MAP 433 : Introduction aux méthodes statistiques. Cours 4

18 Septembre 2015

Aujourd'hui

- 1 M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance
- **2** EMV, asymptotique des Z- et M- estimateurs
 - Approche asymptotique
- 3 Modèles réguliers et information de Fisher
 - Construction de l'information de Fisher
 - Modèle régulier
 - Cadre général et interprétation géométrique

M-estimation

- <u>Situation</u>: on observe X_1, \ldots, X_n de loi \mathbb{P}_{θ} sur \mathbb{R} et $\theta \in \Theta$.
- Principe : Se donner une application $\psi: \Theta \times \mathbb{R} \to \mathbb{R}_+$ telle que, pour tout $\theta \in \Theta \subset \mathbb{R}^d$,

$$\vartheta \leadsto \mathbb{E}_{\theta}\left[\psi(\vartheta,X)\right] = \int \psi(\vartheta,x) \, \mathbb{P}_{\theta}(dx)$$

admet un extremum (maximum ou minimum) en $\vartheta = \theta$.

Definition

On appelle M-estimateur associé à ψ tout estimateur $\widehat{\theta}_n$ satisfaisant

$$\sum_{i=1}^{n} \psi(\widehat{\theta}_{n}, X_{i}) = \max_{\vartheta \in \Theta} \sum_{i=1}^{n} \psi(\vartheta, X_{i}).$$

Au lieu de maximiser, on peut aussi minimiser

Un exemple classique : paramètre de localisation

■ $\Theta = \mathbb{R}$, $\mathbb{P}_{\theta}(dx) = f(x - \theta)dx$, et $\int_{\mathbb{R}} xf(x)dx = 0$, $\int_{\mathbb{R}} x^2 \mathbb{P}_{\theta}(dx) < +\infty$ pour tout $\theta \in \mathbb{R}$. On pose

$$\psi(\vartheta,x)=(\vartheta-x)^2$$

La fonction

$$\vartheta \leadsto \mathbb{E}_{\theta} \left[\psi(\vartheta, X) \right] = \int_{\mathbb{R}} (\vartheta - x)^2 f(x - \theta) dx$$

admet un maximum en $\theta = \mathbb{E}_{\theta} [X] = \int_{\mathbb{D}} x f(x - \theta) dx = \theta.$

■ *M*-estimateur associé :

$$\sum_{i=1}^{n} (X_i - \widehat{\theta}_n)^2 = \min_{\vartheta \in \mathbb{R}} \sum_{i=1}^{n} (X_i - \vartheta)^2.$$

Paramètre de localisation

■ C'est aussi un Z-estimateur associé à $\phi(\vartheta,x)=2(x-\vartheta)$: on résout

$$\sum_{i=1}^{n} (X_i - \vartheta) = 0 \text{ d'où } \widehat{\theta}_n = \overline{X}_n.$$

- Dans cet exemple très simple, tous les points de vue coïncident.
- Si, dans le même contexte, $\int_{\mathbb{R}} x^2 \mathbb{P}_{\theta}(dx) = +\infty$ et f(x) = f(-x), on peut utiliser Z-estimateur avec $\phi(\vartheta, x) = \text{Arctg}(x \vartheta)$. Méthode robuste, mais est-elle optimale? Peut-on faire mieux si f est connue? A suivre...

Lien entre Z- et M- estimateurs

- Pas d'inclusion entre ces deux classes d'estimateurs en général :
 - lacksquare Si ψ non-régulière, M-estimateur \Rightarrow Z-estimateur
 - Si une équation d'estimation admet plusieurs solutions distinctes, Z-estimateur ⇒ M-estimateur (cas d'un extremum local).
- Toutefois, si ψ est régulière, les M-estimateurs sont des Z-estimateurs : si $\Theta \subset \mathbb{R}$ (d=1), en posant

$$\phi(\vartheta, x) = \partial_{\theta} \psi(\vartheta, x),$$

on a

$$\sum_{i=1}^n \partial_{\theta} \psi(\vartheta, X_i)\big|_{\vartheta = \widehat{\theta}_n} = \sum_{i=1}^n \phi(\widehat{\theta}_n, X_i) = 0.$$

Maximum de vraisemblance

- Principe fondamental et incontournable en statistique. Cas particuliers connus depuis le XVIIIème siècle. Définition générale : Fisher (1922).
- Fournit une première méthode systématique de construction d'un *M*-estimateur
- Procédure optimale (dans quel sens?) sous des hypothèses de régularité de la famille $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$.
- Parfois difficile à mettre en oeuvre en pratique → méthodes numériques, statistique computationnelle.

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

Fonction de vraisemblance

■ La famille $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$ est dominée par une mesure σ -finie μ . On se donne, pour $\theta \in \Theta$

$$f(\theta, x) = \frac{d \mathbb{P}_{\theta}}{d\mu}(x), \ x \in \mathbb{R}.$$

Fonction de vraisemblance du *n*-échantillon associée à la famille $\{f(\theta,\cdot), \theta \in \Theta\}$:

$$\theta \rightsquigarrow \mathcal{L}_n(\theta, X_1, \dots, X_n) = \prod_{i=1}^n f(\theta, X_i)$$

• C'est une fonction aléatoire (définie μ -presque partout).

