

# MDI220, Statistique

## Cours 2: Estimation ponctuelle

Anne Sabourin

19 Septembre 2017

1. M- et Z- estimation : exemple et cadre général
2. Maximum de vraisemblance
3. Méthode des moindres carrés
4. Méthode des moments

# Cadre de l'estimation

- $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  un modèle statistique sur l'espace d'observations  $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ .
- $\mathbf{X}$  : les données,  $\mathbf{X} = (X_1, \dots, X_n) \stackrel{\text{i.i.d.}}{\sim} P_\theta$
- **But** : estimer  $g(\theta) \in \mathcal{A}$  une quantité d'intérêt.
- **Estimateur** (Rappel) :  
une fonction  $\mathcal{X}^n \rightarrow \mathcal{A}$ , *i.e.* une **statistique**.

1. M- et Z- estimation : exemple et cadre général

2. Maximum de vraisemblance

3. Méthode des moindres carrés

4. Méthode des moments

## Idée directrice

- M-estimateur : construit en **Minimisant** (par rapport à  $\theta$ ) une fonction qui dépend de  $\theta$  et de  $\mathbf{X}$ .

→  $\hat{\theta}$  est un arg max

- Z-estimateur : construit en annulant (*i.e.* en trouvant d'un **Z**éro) une fonction dépendant de  $\theta$  et de  $\mathbf{X}$ , en faisant varier  $\theta$ .

→  $\hat{\theta}$  est une racine.

## Exemple type de M-estimateur

- Modèle dominé (par la mesure de Lebesgue) sur  $\mathbb{R}$  :  $P_\theta$  a une densité  $p_\theta(x)$ .
- à  $x = (x_1, \dots, x_n)$  fixé,  $t \mapsto p_t(x)$  est la fonction de vraisemblance
- **But** : estimer  $g(\theta) = \theta$ . (donc  $\mathcal{A} = \Theta$ ).
- Supposons que  $\forall x \in \mathcal{X}^n, \exists ! \hat{\theta}(x)$  tel que

$$\forall t \in \Theta, p_t(x) \leq p_{\hat{\theta}(x)}(x).$$

- On pose  $M(x, t) = -\log p(x, t)$ ,  $t \in \Theta$

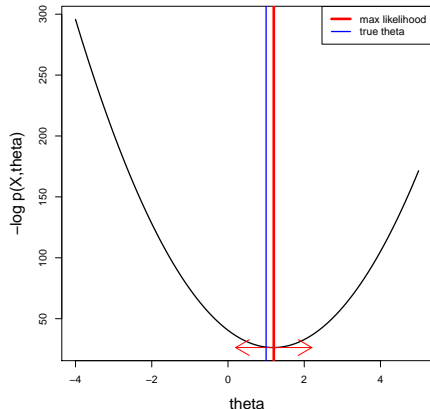
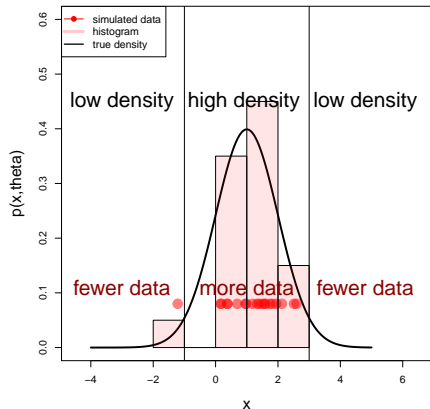
$$\hat{\theta}(x) = \operatorname{argmin}_{t \in \Theta} M(x, t),$$

$M$  est appelée 'fonction de contraste' ou 'contraste'.

- ex : modèle gaussien  $\mathcal{N}(\theta, \sigma^2)$ ,  $\sigma^2$  connu.  $-\log p_t(x) = ?$   $\hat{\theta}(x) = ?$

# Max de vraisemblance : justification heuristique.

Simulation :  $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta = 1, 1)$ ,  $1 \leq i \leq n = 20$ .



$X_i$  plus fréquents où  $p_\theta(x)$  grande  
 $p_\theta(x)dx \approx \mathbb{P}_\theta(X_i \in dx)$

$\theta$  plus vraisemblable si  $p_\theta(X_{1:n})$  grand

## Définition : M-estimateur

Soit  $M : \mathcal{X}^n \times \mathcal{A} \rightarrow \mathbb{R}^+ \cup \{+\infty\}$  un contraste et

$$\operatorname{argmin}_{t \in \mathcal{A}} M(x, t) = \left\{ t \in \mathcal{A} : \forall t', M(x, t') \geq M(x, t) \right\}.$$

Un M-estimateur est une statistique  $\hat{g}(X)$  telle que

$$\hat{g}(X) = \operatorname{argmin}_{t \in \mathcal{A}} M(X, t),$$

pour un contraste  $M$  admettant un unique minimiseur en  $t$ .

