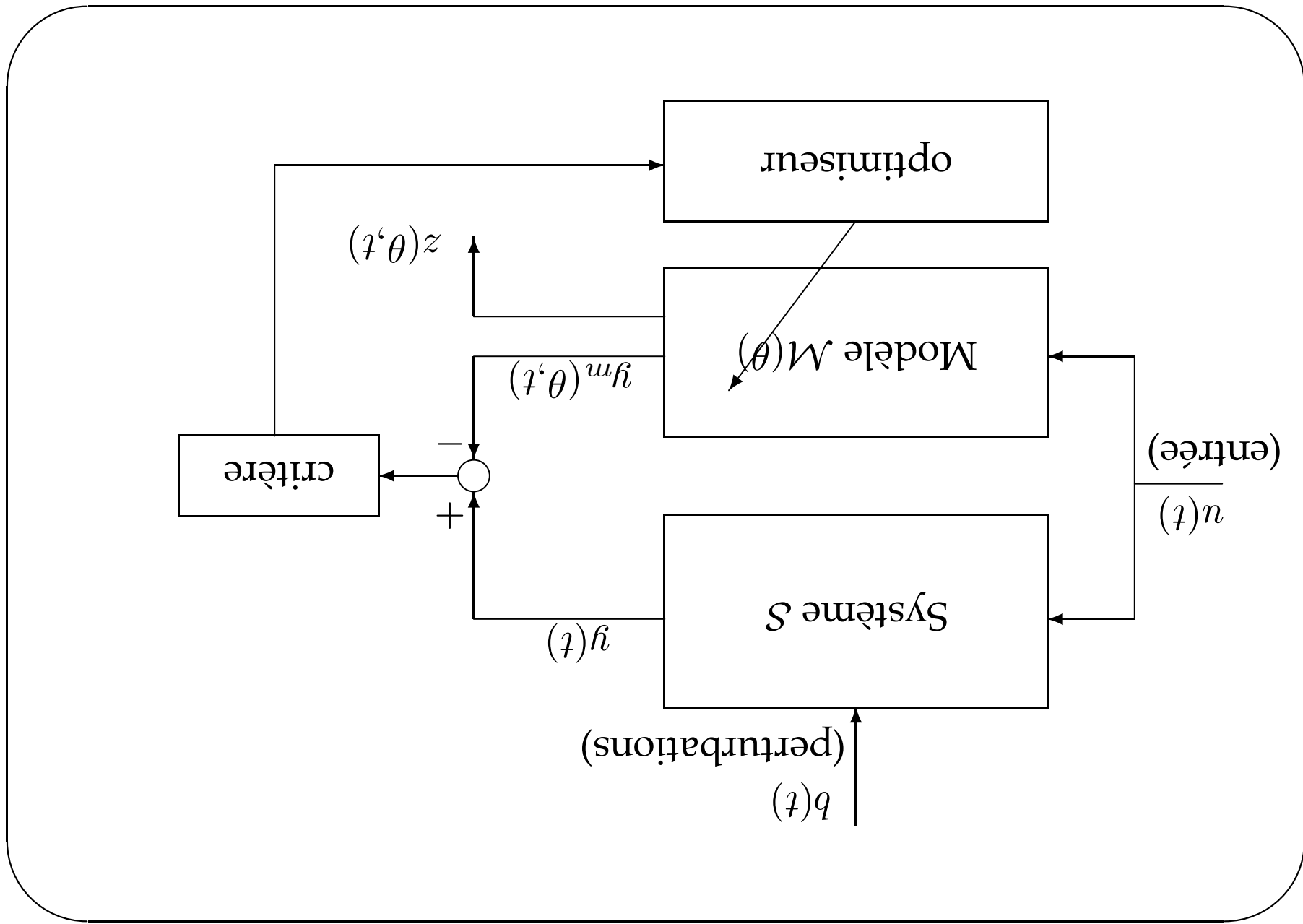


ESTIMATION DE PARAMÈTRES

Luc Pronzato, 2004

- 1 Moindres carrés (L_2)
- 2 Moindres valeurs absolues (L_1)
- 3 Minimum (L_∞)
- 4 Maximum de vraisemblance
- 5 Estimation bayésienne
- 6 Contraintes
- 7 Estimation robuste



0) Généralités

y vecteur de toutes les observations sur S

$y^m(\theta)$ vecteur des réponses de modèle associées pour des

paramètres θ

Choisir θ pour que $y^m(\theta)$ «ressemble» à y

→ minimiser $j(\theta)$, «distance» entre $y^m(\theta)$ et y → $\hat{\theta}(y)$

Rq1 : y aléatoire, donc $\hat{\theta}(y)$ aléatoire aussi (biais, variance, etc.)

Rq2 : Comment optimiser $j(\cdot)$? — voir plus loin

1) Moindres carrés

Distance quadratique (norme L_2)

$$j_{MC}(\theta) = \|\mathbf{y} - \mathbf{y}_m(\theta)\|_2^2 = \sum_{i=1}^N [y(i) - y_m(\theta, i)]^2$$

i est le numéro de la mesure, peut correspondre au temps t_i où on observe, avec $t_i = T$ si système échantillonné avec une période T .

Plus généralement, i correspond à des conditions expérimentales ξ_i

1.1) Moindres carrés pondérés :

$$j_{MCP}(\theta) = [\mathbf{y} - \mathbf{y}_m(\theta)]^T \mathbf{W} [\mathbf{y} - \mathbf{y}_m(\theta)] = \sum_{i=1}^N w_i [y(i) - y_m(\theta, i)]^2$$

avec $\mathbf{W} = \text{diag}(w_i, i = 1, \dots, N)$ matrice de pondération ($w_i \geq 0, \forall i$).
Choix des w_i ?

1.2) Répétition de mesures : n_i mesures $y_j(i)$, $j = 1, \dots, n_i$, pour les mêmes conditions expérimentales ξ_i (au même instant par exemple). Soit m le nb. de conditions expérimentales différentes $(\sum_{i=1}^m n_i = N)$. Posons

$$y_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_j(i)$$

pondérations $n_i w_i$

← moindres carrés pondérés pour les m observations y_i ,

$$= \sum_m^{i=1} n_i w_i [y_i - \hat{y}_i(\theta)]^2 + \text{terme ind. de } \theta$$

$$+ 2 \sum_m^{i=1} \sum_{n_i}^{j=1} w_i [y_j - \hat{y}_j(\theta)] \times [y_i - \hat{y}_i(\theta)]$$

$$= \sum_m^{i=1} \sum_{n_i}^{j=1} w_i [y_j - \hat{y}_j(\theta)]^2 + \sum_{n_i}^{j=1} \sum_m^{i=1} w_i [y_i - \hat{y}_i(\theta)]^2$$