Exemples

■ Exemple 1 : Modèle de Poisson. On observe

$$X_1, \ldots, X_n \sim_{\text{i.i.d.}} \text{Poisson}(\theta),$$

$$\theta \in \Theta = \mathbb{R}_+ \setminus \{0\}$$
 et prenons $\mu(dx) = \sum_{k \in \mathbb{N}} \delta_k(dx)$.

lacksquare La densité de $\mathbb{P}_{ heta}$ par rapport à μ est

$$f(\theta, x) = e^{-\theta} \frac{\theta^x}{x!}, \quad x = 0, 1, 2, \dots$$

■ La fonction de vraisemblance associée s'écrit

$$\theta \rightsquigarrow \mathcal{L}_n(\theta, X_1, \dots, X_n) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{X_i}}{X_i!}$$
$$= \frac{1}{\prod_{i=1}^n X_i!} e^{-n\theta} \theta^{\sum_{i=1}^n X_i}.$$

Principe de maximum de vraisemblance

Exemples

■ Exemple 2 Modèle de Cauchy. On observe

$$X_1, \ldots, X_n \sim_{\text{i.i.d.}} \text{Cauchy},$$

$$\theta \in \Theta = \mathbb{R}$$
 et $\mu(dx) = dx$ (par exemple).

On a alors

$$\mathbb{P}_{\theta}(dx) = f(\theta, x)dx = \frac{1}{\pi(1 + (x - \theta)^2)}dx.$$

La fonction de vraisemblance associée s'écrit

$$\theta \rightsquigarrow \mathcal{L}_n(\theta, X_1, \dots, X_n) = \frac{1}{\pi^n} \prod_{i=1}^n \left(1 + (X_i - \theta)^2\right)^{-1}.$$

Estimateur du maximum de vraisemblance

- On généralise le principe précédent pour une famille de lois et un ensemble de paramètres quelconques.
- <u>Situation</u>: $X_1, \ldots, X_n \sim_{\text{i.i.d.}} \mathbb{P}_{\theta}$, $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$ dominée, $\Theta \subset \mathbb{R}^d$, $\theta \leadsto \mathcal{L}_n(\theta, X_1, \ldots, X_n)$ vraisemblance associée.

Definition

On appelle estimateur du maximum de vraisemblance tout estimateur $\widehat{\theta}_n^{mv}$ satisfaisant

$$\mathcal{L}_n(\widehat{\theta}_n^{\,\mathrm{mv}}, X_1, \dots, X_n) = \max_{\theta \in \Theta} \mathcal{L}_n(\theta, X_1, \dots, X_n).$$

Existence, unicité...

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

Remarques

■ Log-vraisemblance :

$$\theta \leadsto \ell_n(\theta, X_1, \dots, X_n) = n^{-1} \log \mathcal{L}_n(\theta, X_1, \dots, X_n)$$
$$= n^{-1} \sum_{i=1}^n \log f(\theta, X_i).$$

Bien défini si $f(\theta, \cdot) > 0$ μ -pp.

Max. vraisemblance = max. log-vraisemblance.

- L'estimateur du maximum de vraisemblance ne dépend pas du choix de la mesure dominante μ .
- Racine de l'équation de vraisemblance : tout estimateur $\widehat{\theta}_n^{rv}$ vérifiant

$$\nabla_{\theta}\ell_n(\widehat{\theta}_n^{\text{rv}}, X_1, \dots, X_n) = 0.$$

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

Exemple: modèle normal

L'expérience statistique est engendrée par un n-échantillon de loi $\mathcal{N}(\mu, \sigma^2)$, le paramètre est $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+ \setminus \{0\}$.

Vraisemblance

$$\mathcal{L}_n((\mu, \sigma^2), X_1, \dots, X_n) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right).$$

■ Log-vraisemblance

$$\ell_n((\mu, \sigma^2), X_1, \dots, X_n) = -\frac{n}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^n (X_i - \mu)^2.$$

Principe de maximum de vraisemblance

Exemple: modèle normal

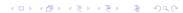
Equation(s) de vraisemblance

$$\begin{cases} \partial_{\mu}\ell_{n}((\mu,\sigma^{2}),X_{1},\ldots,X_{n}) & = & \frac{1}{\sigma^{2}}\sum_{i=1}^{n}(X_{i}-\mu) \\ \\ \partial_{\sigma^{2}}\ell_{n}((\mu,\sigma^{2}),X_{1},\ldots,X_{n}) & = & -\frac{n}{2\sigma^{2}}+\frac{1}{2\sigma^{4}}\sum_{i=1}^{n}(X_{i}-\mu)^{2}. \end{cases}$$

Solution de ces équations (pour $n \ge 2$) :

$$\widehat{\theta}_n^{\text{rv}} = (\overline{X}_n, \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2)$$

et on vérifie que $\widehat{\theta}_n^{\text{rv}} = \widehat{\theta}_n^{\text{mv}}$.



Principe de maximum de vraisemblance

Exemple : modèle de Poisson

Vraisemblance

$$\mathcal{L}_n(\theta, X_1, \dots, X_n) = \frac{1}{\prod_{i=1}^n X_i!} e^{-n\theta} \theta^{\sum_{i=1}^n X_i}.$$

Log-vraisemblance

$$\ell_n(\theta, X_1, \ldots, X_n) = c(X_1, \ldots, X_n) - n\theta + \sum_{i=1}^n X_i \log \theta.$$

Equation de vraisemblance

$$-n + \sum_{i=1}^{n} X_i \frac{1}{\theta} = 0$$
, soit $\widehat{\theta}_n^{\text{rv}} = \frac{1}{n} \sum_{i=1}^{n} X_i = \overline{X}_n$

et on vérifie que $\widehat{\theta}_n^{\text{rv}} = \widehat{\theta}_n^{\text{mv}}$.