- **Notations** : pour  $f : \mathcal{A} \mapsto \mathbb{R}$ ,  
 $\operatorname{argmin}_{\mathcal{A}} f = \operatorname{argmin}_{t \in \mathcal{A}} f(t) = \{t \in \mathcal{A} : \forall t' \in \mathcal{A}, f(t') \geq f(t)\}$ .  
Lorsque  $\operatorname{argmin} f = \{t_0\}$ , on écrit pour simplifier  $\operatorname{argmin} f = t_0$ .
- La définition suppose l'existence et l'unicité du minimum.
- C'est le cas Si  $M$  est strictement convexe en  $t$ .

# Z-estimateur

- Si  $\hat{g}(X)$  est un M-estimateur et si le contraste  $M$  est différentiable p.r.à.  $t$ , on a  $\nabla_t M(X, \hat{g}(X)) = 0$ .  
→  $\hat{g}(X)$  est un **Zéro** de  $\nabla_t M(X, \cdot)$ .

## définition : Z-estimateur

Soit  $\Psi : \mathcal{X}^n \times \mathcal{A} \rightarrow \mathbb{R}^d$  telle que

$$\forall x \in \mathcal{X}^n, \exists ! \hat{g}(x) \text{ tel que } \Psi(x, \hat{g}(x)) = 0.$$

La statistique  $\hat{g}(X)$  est alors appelée Z-estimateur.



## Question

Définitions d'un M- et d'un Z- estimateurs très générales, les propriétés de  $\hat{g}$  dépendent du choix de  $M$  ou  $\Psi$ .

**Comment choisir  $M$  ou  $\Psi$  pour 'bien' estimer  $g(\theta)$  ?**

- Dans ce cours : pas de réponse absolue.
- On donne des exemples de construction et on vérifiera qu'elles ont de bonnes propriétés pour le coût quadratique, à taille d'échantillon  $n$  fixé.
- Propriétés asymptotiques : cf. le cours de statistiques asymptotiques (MACS 203, P2)

1. M- et Z- estimation : exemple et cadre général

2. Maximum de vraisemblance

3. Méthode des moindres carrés

4. Méthode des moments

## justification II (heuristique) de l'estimateur de max de vraisemblance

- $M(x, t) = -\log p_t(x)$ . Si  $\hat{\theta}_{MV}(X) = \operatorname{argmin}_{t \in \Theta} M(X, t)$  est unique,  $\hat{\theta}$  est l'estimateur du maximum de vraisemblance pour le paramètre  $\theta$ .
- raison du bon comportement de  $\hat{\theta}_{MV}$  : Si  $X_{1:n} \stackrel{\text{i.i.d.}}{\sim} P_{\theta_0}$ ,

$$\begin{aligned}\hat{\theta}(X_{1:n}) &= \operatorname{argmin}_t \sum_{i=1}^n -\log p_t(X_i) = \operatorname{argmin}_t \frac{1}{n} \sum_{i=1}^n \log \frac{p_{\theta_0}(X_i)}{p_t(X_i)} \\ &\approx_{n \rightarrow \infty} \operatorname{argmin}_t \mathbb{E}_{\theta_0} \log \frac{p_{\theta_0}(X_1)}{p_t(X_1)} = \underbrace{\int_{\mathcal{X}} p_{\theta_0}(x) \log \frac{p_{\theta_0}(X_1)}{p_t(X_1)} dx}_{KL(P_{\theta_0}, P_t)}\end{aligned}$$

- $KL(P_{\theta_0}, P_t) \geq 0$  mesure la divergence entre  $P_{\theta_0}$  et  $P_t$ .
- cas d'égalité :  $P_t = P_{\theta_0}$  i.e. (si modèle identifiable)  $t = \theta_0$ .
- Justification du ' $\approx_{n \rightarrow \infty}$ ' : cours de stats asymptotiques.

## Max de vraisemblance : Exemple II

- $X_{1:n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}oiss(\theta), \theta \in \Theta = \mathbb{R}_+^*.$
- Pour  $t > 0, -\log p_t(X_{1:n}) = \dots$

**au tableau**

- Résultat :  $\hat{\theta}_{MV}(X_{1:n}) = \frac{1}{n} \sum_{i=1}^n X_i$
- Pourquoi est-ce rassurant ?