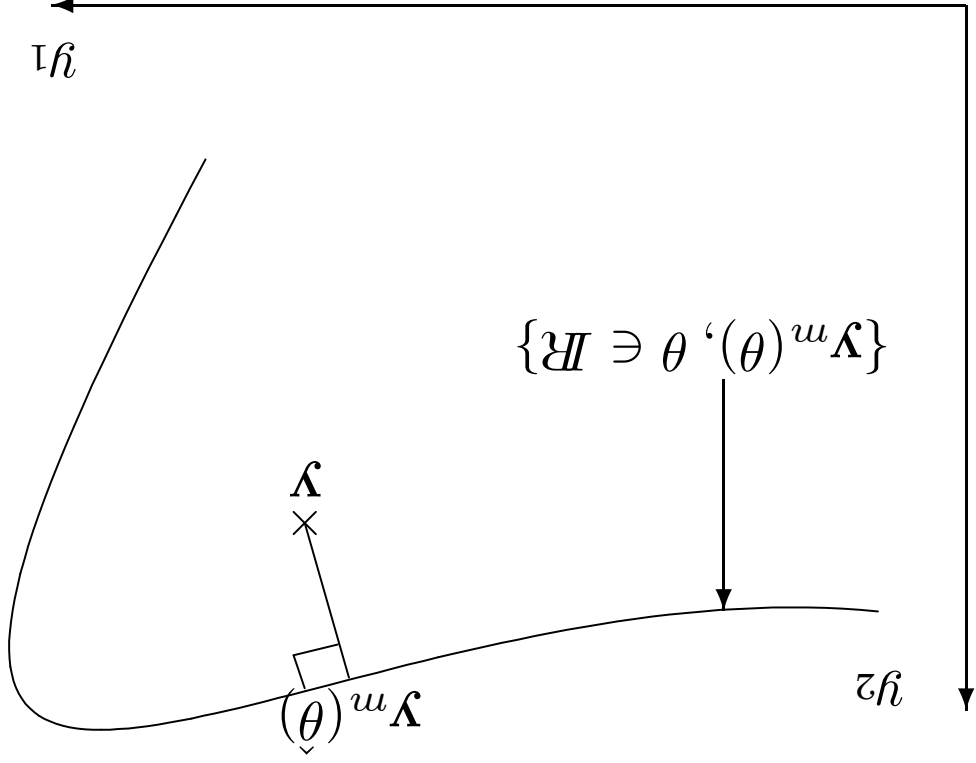
$$= \sum_m^{i=1} \sum_{n_i}^{j=1} w_i [y_j - \hat{y}_j(\theta) + y_i - \hat{y}_i(\theta)]^2$$

$$= \sum_m^{i=1} \sum_{n_i}^{j=1} w_i [y_j - \hat{y}_j(\theta) - (y_i - \hat{y}_i(\theta))]^2 = \hat{J}_{MCP}(\theta)$$

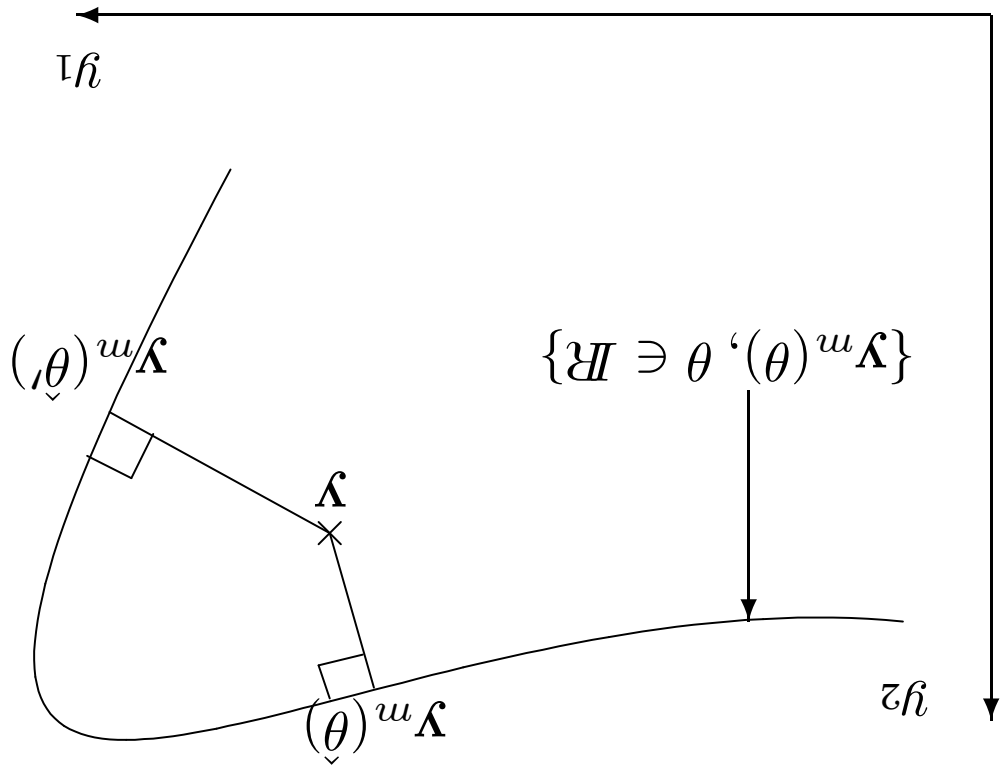
1.3) Minimas locaux :

$\hat{\theta}$ = estimateur des moindres carrés

$y_m(\hat{\theta}) = \text{projection } \perp \text{ de } y \text{ sur la surface des réponses de modèle}$
 $S_r = \{y_m(\theta), \theta \in \mathbb{R}^p\}$.



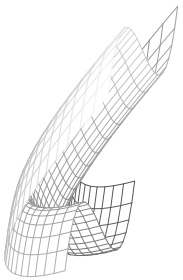
Il peut y avoir des minima locaux si S_r est courbe !



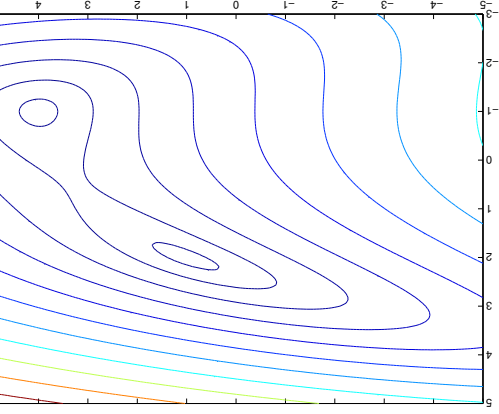
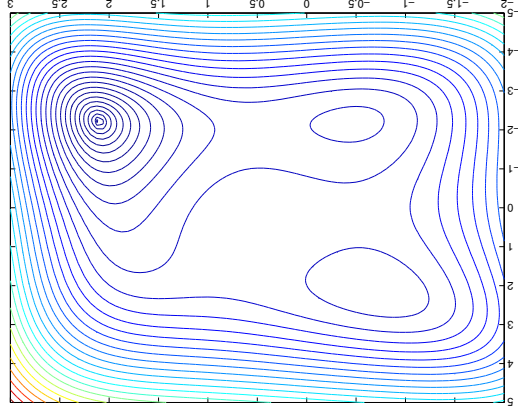
$$y^m(\theta, \mathbf{x}) = \theta_1 x_1 + \theta_2 x_2 + \theta_3^{\frac{1}{2}} (1 - x_1) + \theta_2^{\frac{1}{2}} (1 - x_2)$$

3 observations, en $\mathbf{x}^1 = (1, 0)^{\top}$,
 $\mathbf{x}^2 = (1, 1)^{\top}$ et $\mathbf{x}^3 = (0, 1)^{\top}$

$$\mathbf{y}^m(\theta), \theta_1 \in [-5, 5], \theta_2 \in [-2, 5]$$



isocritères dans le plan
 $(\theta_1, \theta_2), y = (5, -10, 8)^T$



répétition de mesures :

$$\mathbf{x}^1 = (1, 0)^T, \mathbf{x}^2 = \mathbf{x}^3 = (1, 1)^T$$

$$\mathbf{y} = (5, 2, 4)^T$$

Répétition de mesures \rightarrow pertes d'identifiabilité globale :

$$\hat{\theta}^1 = (\hat{\theta}_1, \hat{\theta}_2)^T \text{ et } \hat{\theta}^2 = (\hat{\theta}_1 + 2\hat{\theta}_2 - 1, 1 - \hat{\theta}_2)^T \text{ donnent le même } y^m(\theta)$$

← Supprimer la courbe de S_r

Plusieurs observations pour les mêmes conditions expérimentales :

$$m = p = \dim(\theta)$$

⇒ S_r est plane

⇒ plus de minima locaux !

← On résout en $\theta : y^m(\theta, i) = y(i), i = 1, \dots, p.$

Parfois plusieurs solutions (identifiabilité locale) :

👉 1) chercher toutes les solutions

👉 2) utiliser + de conditions expérimentales

👉 3) \exists peut être des minima locaux : initialiser la recherche aux solutions trouvées en 1.