- -M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

Exemple : modèle de Laplace

L'expérience statistique est engendrée par un n-échantillon de loi de Laplace de paramètre $\theta \in \Theta = \mathbb{R}$. La densité par rapport à la mesure de Lebesgue :

$$f(\theta, x) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \theta|}{\sigma}\right),$$

où $\sigma > 0$ est connu.

Vraisemblance

$$\mathcal{L}_n(\theta, X_1, \dots, X_n) = (2\sigma)^{-n} \exp\left(-\frac{1}{\sigma} \sum_{i=1}^n |X_i - \theta|\right)$$

Log-vraisemblance

$$\ell_n(\theta, X_1, \ldots, X_n) = -n \log(2\sigma) - \frac{1}{\sigma} \sum_{i=1}^n |X_i - \theta|.$$

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

Exemple : modèle de Laplace

Maximiser $\mathcal{L}_n(\theta, X_1, \dots, X_n)$ revient à minimiser la fonction $\theta \leadsto \sum_{i=1}^n \left| X_i - \theta \right|$, dérivable presque partout de dérivée constante par morceaux. Equation de vraisemblance :

$$\sum_{i=1}^n \operatorname{sign}(X_i - \theta) = 0.$$

Soit $X_{(1)} \leq \ldots \leq X_{(n)}$ la statistique d'ordre.

- n pair : $\widehat{\theta}_n^{\text{mv}}$ n'est pas unique; tout point de l'intervalle $\left[X_{\left(\frac{n}{2}\right)}, X_{\left(\frac{n}{2}+1\right)}\right]$ est un EMV.
- n impair : $\widehat{\theta}_{\mathbf{n}}^{\,\mathrm{mv}} = X_{\left(\frac{n+1}{2}\right)}$, l'EMV est unique. Mais $\widehat{\theta}_{\mathbf{n}}^{\,\mathrm{rv}}$ n'existe pas.
- **pour tout** n, la médiane empirique est un EMV.

Principe de maximum de vraisemblance

Exemple : modèle de Cauchy

Vraisemblance

$$\mathcal{L}_n(\theta, X_1, \dots, X_n) = \pi^{-n} \prod_{i=1}^n \frac{1}{1 + (X_i - \theta)^2}.$$

Log-vraisemblance

$$\ell_n(\theta, X_1, \dots, X_n) = -n \log \pi - \sum_{i=1}^n \log \left(1 + (X_i - \theta)^2\right)$$

Equation de vraisemblance

$$\sum_{i=1}^n \frac{X_i - \theta}{1 + (X_i - \theta)^2} = 0$$

pas de solution explicite et admet en général plusieurs solutions.

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance

l'EMV est un M-estimateur

On pose

$$\psi(\vartheta,x) := \log f(\vartheta,x), \ \ \vartheta \in \Theta, \ x \in \mathbb{R}$$

(on suppose que $f(\vartheta, \cdot) > 0$.)

La fonction

$$a \rightsquigarrow \mathbb{E}_{\theta} \left[\psi(\vartheta, X) \right] = \int_{\mathbb{R}} \log f(\vartheta, x) f(\theta, x) \mu(dx)$$

a un maximum en $\theta = \theta$ d'après l'inégalité de convexité.

- M-estimation, rappel du Cours 3
 - Principe de maximum de vraisemblance
 - lacktriangle Le \emph{M} -estimateur associé à ψ maximise la fonction

$$\vartheta \leadsto \sum_{i=1}^n \log f(\vartheta, X_i) = \ell_n(\vartheta, X_1, \dots, X_n)$$

c'est-à-dire la log-vraisemblance. C'est l'estimateur du maximum de vraisemblance.

■ C'est aussi un Z-estimateur si la fonction $\theta \leadsto \log f(\theta, \cdot)$ est régulière, associé à la fonction

$$\phi(\theta, x) = \partial_{\theta} \log f(\theta, x) = \frac{\partial_{\theta} f(\theta, x)}{f(\theta, x)}, \ \theta \in \Theta, x \in \mathbb{R}$$

lorsque $\Theta \subset \mathbb{R}$, à condition que le maximum de log-vraisemblance n'est pas atteint sur la frontière de Θ . (Se généralise en dimension d.)

◆□ → ◆□ → ◆□ → ◆□ → □ □

Asymptotique des Z- et M-estimateurs

- Problème général délicat. Dans ce cours : conditions suffisantes.
- Convergence : critère simple pour les *M*-estimateurs.
- Vitesse de convergence : technique simple pour les Z-estimateurs, à condition de savoir que l'estimateur est convergent.
- Sous des hypothèses de régularité, un *M*-estimateur est un Z-estimateur.

