires, on suppose que $X_i \sim \mathbb{P}_\theta$, $(\mathbb{P}_\theta)_{\theta \in \Theta}$ une famille de mesures de proba.

Expl: $\Theta = \mathbb{R} \times \mathbb{R}_+$, $\theta = (\mu, \sigma^2)$, $\mathbb{P}_{(\mu, \sigma^2)} = \mathcal{N}(\mu, \sigma^2)$

But: trouver θ ou $f(\theta)$ ($f: \Theta \rightarrow \mathbb{R}$), à partir d'une réalisation x_1, \dots, x_n de X_1, \dots, X_n

↳ Estimateurs: fonction de X_1, \dots, X_n qui ne dépend pas de θ .

Expl: i) moyenne empirique: $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

↳ estimateur convergent de $E[X_1]$.

ii) variance empirique: $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$

↳ estimateur convergent de $\text{Var}(X_1)$.

Méthode des moments: Si on peut trouver h, g deux fonctions telles que

$$g(\theta) = h(E[g(X_1)])$$

$$f(\theta) = h(E[g(x_i)])$$

↳ on construit un estimateur de $f(\theta)$

$$\text{via } \hat{f}_n(x_1, \dots, x_n) = h\left(\frac{1}{n} \sum_{i=1}^n g(x_i)\right).$$

Estimateur du maximum de vraisemblance:

Déf: Soit $n \geq 1$, X_1, \dots, X_n un échantillon de loi \mathbb{P}_θ , $\theta \in \Theta$. Soit x_1, \dots, x_n une réalisation de X_1, \dots, X_n , on définit la vraisemblance de x_1, \dots, x_n sachant θ par :

① Si \mathbb{P}_θ est une loi discrète:

$$\begin{aligned} \mathcal{L}(x_1, \dots, x_n | \theta) &= \mathbb{P}((X_1, \dots, X_n) = (x_1, \dots, x_n)) \\ &= \prod_{i=1}^n \mathbb{P}(X_i = x_i) \\ &= \prod_{i=1}^n \mathbb{P}_\theta(x_i) \end{aligned}$$

② Si \mathbb{P}_θ est une loi continue avec densité f_θ ,

$$\mathcal{L}(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f_\theta(x_i)$$

Déf: L'estimateur du maximum de vraisemblance

est l'estimateur n. A d. -

est l'estimateur de θ donné par

$$\text{MLE}(X_1, \dots, X_n) = \arg \max_{\theta \in \Theta} \mathcal{L}(X_1, \dots, X_n | \theta),$$

la valeur de θ pour laquelle

$\mathcal{L}(X_1, \dots, X_n | \theta)$ est maximale.

Ex 1 $\Theta = \mathbb{R} \times \mathbb{R}_+$, $\theta = (\mu, \sigma^2)$, $\mathbb{P}_{(\mu, \sigma^2)} = \mathcal{N}(\mu, \sigma^2)$

$\hookrightarrow \mathbb{P}_{(\mu, \sigma^2)}$ est une loi continue de densité

$$f_{(\mu, \sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Si X_1, \dots, X_n est un échantillon de loi $\mathbb{P}_{(\mu, \sigma^2)}$,
et $x_1, \dots, x_n \in \mathbb{R}$, une réalisation de X_1, \dots, X_n ,

$$\begin{aligned} \mathcal{L}(x_1, \dots, x_n | (\mu, \sigma^2)) &= \prod_{i=1}^n f_{(\mu, \sigma^2)}(x_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right) \end{aligned}$$

maximiser \mathcal{L} est équivalent à maximiser

$\ln(\mathcal{L})$:

$$\begin{aligned} \ln(\mathcal{L}(x_1, \dots, x_n | (\mu, \sigma^2))) &= \sum_{i=1}^n \left(\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{(x_i-\mu)^2}{2\sigma^2} \right) \\ &= \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2} \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \\
&= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2
\end{aligned}$$

Pour $\sigma > 0$ fixé, le maximum en μ est atteint au point qui minimise $\sum_{i=1}^n (x_i - \mu)^2$

$x_i^2 - 2x_i\mu + \mu^2$

\leadsto strict. convexe \Rightarrow unique min atteint au pt critique :

$$\begin{aligned}
\frac{d}{d\mu} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2) &= \sum_{i=1}^n (2\mu - 2x_i) \\
&= 2n\mu - 2 \sum_{i=1}^n x_i \stackrel{!}{=} 0
\end{aligned}$$

$$\hookrightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

min atteint en $\mu = \frac{1}{n} \sum_{i=1}^n x_i \equiv \bar{x}_n$

On doit trouver σ^2 qui maximise

$$-\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

\leadsto même type de dérivation (pt. crit.)

on trouve que le max est atteint en

$$\sigma_{\hat{x}}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \quad \text{la variance empirique}$$

On trouve donc

On trouve donc

$$MLE(X_1, \dots, X_n) = (\bar{X}_n, \hat{\sigma}_n^2) \text{ dans ce cas.}$$

Expl 2

Cas $X_1, \dots, X_n \sim \text{Geo}(p)$ ($\Theta = [0, 1]$, $\theta = p$)
 $\mathbb{P}_p = \text{Geo}(p)$