## Limites de l'estimateur de maximum de vraisemblance

- Souvent pas d'expression explicite pour  $\hat{\theta}_{MV}$
- Alors : recours obligatoire à des méthodes d'optimisation numérique
- → Coûteux en temps de calcul et pas exact.

1. M- et Z- estimation : exemple et cadre général

2. Maximum de vraisemblance

3. Méthode des moindres carrés

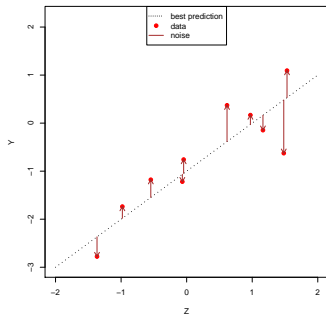
4. Méthode des moments

# Cadre de la régression

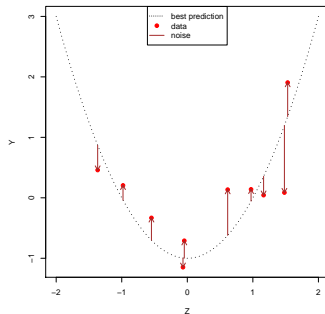
- Observations :  $X_i = (Y_i, z_i)$ ,  $Y_i \in \mathbb{R}$  (aléatoire),  $z_i \in \mathbb{R}^d$  (donnée du problème, non aléatoire), telles que

$$Y_i = \varphi(\theta, z_i) + \epsilon_i, \quad \epsilon_{1:n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \text{ (bruit)} \quad \theta \text{ à estimer}$$

- Cas fréquent : régression linéaire,  
 $\varphi(\theta, z_i) = \langle \theta, \Phi(z_i) \rangle = \sum_{j=1}^d \theta_j \Phi(z_i)_j$ .



$$\varphi(\theta, z) = \theta_0 + \theta_1 z$$



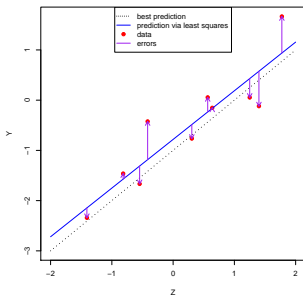
$$\varphi(\theta, z) = \theta_0 + \theta_1 z^2$$

# Estimateur des moindres carrés

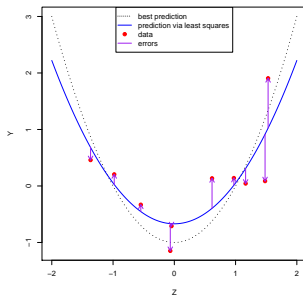
Méthode très ancienne (Gauss).

- Contraste : erreur quadratique entre les  $Y_i$  et leur meilleure prédiction  $\varphi(z_i, t)$  :

$$M(X_{1:n}, t) = \sum_{i=1}^n (\varphi(z_i, t) - Y_i)^2 \quad \hat{\theta}_{MC}(X) = \operatorname{argmin}_{t \in \Theta} M(X_{1:n}, t).$$



$$\varphi(\theta, z) = \theta_0 + \theta_1 z$$



$$\varphi(\theta, z) = \theta_0 + \theta_1 z^2$$



1. M- et Z- estimation : exemple et cadre général
2. Maximum de vraisemblance
3. Méthode des moindres carrés
4. Méthode des moments

# Méthode des moments : principe de substitution

But : estimer  $\theta$ . Supposons que

- on dispose d'une fonction  $h = (h_1, \dots, h_p) : \mathcal{X}^n \rightarrow \mathbb{R}^p$  telle que  $\Phi(\theta) := \mathbb{E}_\theta[h(X)]$  soit calculable (en fonction de  $\theta$ )
- on peut retrouver  $\theta$  à partir de  $\Phi(\theta)$ , i.e.  $\theta \mapsto \Phi(\theta)$  est injective. Alors  $\exists \Phi^{-1} : \text{Im}(\Phi) \subset \mathbb{R}^p \rightarrow \Theta$ .