Convergence des *M*-estimateurs

- Situation: on observe $X_1, ..., X_n$ i.i.d. de loi dans la famille $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$.
- $\psi: \Theta \times \mathbb{R} \to \mathbb{R}$ fonction de contraste.
- Loi des grands nombres :

$$M_n(\vartheta) = \frac{1}{n} \sum_{i=1}^n \psi(\vartheta, X_i)$$

converge en \mathbb{P}_{θ} -probabilité vers

$$M(\vartheta,\theta) = \mathbb{E}_{\theta} \left[\psi(\vartheta, X) \right]$$

qui atteint son maximum en $\vartheta=\theta$

■ « à montrer » :

$$\widehat{\theta}_n = \arg\max_{\vartheta \in \Theta} \, M_{\mathbf{n}}(\vartheta) \xrightarrow{\mathbb{P}_{\theta}} \arg\max_{\vartheta \in \Theta} \mathbb{E}_{\theta} \left[\psi(\vartheta, X) \right] = \theta.$$

EMV, asymptotique des Z- et M- estimateurs

Exemple estimateur de translation

Convergence des *M*-estimateurs

Proposition

Si le M-estimateur $\widehat{\theta}_n$ associé à la fonction de contraste est bien défini et si

- $\sup_{\vartheta \in \Theta} |M_n(\vartheta) M(\vartheta, \theta)| \stackrel{\mathbb{P}_{\theta}}{\longrightarrow} 0$,
- $\forall \varepsilon > 0$, $\sup_{|\vartheta \theta| \ge \varepsilon} M(\vartheta, \theta) < M(\theta, \theta)$ (condition de maximum)

alors

$$\widehat{\frac{\theta}{\theta_n}} \stackrel{\mathbb{P}_{\theta}}{\longrightarrow} \theta$$
.

■ La condition 1 (convergence uniforme) peut être délicate à montrer...

Loi limite des Z-estimateurs

- <u>Situation</u>: on observe $X_1, ..., X_n$ i.i.d. de loi dans la famille $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$, $\Theta \subset \mathbb{R}$.
- $\blacksquare \ \widehat{\theta}_n : Z\text{-estimateur}$ associé à $\phi : \Theta \times \mathbb{R} \to \mathbb{R}$ vérifie

$$\sum_{i=1}^n \phi(\widehat{\theta}_n, X_i) = 0$$

- Si $\widehat{\theta}_n$ est un M-estimateur associé à la fonction de contraste ψ régulière, alors c'est un Z-estimateur associé à la fonction $\phi(\vartheta,x)=\partial_\theta\psi(\vartheta,x)$.
- On suppose $\widehat{\theta}_n$ convergent. Que dire de sa loi limite?

Loi limite des Z-estimateurs : principe

■ Loi des grands nombres

$$Z_n(\vartheta) = \frac{1}{n} \sum_{i=1}^n \phi(\vartheta, X_i) \xrightarrow{\mathbb{P}_{\theta}} Z(\vartheta, \theta) = \mathbb{E}_{\theta} \left[\phi(\vartheta, X) \right]$$

■ Principe. Développement de Taylor autour de θ :

$$0 = Z_n(\widehat{\theta}_n) = Z_n(\theta) + (\widehat{\theta}_n - \theta)Z'_n(\theta) + \frac{1}{2}(\widehat{\theta}_n - \theta)^2 Z''(\widetilde{\theta}_n).$$

On néglige le reste :

$$\sqrt{n}(\widehat{\theta}_n - \theta) \approx \frac{-\sqrt{n}Z_n(\theta)}{Z_n'(\theta)}$$

Loi limite des Z-estimateurs : principe

■ Convergence du numérateur

$$\sqrt{n}Z_n(\theta) = \frac{1}{\sqrt{n}}\sum_{i=1}^n \phi(\theta, X_i) \stackrel{d}{\longrightarrow} \mathcal{N}(0, \mathbb{E}_{\theta}\left[\phi(\theta, X)^2\right])$$

$$\mathsf{si} \,\, \mathbb{E}_{\theta} \left[\phi(\theta, X) \right] = 0 \,\, \mathsf{et} \,\, \mathbb{E}_{\theta} \left[\phi(\theta, X)^2 \right] < +\infty.$$

■ Convergence du dénominateur

$$Z'_n(\theta) = \frac{1}{n} \sum_{i=1}^n \partial_{\theta} \phi(\theta, X_i) \xrightarrow{\mathbb{P}_{\theta}} \mathbb{E}_{\theta} \left[\partial_{\theta} \phi(\theta, X) \right]$$

$$\neq$$
 0 (à supposer).

lacktriangledown + hypothèses techniques pour contrôler le reste (besoin de la convergence de $\widehat{\theta}_n$).

Loi limite des Z-estimateurs

Proposition (Convergence des *Z*-estimateurs)

■ Soit Θ un ouvert de \mathbb{R} . Pour tout $\theta \in \Theta$, $\widehat{\theta}_n \stackrel{\mathbb{P}_\theta}{\to} \theta$, $\mathbb{E}_\theta \left[\phi(\theta, X)^2 \right] < +\infty$ et

$$\mathbb{E}_{\theta}\left[\phi(\theta,X)\right]=0,\;\mathbb{E}_{\theta}\left[\partial_{\theta}\phi(\theta,X)\right]\neq0.$$

■ (Contrôle reste) pour tout $\theta \in \Theta$, pour tout θ dans un voisinage de θ ,

$$|\partial_{\theta}^2 \phi(\vartheta, x)| \le g(x), \ \mathbb{E}_{\theta} [g(X)] < +\infty.$$

Alors

$$\sqrt{n}(\widehat{\underline{\theta}_n} - \theta) \stackrel{d}{\longrightarrow} \mathcal{N}\Big(0, \frac{\mathbb{E}_{\theta}[\phi(\theta, X)^2]}{\big(\mathbb{E}_{\theta}[\partial_{\theta}\phi(\theta, X)]\big)^2}\Big).$$



EMV, asymptotique des Z- et M- estimateurs
Approche asymptotique

Approche asymptotique

■ Hypothèse simplificatrice : $\theta \in \Theta \subset \mathbb{R}$. On se restreint aux estimateurs asymptotiquement normaux c'est-à-dire vérifiant

$$\sqrt{n}(\widehat{\theta}_n - \theta) \stackrel{d}{\longrightarrow} \mathcal{N}(0, \nu(\theta))$$

cf. théorèmes limites obtenus pour les Z-, M-estimateurs.