2) Moindres valeurs absolues (norme L_1)

$$j_{MVA}(\theta) = \|\mathbf{y} - \mathbf{y}^m(\theta)\|_1 = \sum_{i=1}^n |y(i) - y^m(\theta, i)|$$

$$j_{MVAP}(\theta) = \|\mathbf{W}[\mathbf{y} - \mathbf{y}^m(\theta)]\|_1 = \sum_{i=1}^n w_i |y(i) - y^m(\theta, i)|$$

Pas différentiable partout !

Optimum pas forcément unique (même si le modèle est s.g.i.)

Exemple : $y^m(\theta, 1) = \theta_1 + \theta_2^2$, $y^m(\theta, 2) = y^m(\theta, 3) = \theta_1 + \theta_2$, $\mathbf{y} = (5, 2, 4)$.
 $j_{MVA}(\theta) = 2 \forall \theta$ tel que $\theta_1 + \theta_2^2 = 5$, $2 \leq \theta_1 + \theta_2 \leq 4$

Pas recommandé tel quel — voir + loin

Modèles LP : $y_m(\theta, i) = \mathbf{r}_i^\top \theta + c_i$

Le problème revient à minimiser $\sum_{i=1}^N \alpha_i$, sous les contraintes

$$y(i) - \mathbf{r}_i^\top \theta - c_i \leq \alpha_i \text{ et } y(i) - \mathbf{r}_i^\top \theta - c_i \geq -\alpha_i, i = 1, \dots, N$$

$N + p$ variables (α et θ), $2N$ contraintes.

Les contraintes et la fonction sont linéaires en les variables \rightarrow
programmation linéaire

Répétition de mesures : Pas toujours comme en 1.2, mais sous certaines conditions ($m = d$) \rightarrow moindres valeurs absolues pondérées, pour les observations

$$y_i = \text{med}\{y_j(i), j = 1, \dots, n_i\}$$

avec les pondérations $n_i w_i$.

3) Minimax (norme L^∞)

$$j_{MM}(\theta) = \|\mathbf{y} - \mathbf{y}_m(\theta)\|_\infty = \max_{i=1, \dots, N} |y(i) - y_m(\theta, i)|$$

$$j_{MMP}(\theta) = \|\mathbf{W}[\mathbf{y} - \mathbf{y}_m(\theta)]\|_\infty = \max_{i=1, \dots, N} w_i |y(i) - y_m(\theta, i)|$$

Modèles LP : $y_m(\theta, i) = \mathbf{r}_i^\top \theta + c_i$

Minimiser α , sous les contraintes

$$y(i) - \mathbf{r}_i^\top \theta - c_i \leq \alpha \text{ et } y(i) - \mathbf{r}_i^\top \theta - c_i \geq -\alpha, \quad i = 1, \dots, N$$

$1 + p$ variables (α et θ), $2N$ contraintes.

← programmation linéaire

Répétition de mesures : Presque comme en 1.2, minimax pondéré, pour les observations

$$y_i = \frac{\max_{j=1, \dots, n_i} \{y_j(i)\} + \min_{j=1, \dots, n_i} \{y_j(i)\}}{2}$$

avec les pondérations $n_i w_i$

(La raison :

$$\max_{j=1, \dots, n_i} |y_j(i) - y_m(\theta, i)| = |y(i) - y_m(\theta, i)|$$

$$+ \frac{\max_{j=1, \dots, n_i} \{y_j(i)\} - \min_{j=1, \dots, n_i} \{y_j(i)\}}{2}$$

4) Maximum de vraisemblance

4.1) Généralités

y variable aléatoire, générée par un modèle déterministe, de paramètres θ (inconnus), perturbé par $(\epsilon_k)_k$ (suite de v.a.i., éventuellement filtrée par un partie du système)

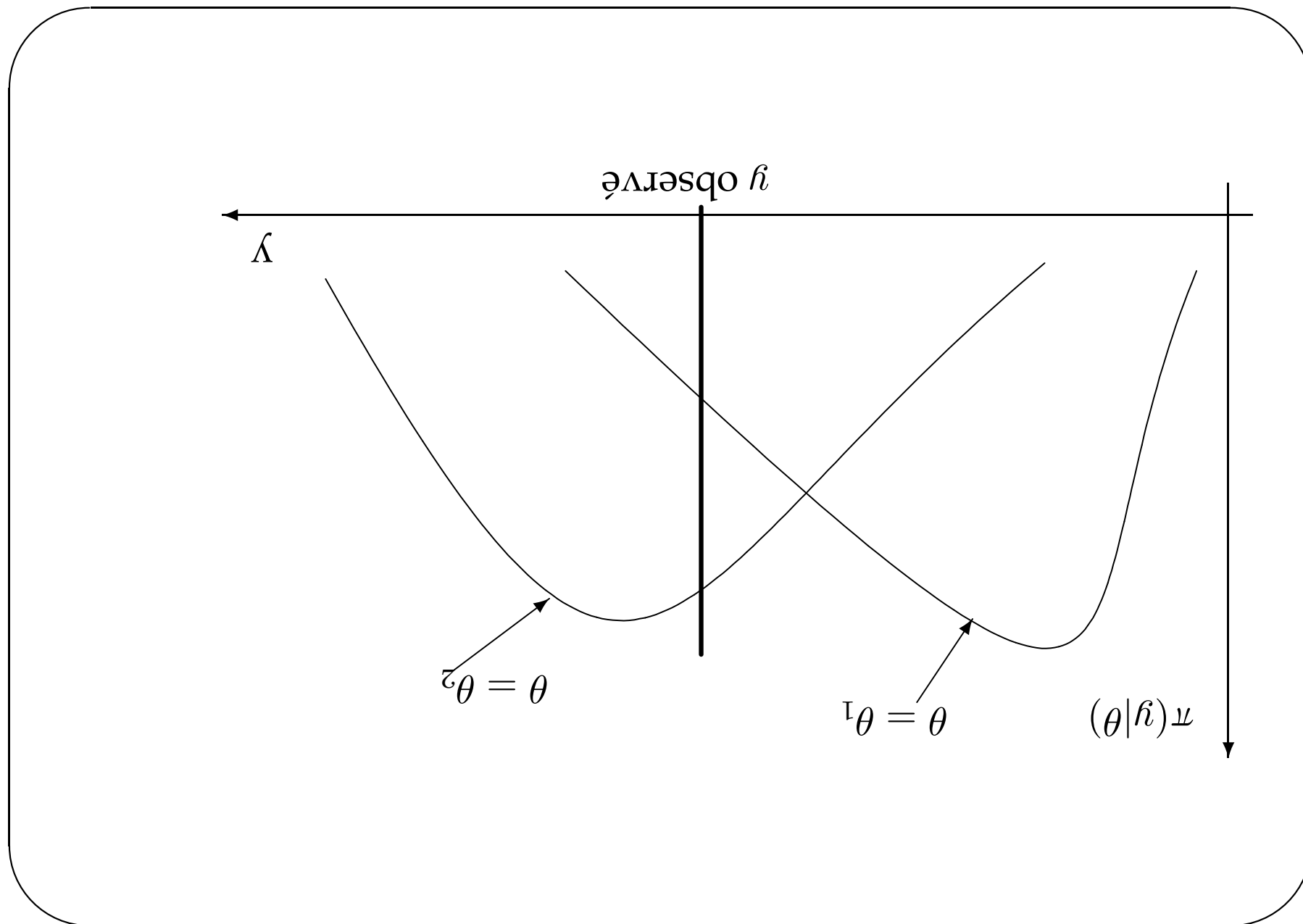
vraisemblance de $y : \pi(y|\theta)$

Ici, y est fixé (observations) : $\hat{\theta}_{MV}$ maximise $\pi(y|\theta)$

(en pratique on maximise la log-vraisemblance $\log \pi(y|\theta)$ —

pourquoi?)