Pour $x_1, \dots, x_n \in \mathbb{N}^*$ une réal. de X_1, \dots, X_n ,
on calcule la vraisemblance :

$$\begin{aligned} \mathcal{L}(x_1, \dots, x_n | p) &= \prod_{i=1}^n \mathbb{P}(X_i = x_i) \\ &= \prod_{i=1}^n (p(1-p)^{x_i-1}) \\ &= \left(\frac{p}{1-p}\right)^n \prod_{i=1}^n (1-p)^{x_i} = \left(\frac{p}{1-p}\right)^n (1-p)^{\sum_{i=1}^n x_i} \end{aligned}$$

on cherche la valeur p pour laquelle
 $\mathcal{L}(x_1, \dots, x_n | p)$ est maximale. On va
maximiser $\ln(\mathcal{L})$ à la place, et on va
poser

$$p = \frac{e^\beta}{1+e^\beta} \Rightarrow 1-p = \frac{1}{1+e^\beta}, \quad \frac{p}{1-p} = e^\beta$$

$$\beta \in \mathbb{R}.$$

$$\ln(\mathcal{L}(x_1, \dots, x_n | p)) = n \ln\left(\frac{p}{1-p}\right) + \ln(1-p) \left(\sum_{i=1}^n x_i\right)$$

$$\begin{aligned} \ln(\mathcal{L}(x_1, \dots, x_n | p)) &= \ln\left(\frac{1}{1-p}\right) + \ln(1-p)^{\sum_{i=1}^n x_i} \\ &= n\beta - \ln(1+e^\beta) \sum_{i=1}^n x_i \end{aligned}$$

↳ méthode de pt critique, on dérive p.r. à β :

$$n - \frac{e^\beta}{1+e^\beta} \sum_{i=1}^n x_i = 0$$

$$\Leftrightarrow n = p \sum_{i=1}^n x_i \quad \Leftrightarrow p = \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i} = \frac{1}{\bar{x}_n}$$

On a $\text{MLE}(X_1, \dots, X_n) = \frac{1}{\bar{X}_n}$ dans ce cas.

Intervalles de confiance:

Déf: Soit X_1, \dots, X_n un échantillon de loi $\mathbb{P}_\theta, \theta \in \Theta$,

et soit $f(\theta)$ une quantité à estimer. Un

intervalle de confiance est un intervalle

$I(X_1, \dots, X_n)$ qui ne dépend pas de θ .

Pour $\alpha \in (0, 1)$, $I(X_1, \dots, X_n)$ est un intervalle de confiance de niveau $1-\alpha$ pour $f(\theta)$ si

$$\mathbb{P}(f(\theta) \in I(X_1, \dots, X_n)) = 1 - \alpha \quad \forall \theta \in \Theta.$$

(si $\mathbb{P}(f(\theta) \in I(X_1, \dots, X_n)) \geq 1 - \alpha$, on parlera d'intervalle de confiance pas excès.)

de confiance pas excès.)

Expl: X_1, \dots, X_n est un échantillon de loi $\mathcal{N}(\mu, 1)$
($\Theta = \mathbb{R}$, $\theta = \mu$, $\mathbb{P}_\mu = \mathcal{N}(\mu, 1)$).

On cherche un intervalle de confiance au niveau $1-\alpha$ pour $\mu = E[\bar{X}_n]$.

\bar{X}_n est. de $\mu \rightsquigarrow$ on peut chercher un intervalle de la forme

$$I(X_1, \dots, X_n) = [\bar{X}_n - \delta, \bar{X}_n + \delta]$$

pour $\delta > 0$. On cherche la valeur de δ qui donne un niveau de confiance $1-\alpha$.

$$\begin{aligned} \mathbb{P}(\mu \in I(X_1, \dots, X_n)) &= \mathbb{P}(\mu \in [\bar{X}_n - \delta, \bar{X}_n + \delta]) \\ &= \mathbb{P}(|\bar{X}_n - \mu| \leq \delta) \\ &= 1 - \mathbb{P}(|\bar{X}_n - \mu| > \delta) \quad \otimes \end{aligned}$$

$\rightsquigarrow X_1, \dots, X_n \sim \mathcal{N}(\mu, 1)$ indép. on a

$$\begin{aligned} \bar{X}_n &= \frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\frac{n\mu}{n}, \frac{n}{n^2}\right) \quad (\text{lemme 2.6.7}) \\ &= \mathcal{N}\left(\mu, \frac{1}{n}\right) \end{aligned}$$

$$\Rightarrow \bar{X}_n - \mu \sim \mathcal{N}(0, \frac{1}{n})$$

$$\Rightarrow \sqrt{n}(\bar{X}_n - \mu) \sim \mathcal{N}(0, 1)$$

$$\textcircled{*} = 1 - \mathbb{P}(\sqrt{n}|\bar{X}_n - \mu| > \sqrt{n}S)$$

$$= 1 - \mathbb{P}(|Z| > \sqrt{n}S) \quad \text{avec } Z \sim \mathcal{N}(0, 1).$$

On voulait $\mathbb{P}(\mu \in I(x_1, \dots, x_n))$

$$= 1 - \mathbb{P}(|Z| > \sqrt{n}S) = 1 - \alpha$$

\Rightarrow on cherche $S > 0$ t.q.

$$\mathbb{P}(|Z| > \sqrt{n}S) = \alpha$$

"

$$2\mathbb{P}(Z > \sqrt{n}S) = 2\left(\frac{1}{2} - \mathbb{P}(Z \in [c, \sqrt{n}S])\right)$$

\leadsto il existe des tables des nombres a, b tels que

$$\mathbb{P}(Z \in [c, a]) = b$$

Si on prend $\alpha = 10\% = 0.1$, on cherche

S t.q.

$$\mathbb{P}(Z \in [c, \sqrt{n}S]) = 45\%$$

$$\Leftrightarrow S\sqrt{n} = 1.64 \quad \Leftrightarrow S = \frac{1.64}{\sqrt{n}}$$