## principe de substitution

Remplacer  $\Phi(\theta) = \mathbb{E}_\theta(h(X))$  (inconnu car  $\theta$  inconnu) par

$$\Phi_n(X_{1:n}) := \frac{1}{n} \sum_{i=1}^n h(X_i)$$

- (si  $\Phi_n(X) \in \text{Im}(\Phi)$ ), on pose  $\hat{\theta}(X_{1:n}) = \Phi^{-1} \circ \Phi_n(X_{1:n})$
- cas général : on minimise le contraste  $M(X_{1:n}, t) = \|\Phi_n(X_{1:n}) - \Phi(t)\|$

# Principe de substitution et minimisation de contraste

- On définit le contraste

$$M(X_{1:n}, t) = \|\Phi_n(X_{1:n}) - \Phi(t)\|$$

- Si  $\exists!$  minimiseur, l'estimateur par la méthode des moments est

$$\hat{\theta}(X_{1:n}) = \operatorname{argmin}_t M(X_{1:n}, t).$$

## Lemme : Condition suffisante pour que $\exists!$ minimiseur

Sous l'hypothèse d'injectivité de  $\theta \mapsto \Phi(\theta)$ , s'il existe  $t^*$  tel que  $M(X_{1:n}, t^*) = 0$ , alors  $t^*$  est l'unique minimiseur de  $M$

- Sous l'hypothèse d'injectivité, si  $\Phi_n(X_{1:n}) \in \operatorname{Im}(\Phi)$ , le lemme s'applique

## Exemple I : paramètre d'une loi Gamma

- $\theta = (\alpha, \lambda) := (\theta_1, \theta_2), \alpha > 0, \lambda > 0$ .
- $X_{1:n} \stackrel{\text{i.i.d.}}{\sim} P_\theta = \mathcal{Gamma}(\alpha, \lambda)$ .
- Modèle dominé par la mesure de Lebesgue, densité

$$p_{(\alpha, \lambda)}(x) = \mathbb{1}_{x>0} \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}$$

- On choisit  $h(X) = (X, X^2)$  (méthode des *moments*)
- On montre que  $\Phi(\theta) := \mathbb{E}_\theta(h(X)) = \left(\frac{\theta_1}{\theta_2}, \frac{\theta_1(1+\theta_1)}{\theta_2^2}\right) := (m_1, m_2)$
- sur  $\text{Im}(\Phi) = \{(m_1, m_2) : m_1 > 0, m_2 > m_1^2\}$ ,

$$\Phi^{-1}(m) = \left(\frac{m_1^2}{m_2 - m_1^2}, \frac{m_1}{m_2 - m_1^2}\right).$$

## Exemple I : paramètre d'une loi Gamma (suite)

- Contraste :  $M(X_{1:n}, \alpha, \lambda) = \|\Phi_n(X_{1:n}) - \Phi(\alpha, \lambda)\|$  avec

$$\Phi_n(X_{1:n}) = \left( \frac{1}{n} \sum_i X_i, \frac{1}{n} \sum_i X_i^2 \right)$$

- on montre que  $\Phi_n(X_{1:n}) \in \text{Im}(\Phi) \rightarrow$  le lemme s'applique
- on obtient

$$\hat{\theta}_M(X_{1:n}) = \Phi^{-1}(\Phi_n(X_{1:n})) = \left( \frac{\overline{X_n}^2}{\widehat{\sigma}_n^2}, \frac{\overline{X_n}}{\widehat{\sigma}_n^2} \right)$$

$$\text{avec } \overline{X_n} = \frac{1}{n} \sum_i X_i, \widehat{\sigma}_n^2 = \frac{1}{n} \sum_i X_i - \overline{X_n}^2.$$

## Exemple II : paramètres d'une loi normale

- $\theta = (\mu, \sigma^2), X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$
- on choisit  $h(X) = (X, X^2)$ . On a immédiatement

$$\Phi(\theta) = \mathbb{E}_\theta(h(X)) = (\mu, \mu^2 + \sigma^2); \quad \text{Im}(\Phi) = \{(m_1, m_2) : m_2 > m_1^2\}$$
$$\Phi^{-1}(m) = (m_1, m_2 - m_1^2)$$

- On vérifie que  $\Phi_n(X_{1:n}) = \left(\frac{1}{n} \sum_i X_i, \frac{1}{n} \sum_i X_i^2\right) \in \text{Im}(\Phi)$  (comme précédemment)
- On peut poser  $\hat{\theta}_M(X_{1:n}) = \Phi^{-1}(\Phi_n(X))$ ;

$$\hat{\theta}_M(X_{1:n}) = (\hat{\mu}_M, \hat{\sigma}_M^2) = \left(\frac{1}{n} \sum_i X_i, \frac{1}{n} \sum_i X_i^2 - \left(\frac{1}{n} \sum_i X_i\right)^2\right)$$

(moyenne et variance empiriques)