Suivant la loi des perturbations \rightarrow divers critères d'estimation



4.2) Exemples (on fait toujours N observations, $i = 1, \dots, N$)

Ex1: $y(i) = \epsilon_i \sim \mathcal{N}(\mu, \sigma^2), \theta = (\mu, \sigma^2)$

$\hat{\theta}_{MV}$?

Biais $(E\{\hat{\theta}(y)|\theta\} - \theta)$?

Variance $(E\{\hat{\theta}(y) - E\{\hat{\theta}(y)|\theta\}\}[\hat{\theta}(y) - E\{\hat{\theta}(y)|\theta\}]^T | \theta)$?

Ex2: $y(i) = y_m(\theta, i) + \epsilon_i, (\epsilon_i)_i \text{ ind. } \epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$

a) σ_i^2 connus $\rightarrow \hat{\theta}_{MV}$?

b) $\forall i, \sigma_i^2 = \sigma^2$ (inconnu), $\hat{\theta}_{MV}$ et $\hat{\sigma}_{MV}^2$?

c) σ_i^2 inconnus, comment faire ?

Ex3: Sortie vectorielle

$y(i) = \mathbf{y}_m(\theta, i) + \epsilon_i, \mathbf{y}_i \in \mathbb{R}^n, (\epsilon_i)_i \text{ ind.}, \epsilon_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$

$\hat{\theta}_{MV}$ et $\hat{\Sigma}_{MV}$?

Ex4 : Bruit de Laplace, moyenne 0, variance s^2

$$y(i) = y_m(\theta, i) + \epsilon_i, (\epsilon_i)_i \text{ ind.},$$

$$\pi(\epsilon_i) = \frac{1}{\sqrt{2s}} \exp\left(-\frac{|\epsilon_i|}{s}\right)$$

a) σ_i connus $\rightarrow \hat{\theta}_{MV}$?

b) $\forall i, \sigma_i = \sigma$ (inconnu), $\hat{\theta}_{MV}$ et $\hat{\sigma}_{MV}$?

Ex5 : Bruit uniforme

$$y(i) = y_m(\theta, i) + \epsilon_i, (\epsilon_i)_i \text{ ind.},$$

$$\pi(\epsilon_i) = \begin{cases} \frac{1}{2A} & \text{si } |\epsilon_i| \leq A \\ 0 & \text{sinon} \end{cases}$$

a) A connu, $\hat{\theta}_{MV}$?

b) A inconnu, $\hat{\theta}_{MV}$, \hat{A}_{MV} ?

4.3) Complexité

Plusieurs structures de modèles M_i en compétition (par exemple, avec un nb. croissant de paramètres p_i)

→ estimer en même temps quelle est la bonne structure?
→ quelle est la valeur de ses paramètres?

Akaike Information Criterion : $j_{AIC}(\theta, i) = \frac{1}{N} [-\log \pi(\mathbf{y}|\theta) + p_i]$

(à minimiser)

Max de vraisemblance pour chaque structure (→ $\hat{\theta}_i$), choisir M_i avec $j_{AIC}(\hat{\theta}_i, i)$ minimum

Autres critères : toujours la même idée → pénaliser les structures complexes

4.4) Erreur de prédiction

Pour écrire $\pi(y|\theta)$, on utilise l'indépendance (\rightarrow produit de lois individuelles, $\log \pi(y|\theta) = \sum$)

\rightarrow toujours revenir à $(\epsilon_i)_i$

\rightarrow construire des variables $e(\theta, i)$ (**erreurs de prédiction**) qui

correspondent à ϵ_i sous l'hypothèse de données générées par $\mathcal{M}(\theta)$.

Ex : Box & Jenkins

$$y = F(\theta, q)u + G(\theta, q)\epsilon \rightarrow e(\theta, i) = G^{-1}(\theta, q)[y(i) - F(\theta, q)u(i)]$$

Rq 1 : Système LI \rightarrow on peut toujours propager l'erreur pour la rendre additive en sortie

Rq 2 : Th. de factorisation spectrale $\rightarrow \forall$ spectre de bruit sur y , on peut toujours le considérer comme un bruit blanc (suite de v.a.i.) filtré par $G(\theta, q)$, rationnel en q , stable et d'inverse stable.

4.5) Propriétés

Pourquoi le maximum de vraisemblance?

Hypothèses :

- ✓ données générées par $\mathcal{M}(\theta)$
- ✓ perturbations $(\epsilon_i)_i$ i.i.d. (quelque part)
- ✓ $\{y/\pi(y|\theta) > 0\}$ indépendant de θ
- ✓ $\frac{\partial^2 \log \pi(y|\theta)}{\partial \theta \partial \theta^\top}$ existe, continue en θ uniformément en y
- ✓ $\mathbb{E} \left\{ \left| \frac{\partial \theta_i}{\partial \log \pi(y|\theta)} \right| \mid \theta \right\} < \infty, \forall i = 1, \dots, p, \forall \theta$
- ✓ $\mathbb{E} \left\{ \left| \frac{\partial^2 \log \pi(y|\theta)}{\partial \theta_i \partial \theta_j} \right| \mid \theta \right\} < \infty, \forall (i, j) = 1, \dots, p, \forall \theta$

... alors

$$\sqrt{N}(\hat{\theta}_N - \theta) \underset{N \rightarrow \infty}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{M}^{-1}(\bar{\theta}, \bar{\xi}))$$

avec $\mathbf{M}(\bar{\theta}, \bar{\xi})$ la **matrice d'information de Fisher** (moyenne, par observation)

$$\mathbf{M}(\bar{\theta}, \bar{\xi}) = \mathbb{E} \left\{ \frac{1}{N} \frac{\partial \log \pi(\mathbf{y} | \bar{\theta})}{\partial \bar{\theta}} \frac{\partial \log \pi(\mathbf{y} | \bar{\theta})}{\partial \bar{\theta}^\top} \middle| \bar{\theta} \right\} = - \mathbb{E} \left\{ \frac{1}{N} \frac{\partial^2 \log \pi(\mathbf{y} | \bar{\theta})}{\partial \bar{\theta} \partial \bar{\theta}^\top} \middle| \bar{\theta} \right\}$$

(ξ caractérise les conditions expérimentales utilisées)

asymptotiquement sans biais

asymptotiquement efficace (retour à l'exemple 1 de 4.2)

Pour $y^{(i)} = y_m(\theta, i) + \epsilon_i$ avec $(\epsilon_i)_i$ i.i.d. (si ce n'est pas le cas, on se ramène à des erreurs de prédiction, et on considère $\pi[e(\theta, i)|\theta]$) :

$$\sum_{i=1}^N \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} = \frac{\partial \theta}{\partial \log \pi[\mathbf{y}|\theta]}$$

et dans $\mathbf{M}(\theta, \xi)$,

$$\mathbb{E} \left\{ \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \frac{\partial \theta}{\partial \log \pi[y^{(j)}|\theta]} \right\}_{\mathbb{E}} = \mathbb{E} \left\{ \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \right\}_{\mathbb{E}} \times \mathbb{E} \left\{ \frac{\partial \theta}{\partial \log \pi[y^{(j)}|\theta]} \right\}_{\mathbb{E}}$$

$$\int \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \pi[y^{(i)}|\theta] dy^{(i)} = \mathbb{E} \left\{ \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \right\}_{\mathbb{E}}$$

$$= \int \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \pi[y^{(i)}|\theta] dy^{(i)} = 0$$

$$\text{et } \mathbf{M}(\theta, \xi) = \mathbb{E} \left\{ \frac{1}{N} \sum_{i=1}^N \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \frac{\partial \theta}{\partial \log \pi[y^{(i)}|\theta]} \right\}_{\mathbb{E}}$$

Inégalité de Cramér-Rao :