■ Si $\widehat{\theta}_{n,1}$ et $\widehat{\theta}_{n,2}$ as. normaux de variance asymptotique $v_1(\theta) \leq v_2(\theta)$, alors la précision de $\widehat{\theta}_{n,1}$ est asymptotiquement meilleure que celle de $\widehat{\theta}_{n,2}$ au point θ :

$$\widehat{\theta}_{n,1} = \theta + \sqrt{\frac{v_1(\theta)}{n}} \xi^{(n)}$$

$$\widehat{\theta}_{n,2} = \theta + \sqrt{\frac{v_2(\theta)}{n}} \zeta^{(n)}$$

où
$$\xi^{(n)}$$
 et $\zeta^{(n)} \xrightarrow{d} \mathcal{N}(0,1)$.

Approche asymptotique

Comparaison d'estimateurs : cas asymptotique

■ Si $v_1(\theta) < v_2(\theta)$, et si $\theta \rightsquigarrow v_i(\theta)$ est continue, on pose

$$C_{n,\alpha}(\widehat{\theta}_{n,i}) = \left[\widehat{\theta}_{n,i} \pm \sqrt{\frac{v_i(\widehat{\theta}_{n,i})}{n}} \Phi^{-1}(1 - \alpha/2)\right], \quad i = 1, 2$$

où $\alpha \in (0,1)$ et $\Phi(\cdot)$ est la fonction de répartition de la loi normale standard.

■ $C_{n,\alpha}(\widehat{\theta}_{n,i})$, i=1,2 sont deux intervalles de confiance asymptotiquement de niveau $1-\alpha$ et on a

$$\frac{|\mathcal{C}_{n,\alpha}(\widehat{\theta}_{n,1})|}{|\mathcal{C}_{n,\alpha}(\widehat{\theta}_{n,2})|} \xrightarrow{\mathbb{P}_{\theta}^{n}} \sqrt{\frac{v_{1}(\theta)}{v_{2}(\theta)}} < 1.$$

■ La notion de longueur minimale possible d'un intervalle de confiance est en général difficile à manipuler.

Approche asymptotique

Conclusion provisoire

- Il est difficile en général de comparer des estimateurs.
- Cadre asymptotique + normalité asymptotique \rightarrow comparaison de la variance asymptotique $\theta \rightsquigarrow \nu(\theta)$.
- Sous des hypothèses de régularité du modèle $\{\mathbb{P}^n_{\theta}, \theta \in \Theta\}$

alors

- Il existe une variance asymptotique $v^*(\theta)$ minimale parmi les variances de la classe des M-estimateurs as. normaux.
- Cette fonction est associée à une quantité d'information intrinsèque au modèle.
- La variance asymptotique de l'EMV est $v^*(\theta)$.
- Ceci règle partiellement le problème de l'optimalité.

Régularité d'un modèle statistique et information

■ Cadre simplificateur : modèle de densité

$$X_1,\ldots,X_n$$
 i.i.d. de loi \mathbb{P}_{θ}

dans la famille $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$ avec $\Theta \subset \mathbb{R}$ pour simplifier.

Notation :

$$f(\theta, x) = \frac{d \mathbb{P}_{\theta}}{d \mu}(x), \ \ x \in \mathbb{R}, \theta \in \Theta.$$

■ Hypothèse : la quantité

$$\left|\mathbb{I}(heta) = \mathbb{E}_{ heta}\left[\left(\partial_{ heta} \log f(heta, X)
ight)^2
ight]
ight|$$

est bien définie.

Construction de l'information de Fisher

Information de Fisher

Definition

- $\mathbb{I}(\theta) = \mathbb{E}_{\theta} \left[\left(\partial_{\theta} \log f(\theta, X) \right)^2 \right]$ s'appelle l'information de Fisher de la famille $\{ \mathbb{P}_{\theta}, \theta \in \Theta \}$ au point θ . Elle ne dépend pas de la mesure dominante μ .
- Le cadre d'intérêt est celui où

$$0 < \mathbb{I}(\theta) < +\infty$$
.

■ $\mathbb{I}(\theta)$ quantifie « l'information » qu'apporte chaque observation X_i sur le paramètre θ .

Remarque : on a $\mathbb{P}_{\theta}\left[f(\theta,X)>0\right]=1$, donc la quantité $\log f(\theta,X)$ est bien définie.

- Modèles réguliers et information de Fisher
 - Construction de l'information de Fisher

Information dans quel sens? Origine de la notion

- Supposons l'EMV $\widehat{\theta}_n^{mv}$ bien défini et convergent.
- Supposons l'application $(\theta, x) \rightsquigarrow f(\theta, x)$ possédant toutes les propriétés de régularité et d'intégrabilité voulues.
- Alors

$$\boxed{\sqrt{n}\big(\,\widehat{\boldsymbol{\theta}}_{\mathsf{n}}^{\,\mathsf{mv}}\,\!-\!\!\boldsymbol{\theta}\big) \stackrel{d}{\longrightarrow} \mathcal{N}\Big(\boldsymbol{0},\frac{1}{\underline{\mathbb{I}(\boldsymbol{\theta})}}\Big)}$$

en loi sous \mathbb{P}_{θ} , où encore

$$\widehat{ heta}_{\mathsf{n}}^{\;\mathsf{mv}} \overset{d}{pprox} heta + rac{1}{\sqrt{n\mathbb{I}(heta)}} \, \mathcal{N}(0,1)$$

en loi sous \mathbb{P}_{θ} .