Pour tout estimateur $\hat{\theta}$ non-biaisé ($E\{\hat{\theta}(y)|\theta\} = \theta$)

$$\text{Cov}(\hat{\theta}) \succeq \mathbf{M}^{-1}(\theta, \xi) / N = \left[E \left\{ \frac{\partial \log \pi(y|\theta)}{\partial \log \pi(y|\theta)} \frac{\partial \log \pi(y|\theta)}{\partial \theta} \middle| \theta \right\} \right]^{-1}$$

($A \succeq B \Leftrightarrow A - B$ définie non-négative)

Démonstration :

$$\nabla \text{ Calculer } \mathbf{A} = E\{\mathbf{v}(\theta, \mathbf{y}) \mathbf{v}^\top(\theta, \mathbf{y}) | \theta\} \text{ avec } \mathbf{v}(\theta, \mathbf{y}) = \begin{pmatrix} \hat{\theta}(\mathbf{y}) - \theta \\ \frac{\partial \theta}{\partial \log \pi(\mathbf{y}|\theta)} \end{pmatrix}$$

$$\leftarrow \mathbf{A} = \begin{pmatrix} \text{Cov}(\hat{\theta}) & \mathbf{I}_p \\ \mathbf{I}_p & \text{NM}(\theta, \xi) \end{pmatrix} \succeq \mathbf{O}$$

$$\nabla \text{ Calculer } \begin{pmatrix} \mathbf{I}_p & -\mathbf{M}^{-1}(\theta, \xi) / N \\ \mathbf{I}_p & -\mathbf{M}^{-1}(\theta, \xi) / N \end{pmatrix} \mathbf{A} \succeq \mathbf{O}$$

$$x p \frac{\partial}{\partial y} \int_q^v = \frac{\partial}{\partial y} \int_q^v x p f(x, y)$$

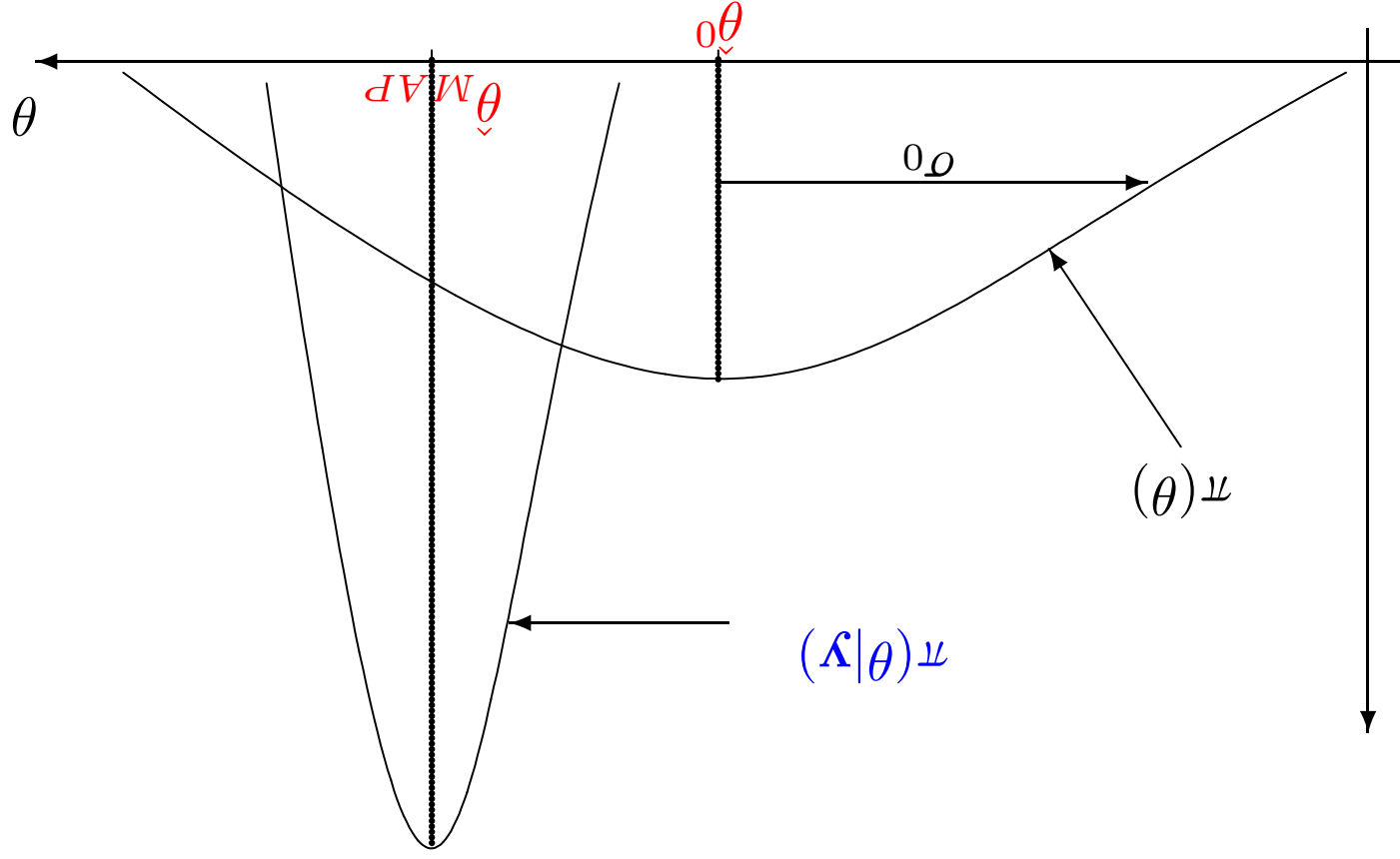
alors

$$\left. \begin{array}{l} f : [a, b] \times S \rightarrow \mathbb{R}, \quad -\infty < a < b \leq \infty, \quad S \subseteq \mathbb{R} \\ \exists \frac{\partial f(x, y)}{\partial y}, \quad \forall x \in [a, b], \quad \forall y \in S \\ \forall y \in S, \quad f(x, y) \text{ et } \frac{\partial f(x, y)}{\partial y} \text{ sont int\^egrables sur } [a, b] \\ \forall x \in [a, b], \quad \forall y \in S, \quad \left| \frac{\partial f(x, y)}{\partial y} \right| \leq g(x), \text{ avec } g(x) \text{ int\^egrable} \end{array} \right\} \text{S!}$$

On utilise la d\^erivation sous le signe $\int \dots$ rappel :

5) Estimation bayésienne

θ : v.a., densité connue $\pi(\theta)$, après observations $y \rightarrow \pi(\theta|y)$ (loi a posteriori)



Règle de Bayes :

$$\pi(\theta|y) = \frac{\pi(\theta, y)}{\pi(y)} = \frac{\pi(y|\theta)\pi(\theta)}{\pi(y)}$$

Estimateur du maximum *a posteriori* : $\hat{\theta}_{MAP}$ maximise $\pi(\theta|y)$
 ← maximise $\overbrace{\log \pi(y|\theta)}^{\text{log-vraisemblance}}$ + $\overbrace{\log \pi(\theta)}^{\text{a priori}}$

Mêmes propriétés asymptotiques que max. de vraisemblance

Si loi *a priori* normale $\mathcal{N}(\hat{\theta}_0, \Omega)$: $\hat{\theta}_{MAP}$ maximise

$$\log \pi(y|\theta) - \frac{1}{2}(\theta - \hat{\theta}_0)^\top \Omega^{-1}(\theta - \hat{\theta}_0)$$

Si en plus, modèle LP $y(\theta) = \mathbf{R}\theta + \epsilon$, avec $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$,
 $\hat{\theta}_{MAP} = (\mathbf{R}^\top \Sigma^{-1} \mathbf{R} + \Omega^{-1})^{-1} (\mathbf{R}^\top \Sigma^{-1} y + \Omega^{-1} \hat{\theta}_0)$

Coût $j(\hat{\theta}|\bar{\theta})$ d'estimer θ par $\hat{\theta}$ si vraie valeur $\bar{\theta}$ (commande optimale...)

Estimateur du risque minimum :

$$\begin{aligned}
 j_{RM}(\theta) &= \int j(\theta|\bar{\theta})\pi(\bar{\theta}|\mathbf{y})d\bar{\theta} \\
 &= \int \frac{1}{\pi(\mathbf{y})} j(\theta|\bar{\theta})\pi(\bar{\theta}|\mathbf{y})\pi(\bar{\theta})d\bar{\theta} \\
 &= \text{const.} \int \overbrace{j(\theta|\bar{\theta})}^{\text{coût perturbations}} \overbrace{\pi(\bar{\theta}|\mathbf{y})}^{\text{a priori}} d\bar{\theta}
 \end{aligned}$$

Si $j(\theta|\bar{\theta})$ quadratique en θ

$$j(\theta|\bar{\theta}) = (\theta - \bar{\theta})^T \mathbf{Q}(\theta - \bar{\theta}), \mathbf{Q} \text{ sym. } \succ \mathbf{0}$$

Stationnarité du critère $j_{RM}(\theta)$:

$$\frac{\partial j_{RM}(\theta)}{\partial \theta} \Big|_{\hat{\theta}_{RM}} = \mathbf{0} = 2\mathbf{Q} \int (\hat{\theta}_{RM} - \bar{\theta}) \pi(\bar{\theta}|\mathbf{y}) d\bar{\theta}$$

$\Rightarrow \hat{\theta}_{RM} = \int \bar{\theta} \pi(\bar{\theta}|\mathbf{y}) d\bar{\theta} = \mathbb{E}\{\theta|\mathbf{y}\}$ moyenne a posteriori de θ

... et si la loi a posteriori est symétrique, $\hat{\theta}_{RM} = \mathbb{E}\{\theta|\mathbf{y}\} = \hat{\theta}_{MAP}$

6) Contraintes

Égalités ($c_e(\theta) = 0$) ou inégalités ($c_i(\theta) \leq 0$)

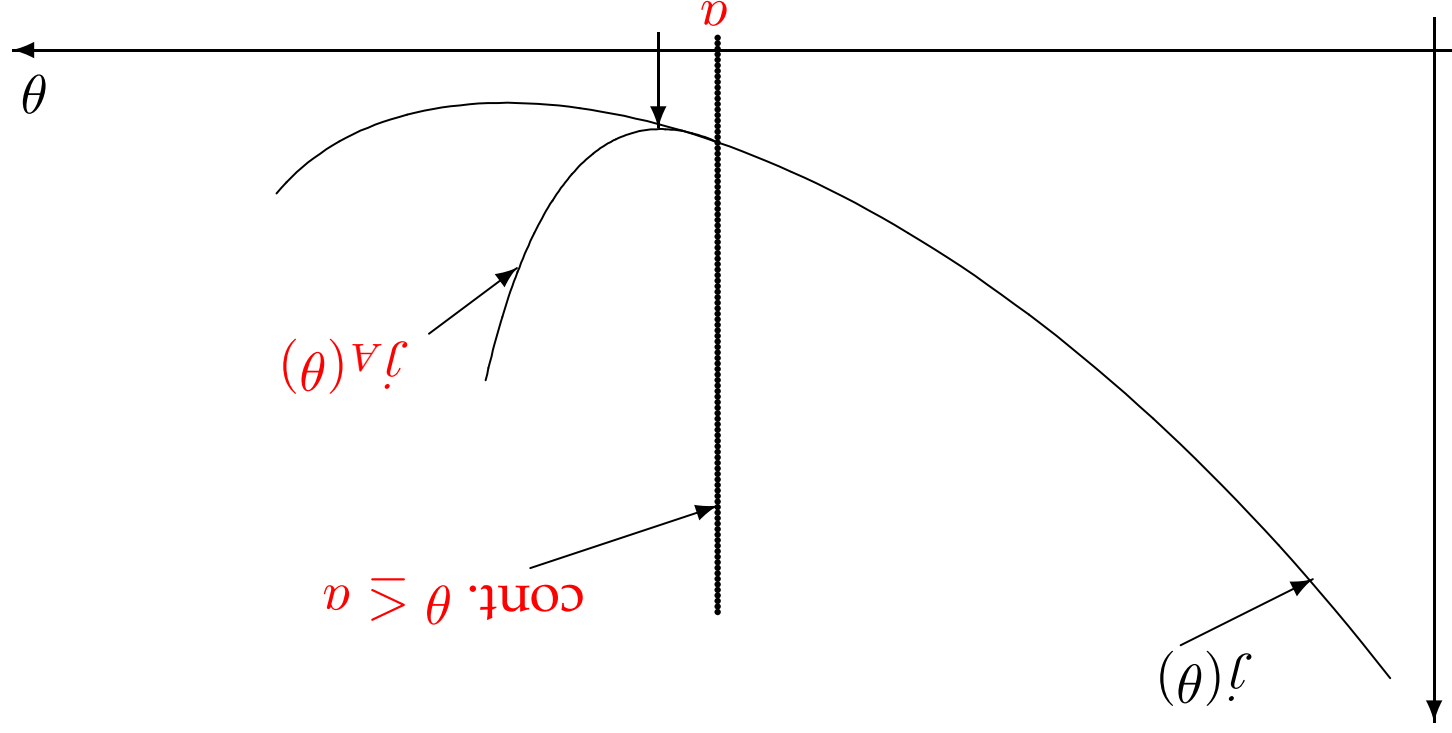
On les prend en compte soit par le critère, soit par l'algorithme d'optimisation

Par le critère : **pénalisation** (on suppose $j(\theta)$ à minimiser)

6.1) Inégalités

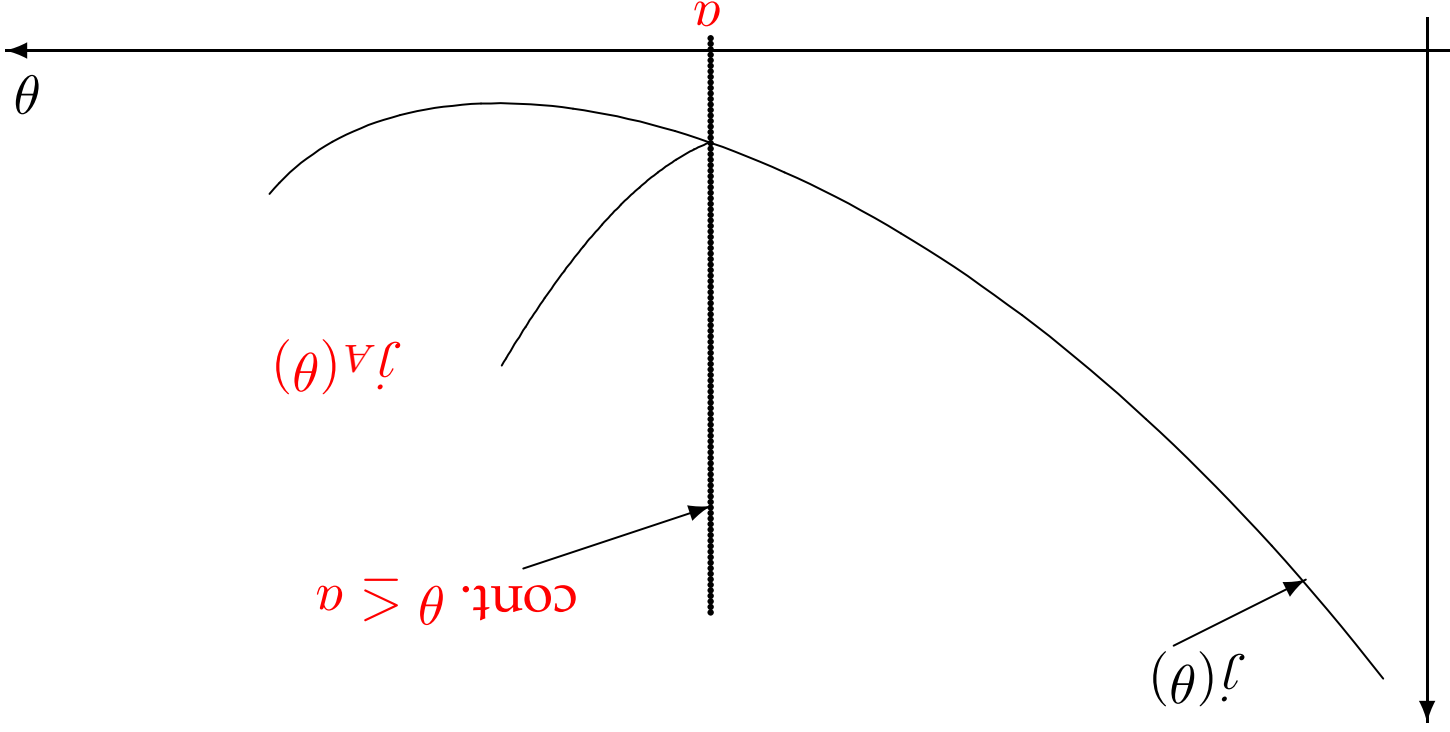
On définit $c_+^i(\theta) = \max\{c_i(\theta), 0\}$