Construction de l'information + jeu d'hypothèses attenant

- Heuristique : on établira un jeu d'hypothèses justifiant a posteriori le raisonnement.
- Etape 1 : l'EMV $\widehat{\theta}_{n}^{\text{mv}}$ converge :

$$\widehat{\theta}_{\mathsf{n}}^{\;\mathsf{mv}} \stackrel{\mathbb{P}_{\theta}}{\longrightarrow} \theta$$

via le théorème de convergence des M-estimateurs.

■ Etape 2 : l'EMV $\widehat{\theta}_{\mathbf{n}}^{\,\mathrm{mv}}$ est un \mathbf{Z} -estimateur :

$$0 = \partial_{\vartheta} \left(\sum_{i=1}^{n} \log f(\vartheta, X_{i}) \right)_{\vartheta = \widehat{\theta}_{\mathbf{n}}^{mv}}.$$

Construction de $\mathbb{I}(\theta)$ cont.

Etape 3 : développement asymptotique autour de θ :

$$0 \approx \sum_{i=1}^{n} \partial_{\theta} \log f(\theta, X_{i}) + (\widehat{\theta}_{n}^{mv} - \theta) \sum_{i=1}^{n} \partial_{\theta}^{2} \log f(\theta, X_{i}),$$

soit

$$\widehat{\theta}_{\mathsf{n}}^{\mathsf{mv}} - \theta \approx -\frac{\sum_{i=1}^{n} \partial_{\theta} \log f(\theta, X_{i})}{\sum_{i=1}^{n} \partial_{\theta}^{2} \log f(\theta, X_{i})}$$

■ Etape 4 : le numérateur. Normalisation et convergence de $\frac{1}{\sum_{i=1}^{n} \partial_{\theta} \log f(\theta, X_i)}$?

MAP 433: Introduction aux méthodes statistiques. Cours 4

- Modèles réguliers et information de Fisher

Construction de l'information de Fisher

Numérateur

Lemme

On a

$$\mathbb{E}_{\theta}\left[\frac{\partial_{\theta}\log f(\theta,X)}{\partial\theta}\right]=0.$$

Démonstration.

$$\mathbb{E}_{\theta} \left[\frac{\partial_{\theta} \log f(\theta, X)}{\partial_{\theta} \log f(\theta, x)} \right] = \int_{\mathbb{R}} \partial_{\theta} \log f(\theta, x) f(\theta, x) \mu(dx)$$

$$= \int_{\mathbb{R}} \frac{\partial_{\theta} f(\theta, x)}{f(\theta, x)} f(\theta, x) \mu(dx)$$

$$= \int_{\mathbb{R}} \partial_{\theta} f(\theta, x) \mu(dx)$$

- Modèles réguliers et information de Fisher

Construction de l'information de Fisher

Dénominateur

De même $\int_{\mathbb{R}} \partial_{\theta}^2 f(\theta, x) \mu(dx) = 0$. Conséquence :

$$\boxed{\mathbb{I}(\theta) = \mathbb{E}_{\theta} \left[\left(\partial_{\theta} \log f(\theta, X) \right)^{2} \right] = -\mathbb{E}_{\theta} \left[\partial_{\theta}^{2} \log f(\theta, X) \right]}$$

En effet

$$\begin{split} &\mathbb{E}_{\theta} \left[\partial_{\theta}^{2} \log f(\theta, X) \right] \\ &= \int_{\mathbb{R}} \frac{\partial_{\theta}^{2} f(\theta, x) f(\theta, x) - \left(\partial_{\theta} f(\theta, x) \right)^{2}}{f(\theta, x)^{2}} f(\theta, x) \mu(dx) \\ &= \int_{\mathbb{R}} \partial_{\theta}^{2} f(\theta, x) \mu(dx) - \int_{\mathbb{R}} \frac{\left(\partial_{\theta} f(\theta, x) \right)^{2}}{f(\theta, x)} \mu(dx) \\ &= 0 - \int_{\mathbb{R}} \left(\frac{\partial_{\theta} f(\theta, x)}{f(\theta, x)} \right)^{2} f(\theta, x) \mu(dx) = - \mathbb{E} \left[\left(\partial_{\theta} \log f(\theta, X) \right)^{2} \right]. \end{split}$$

- Modèles réguliers et information de Fisher
- Construction de l'information de Fisher

Conséquences

Les $\partial_{\theta} \log f(\theta, X_i)$ sont i.i.d. et $\mathbb{E}_{\theta} \left[\partial_{\theta} \log f(\theta, X) \right] = 0$. TCL :

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \partial_{\theta} \log f(\theta, X_{i}) \stackrel{d}{\longrightarrow} \mathcal{N}(0, \mathbb{E}_{\theta} \left[\left(\partial_{\theta} \log f(\theta, X) \right)^{2} \right])$$
$$= \mathcal{N}(0, \mathbb{I}(\theta)).$$

Les $\partial_{\theta}^2 \log f(\theta, X_i)$ sont i.i.d. LGN :

$$\frac{1}{n} \sum_{i=1}^{n} \frac{\partial_{\theta}^{2} \log f(\theta, X_{i})}{\sum_{i=1}^{n} \mathbb{E}_{\theta} \left[\frac{\partial_{\theta}^{2} \log f(\theta, X)}{\sum_{i=1}^{\text{conséquence}} -\mathbb{I}(\theta). \right]}$$

- Modèles réguliers et information de Fisher
 - Construction de l'information de Fisher

Conclusion

■ En combinant les deux estimations + lemme de Slutsky :

$$\begin{split} \sqrt{n}(\widehat{\theta}_{\mathsf{n}}^{\,\mathsf{mv}} - \theta) &\approx -\frac{\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \partial_{\theta} \log f(\theta, X_{i})}{\frac{1}{n} \sum_{i=1}^{n} \partial_{\theta}^{2} \log f(\theta, X_{i})} \\ &\xrightarrow{d} \frac{\mathcal{N}(0, \mathbb{I}(\theta))}{\mathbb{I}(\theta)} \\ &\stackrel{\mathsf{loi}}{=} \mathcal{N}\Big(0, \frac{1}{\mathbb{I}(\theta)}\Big). \end{split}$$

Le raisonnement est rigoureux dès lors que : (i) on a la convergence de θ̂_n^{mv}, (ii) on peut justifier le lemme et sa conséquence, (iii) I(θ) est bien définie et non dégénérée et (iv) on sait contrôler le terme de reste dans le développement asymptotique, partie la plus difficile.

└ Modèle régulier

Modèle régulier

Definition

La famille de densités $\{f(\theta,\cdot), \theta \in \Theta\}$, par rapport à la mesure dominante $\mu, \Theta \subset \mathbb{R}$, est régulière si

- Θ ouvert et $\{f(\theta,\cdot)>0\}=\{f(\theta',\cdot)>0\}$, $\forall \theta,\theta'\in\Theta$.
- μ -p.p. $\theta \leadsto f(\theta, \cdot)$, $\theta \leadsto \log f(\theta, \cdot)$ sont C^2 .
- $\forall \theta \in \Theta, \exists \mathcal{V}_{\theta} \subset \Theta \text{ t.q. pour } a \in \mathcal{V}_{\theta}$

$$|\partial_a^2 \log f(a,x)| + |\partial_a \log f(a,x)| + (\partial_a \log f(a,x))^2 \le g(x)$$

οù

$$\int_{\mathbb{R}} g(x) \sup_{a \in \mathcal{V}(\theta)} f(a, x) \mu(dx) < +\infty.$$

L'information de Fisher est non-dégénérée :

└ Modèle régulier

Résultat principal

Proposition

Si l'expérience engendrée par l'observation $X_1, \ldots, X_n \sim_{i.i.d.} \mathbb{P}_{\theta}$ est associée à une famille de probabilités $\{\mathbb{P}_{\theta}, \theta \in \Theta\}$ sur \mathbb{R} régulière au sens de la définition précédente, alors

$$\sqrt{n} \Big(\widehat{\theta}_{\mathsf{n}}^{\;\mathsf{mv}} - \theta \Big) \overset{d}{\longrightarrow} \mathcal{N} \Big(0, \frac{1}{\mathbb{I}(\theta)} \Big).$$

■ $Si \ \widehat{\theta}_n$ est un Z-estimateur régulier asymptotiquement normal de variance $v(\theta)$, alors

$$\forall heta \in \Theta, \ \ v(heta) \geq rac{1}{\mathbb{I}(heta)}.$$



└ Modèle régulier

Preuve de la proposition

- Le premier point consiste à rendre rigoureux le raisonnement précédent. Point délicat : le contrôle du terme de reste.
- Optimalité de la variance de l'EMV parmi celle des Z-estimateurs : on a vu que si $\widehat{\theta}_n$ est un Z-estimateur régulier associé à la fonction ϕ , alors, sa variance asymptotique $v(\theta) = v_{\phi}(\theta)$ vaut

$$onumber v_\phi(heta) = rac{\mathbb{E}_ heta\left[\phi(heta,X)^2
ight]}{\left(\mathbb{E}_ heta\left[\partial_ heta\phi(heta,X)
ight]
ight)^2}.$$

A montrer : pour toute fonction ϕ :

$$\boxed{\frac{\mathbb{E}_{\theta}\left[\phi(\theta,X)^{2}\right]}{\left(\mathbb{E}_{\theta}\left[\partial_{\theta}\phi(\theta,X)\right]\right)^{2}} \geq \frac{1}{\mathbb{I}(\theta)}}.$$

└─ Modèle régulier

Preuve de l'inégalité

Par construction

$$\partial_a \mathbb{E}_{\theta} \left[\phi(a, X) \right]_{\big| a = \theta} = 0.$$

• (avec $\dot{\phi}(\theta, x) = \partial_{\theta}\phi(\theta, x)$)

$$\begin{split} 0 &= \int_{\mathbb{R}} \left[\dot{\phi}(\theta, x) f(\theta, x) + \phi(\theta, x) \partial_{\theta} f(\theta, x) \right] \mu(dx) \\ &= \int_{\mathbb{R}} \left[\dot{\phi}(\theta, x) f(\theta, x) + \phi(\theta, x) \partial_{\theta} \log f(\theta, x) f(\theta, x) \right] \mu(dx). \end{split}$$