Pénalisation extérieure :
 $j_A(\theta) = j(\theta) + (A/2)\|c_+^i(\theta)\|_2^2$: contrainte violée à l'optimum, VA



Problème : la tangente de la pénalisation $(A/2)\|c_+^i(\theta)\|_2^2$ en $\theta = a$ est nulle

← utiliser une pénalisation exacte : $j_A(\theta) = j(\theta) + A\|c_+^i(\theta)\|$



Pénalisation intérieure : reparamétriser

Ex : $\theta \leq a \mapsto \theta = a - \exp(\theta')$ ou $\theta = a - \theta'^2 \dots$

6.1) Égalités

Pénalisation : comme pour contraintes inégales

$$j_A(\theta) = j(\theta) + (A/2)\|c^e(\theta)\|_2^2$$

ou $j_A(\theta) = j(\theta) + A\|c^e(\theta)\|$ (pénalisation exacte)

$$\text{Lagangien : } \mathcal{L}(\theta, \mathbf{d}) = j(\theta) + \mathbf{d}^\top \mathbf{c}^e(\theta)$$

← minimum sur θ maximum sur \mathbf{d}

$$\text{Lagangien augmenté : } \mathcal{L}(\theta, \mathbf{d}) = j(\theta) + \mathbf{d}^\top \mathbf{c}^e(\theta) + (A/2)\|c^e(\theta)\|_2^2$$

Alterner

des minimisations sur $\theta : \hat{\theta}^k = \arg \min_{\theta} \mathcal{L}(\theta, \mathbf{d}^{k-1})$

des maximisations sur $\mathbf{d} : \mathbf{d}^k = \mathbf{d}^{k-1} + A c^e(\hat{\theta}^k)$ (gradient à pas fixe A)

7) Estimation robuste

7.1) Robustesse par rapport à la loi des erreurs : variance

Si $y(i) = y_m(\theta, i) + \epsilon_i$, avec $(\epsilon_i)_i$ i.i.d. $\pi(\epsilon)$, matrice d'inf. de Fisher :

$$\mathbf{M}(\theta, \xi) = \mathbf{E} \left\{ \frac{1}{N} \sum_{i=1}^N \frac{\partial \log \pi[y(i)|\theta]}{\partial \log \pi[y(i)|\theta]} \frac{\partial \theta}{\partial \theta^\top} \middle| \theta \right\} = \int \left[\frac{1}{N} \sum_{i=1}^N \frac{\partial \log \pi(e)}{\partial \log \pi(e)} \right]_2 \pi(e) de \left[\frac{\partial \theta}{\partial y_m(\theta, i)} \frac{\partial \theta}{\partial \theta^\top} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\partial \theta}{\partial y_m(\theta, i)} \frac{\partial \theta}{\partial \theta^\top} = \mathcal{I}(\pi)$$

avec $\mathcal{I}(\pi)$ l'information de Fisher de la loi $\pi : \mathcal{I}(\pi) = \int \left(\frac{d\pi(e)}{d\theta} \right)_2 \frac{\pi(e)}{1} de$

Propriétés asymptotiques de $\hat{\theta}_{MV}$: grande précision → «grande»
matrice $M(\theta, \xi)$

Approche minimum : MV pour la pire distribution dans une classe
donnée Π

$$\hat{\pi} = \arg \min_{\pi \in \Pi} \mathcal{I}(\pi)$$

Ex1 : $\Pi = \{\pi/\mathcal{I}(\pi) \text{ existe et } \pi(0) \geq 1/(s\sqrt{2}) > 0\}$

→ $\hat{\pi}$ = loi de Laplace, de moyenne 0 et variance s^2

Ex2 : $\Pi = \{\pi/\mathcal{I}(\pi) \text{ existe et } \text{var}(\pi) \leq \sigma^2\}$

→ $\hat{\pi}$ = loi normale, de moyenne 0 et variance σ^2

Observations

$$y(i) = \bar{\theta}_1 \exp(-\bar{\theta}_2 t_i) + \epsilon_i,$$

$$\bar{\theta} = (1, 2)^\top,$$

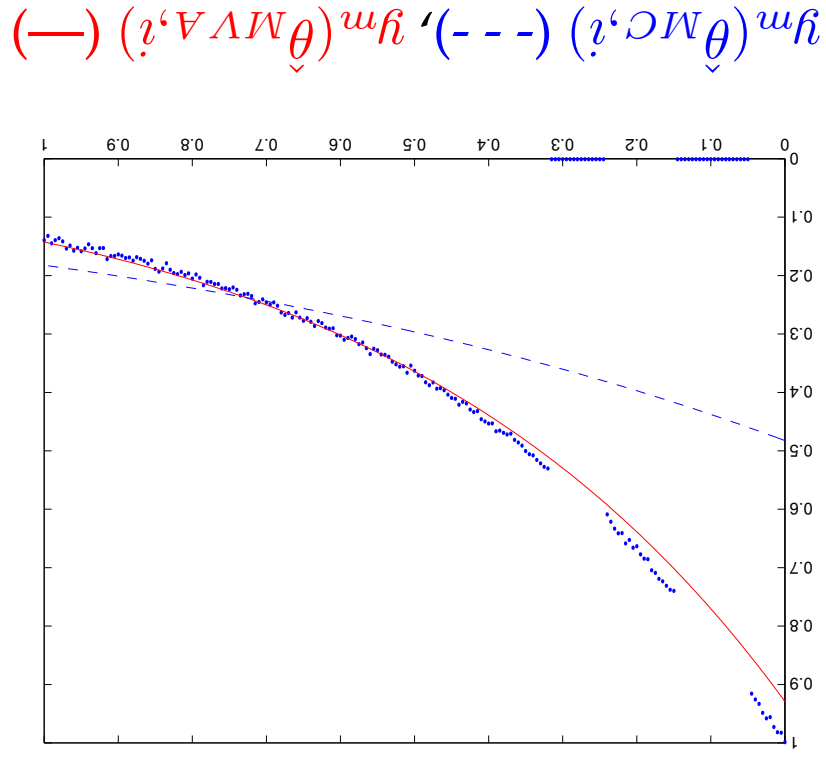
ϵ_i i.i.d. $\mathcal{N}(0, \sigma^2),$

$$\sigma = 5 \cdot 10^{-3}$$

t_i : 200 points dans $[0, 1]$

Données aberrantes:

$y(i) = 0$ de $i = 11$ à $i = 30$, puis de $i = 50$ à $i = 64$

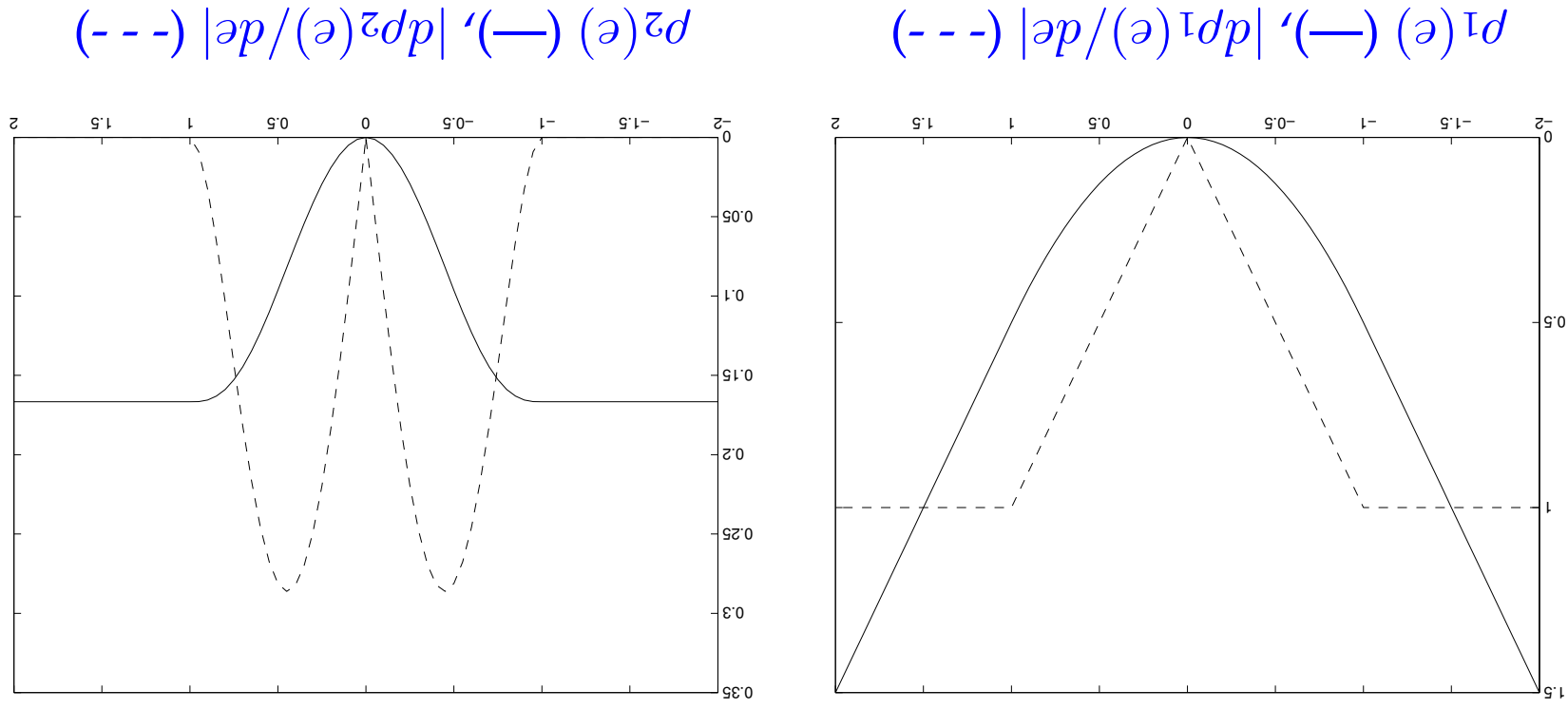


7.2) M-estimateurs

Moindres carrés, moindre valeurs absolues : cas particulier de

$$j(\theta) = \sum_{i=1}^N \rho[e(\theta, i)]$$

(avec $e(\theta, i) = y(i) - y_m(\theta, i)$ si erreur de sortie)



Estimateur non-redescendant : Huber

$$p_1(e) = \begin{cases} e^2/2 & \text{si } |e| \leq \delta \\ \delta|e| - \delta^2/2 & \text{sinon} \end{cases}$$

Estimateur redescendant : Tukey

$$p_2(e) = \begin{cases} \frac{1}{2} \left(e^2 - \frac{\delta}{e^4} + \frac{e}{3\delta^4} \right) & \text{si } |e| \leq \delta \\ \delta^2/6 & \text{sinon} \end{cases}$$

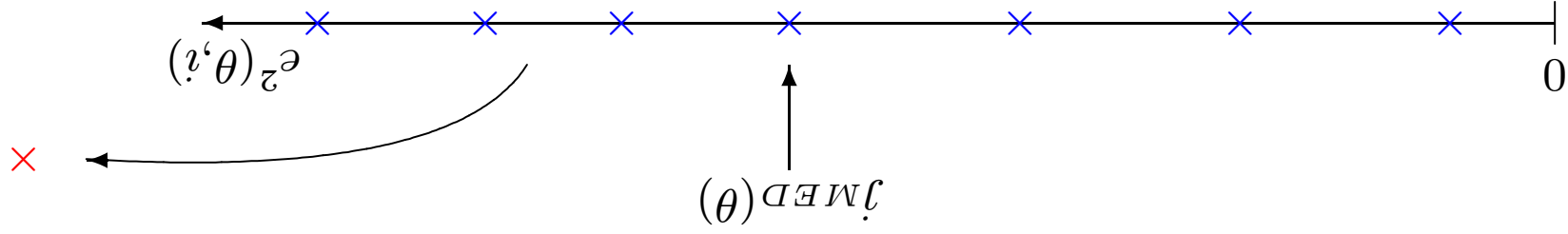
7.3) Point de rupture : données aberrantes → biais

y données correctes → $y'(\alpha)$ avec $\alpha\%$ de données aberrantes

$$\text{biais} = \mathbf{b}[y, \alpha, \hat{\theta}(\cdot)] = \max_{y'(\alpha)} \|\hat{\theta}(y) - \hat{\theta}[y'(\alpha)]\|$$

Mediane du carré des erreurs :

$$j_{MED}(\theta) = \text{med}\{e_2^c(\theta, i), i = 1, \dots, N\}$$



ou somme des h + petits carrés d'erreurs :

$$e_2^c(\theta, 1) \leq e_2^c(\theta, 2) \leq \dots \leq e_2^c(\theta, N), \text{ puis}$$

$$j_{Sh}(\theta) = \sum_{i=1}^h e_2^c(\theta, i)$$

$$\text{PR}[y, \hat{\theta}_{MED}(\cdot)] = [N/2 - p + 2]/N$$

$$h = [N/2] + 1 \rightarrow \text{PR}[y, \hat{\theta}_{Sh}(\cdot)] = ([N/2] - p + 2)/N$$

$$h = [N/2] + [p + 1]/2 \rightarrow \text{PR}[y, \hat{\theta}_{Sh}(\cdot)] = ([N - p]/2 + 1)/N = \text{Borne sup!}$$

$$\text{et } \lim_{N \rightarrow \infty} \text{PR}[y, \hat{\theta}_{MED}(\cdot)] = \lim_{N \rightarrow \infty} \text{PR}[y, \hat{\theta}_{Sh}(\cdot)] = 1/2$$

7.4) Application en traitement d'images

T_θ : Image $A \rightarrow$ Image $T_\theta(A)$, avec θ paramètres de translation, rotation, facteur d'échelle, changement de niveaux de gris...

2 images A et B , calibration automatique : transformer B pour qu'elle ressemble le + possible à A (imagerie médicale)

Considérer B comme $T_\theta(A)$ + bruit, et estimer θ qui maximise un critère de ressemblance entre A et B

Problème : \exists différences « significatives », non explicables par T_θ !
 \Leftrightarrow données aberrantes !

Transformer A et B en un signal à une dimension (par ex., balayage horizontal),

$$\hat{\theta} = \arg \max_{\theta} \text{Nb. de changements de signe dans } [T_\theta(A) - B]$$

Rapide à calculer, robuste vis-à-vis de données aberrantes