Conclusion

$$egin{aligned} \mathbb{E}_{ heta}\left[\dot{\phi}(heta, \mathsf{X})
ight] = -\mathbb{E}_{ heta}\left[\phi(heta, \mathsf{X})\partial_{ heta}\log f(heta, \mathsf{X})
ight] \end{aligned}$$

- Modèles réguliers et information de Fisher

└─ Modèle régulier

Preuve de l'inégalité (fin)

On a

$$\mathbb{E}_{\theta} \left[\dot{\phi}(\theta, X) \right] = - \mathbb{E}_{\theta} \left[\phi(\theta, X) \partial_{\theta} \log f(\theta, X) \right]$$

Cauchy-Schwarz :

$$\left(\mathbb{E}_{\theta}\left[\dot{\phi}(\theta, X)\right]\right)^{2} \leq \mathbb{E}_{\theta}\left[\phi(\theta, X)^{2}\right] \mathbb{E}_{\theta}\left[\left(\partial_{\theta} \log f(\theta, X)\right)^{2}\right],$$

c'est-à-dire

$$v_{\phi}(heta)^{-1} = rac{ig(\mathbb{E}_{ heta} \left[\dot{\phi}(heta, X)
ight] ig)^2}{\mathbb{E}_{ heta} \left[\phi(heta, X)^2
ight]} \leq \mathbb{I}(heta).$$

Cadre général et interprétation géométrique

Information de Fisher dans un modèle général

Definition

■ Situation : suite d'expériences statistiques

$$\mathcal{E}^{n} = \left(\mathfrak{Z}^{n}, \mathcal{Z}^{n}, \{\mathbb{P}_{\theta}^{n}, \theta \in \Theta\}\right)$$

dominées par μ_n , associées à l'observation $Z^{(n)}$,

$$f_n(\theta,z) = \frac{d \mathbb{P}^n_{\theta}}{du^n}(z), \ z \in \mathfrak{Z}^n, \theta \in \Theta \subset \mathbb{R}.$$

Information de Fisher (si elle existe) de l'expérience au point θ :

$$\mathbb{I}(\theta \mid \mathcal{E}_n) = \mathbb{E}_{\theta}^n \left[\left(\partial_{\theta} \log f_n(\theta, Z^{(n)}) \right)^2 \right]$$

Le cas multidimensionnel

- Même contexte que précédemment, avec $\Theta \subset \mathbb{R}^d$, et $d \geq 1$.
- Matrice d'information de Fisher

$$\mathbb{I}(\theta) = \mathbb{E}_{\theta} \left[\nabla_{\theta} \log f(\theta, Z^n) \nabla_{\theta} \log f(\theta, Z^n)^T \right]$$

matrice symétrique positive.

■ Si $\mathbb{I}(\theta)$ définie et si \mathcal{E}^n modèle de densité, en généralisant à la dimension d les conditions de régularité, on a

$$\sqrt{n}\big(\,\widehat{\boldsymbol{\theta}}_n^{\,\text{mv}}\, - \!\boldsymbol{\theta}\big) \stackrel{d}{\longrightarrow} \mathcal{N}\Big(\boldsymbol{0}, \underline{\mathbb{I}(\boldsymbol{\theta})}^{-1}\Big).$$

Cadre général et interprétation géométrique

Interprétation géométrique

• On pose $\mathbb{D}(\vartheta,\theta) = \mathbb{E}_{\theta} \left[\log f(\vartheta,X) \right]$. On a vu (inégalité d'entropie) que

$$\mathbb{D}(\vartheta,\theta) = \int_{\mathbb{R}} \log f(\vartheta,x) f(\theta,x) \mu(dx)$$

$$\leq \int_{\mathbb{R}} \log f(\theta,x) f(\theta,x) \mu(dx) = \mathbb{D}(\theta,\theta).$$

On a

$$\boxed{\mathbb{I}(\theta) = \partial_{\vartheta}^{2} \mathbb{D}(\vartheta, \theta)_{\big|\vartheta = \theta}.}$$

- Si $\mathbb{I}(\theta)$ est « petite », le rayon de courbure de $\vartheta \leadsto \mathbb{D}(\vartheta, \theta)$ est grand dans un voisinage de θ : la stabilisation d'un maximum empirique (l'EMV) est plus difficile, rendant moins précis l'estimation.
- Si $\mathbb{I}(\theta)$ est « grande », le rayon de courbure est petit et le

- Modèles réguliers et information de Fisher
 - Cadre général et interprétation géométrique

Efficacité à un pas

- Dans un modèle régulier, le calcul numérique de l'EMV peut être difficile à réaliser.
- Si l'on dispose d'un estimateur $\widehat{\theta}_n$ asymptotiquement normal et si les évaluations

$$\ell'_n(\theta) = \frac{1}{n} \sum_{i=1}^n \partial_{\theta} \log f(\theta, X_i), \quad \ell''_n(\theta) = \frac{1}{n} \sum_{i=1}^n \partial_{\theta}^2 \log f(\theta, X_i)$$

sont faciles, alors on peut corriger $\widehat{\theta}_n$ de sorte d'avoir le même comportement asymptotique que l'EMV :

$$\widetilde{\theta}_n = \widehat{\theta}_n - \frac{\ell'_n(\widehat{\theta}_n)}{\ell''_n(\widehat{\theta}_n)}$$
 (algorithme de Newton)

satisfait

$$\left|\sqrt{n}(\widetilde{\theta}_n-\theta)\stackrel{d}{\longrightarrow}\mathcal{N}\left(0,\frac{1}{\mathbb{I}(\widehat{\theta})}\right)\right|$$