

MDI220, Statistique

Cours 2: Estimation ponctuelle

Anne Sabourin

Septembre 2019

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

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

Maximum de vraisemblance : définition

- Modèle dominé (par la mesure de Lebesgue) sur \mathbb{R} : P_θ a une densité $p(x; \theta)$.
- à $x = (x_1, \dots, x_n)$ fixé, $t \mapsto p(x; t)$ est la fonction de vraisemblance
- **But** : estimer $g(\theta) = \theta$. (donc $\mathcal{A} = \Theta$).

Définition : Estimateur de maximum de vraisemblance

Supposons que $\forall x \in \mathcal{X}^n, \exists ! \hat{\theta}_{MV}(x)$ tel que

$$\forall t \in \Theta, p(x; t) \leq p(x; \hat{\theta}_{MV}).$$

L'estimateur $\hat{\theta}_{MV}$ est appelé *Estimateur du maximum de vraisemblance*

- **remarque** $\hat{\theta}_{MV}$ dépend des données $x = (x_1, \dots, x_n)$ (et seulement des données), c'est donc bien un *estimateur*

L'estimateur du max de vraisemblance est obtenu en résolvant un problème de minimisation

- On pose $M(x, t) = -\log p(x; t)$, $t \in \Theta$ (une fonction à minimiser par rapport à t). M est appelée 'fonction de contraste' ou 'contraste'.
- **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 $\operatorname{argmin} f = t_0$.
- Avec ces notations, on a

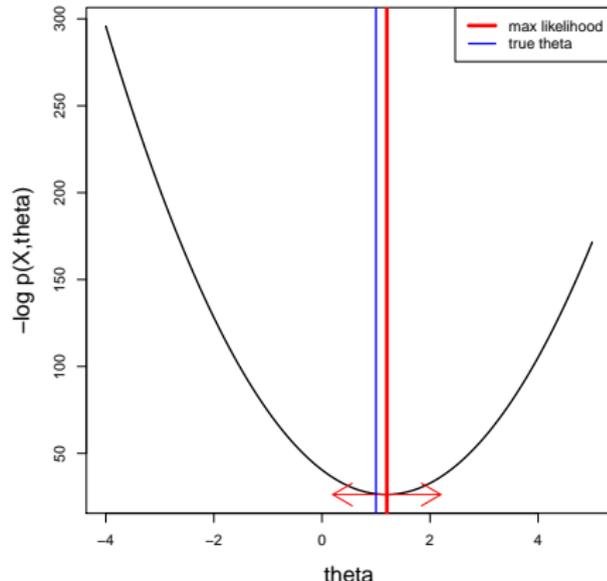
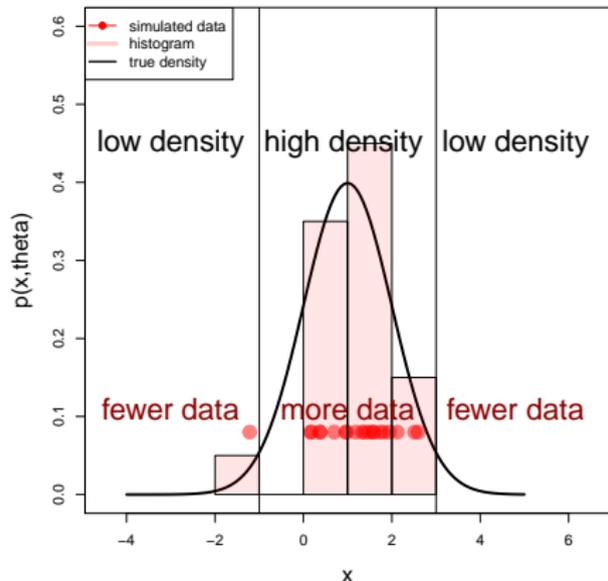
$$\hat{\theta}(x) = \operatorname{argmin}_{t \in \Theta} M(x, t),$$

- **exemple** : modèle gaussien $\mathcal{N}(\theta, \sigma^2)$, σ^2 connu.

$$-\log p(x; \theta) =? \quad \hat{\theta}(x) = ?$$

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

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

justification II de l'estimateur de max de vraisemblance

$\hat{\theta}_{MV}(X) = \operatorname{argmin}_{t \in \Theta} \{-\log p(X; t)\}$. On suppose $X = X_{1:n} \stackrel{\text{i.i.d.}}{\sim} P_{\theta_0}$,

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

- $KL(P_{\theta_0}, P_t) \geq 0$ mesure la divergence entre P_{θ_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

- $X_{1:n} \stackrel{\text{i.i.d.}}{\sim} \text{Poiss}(\theta), \theta \in \Theta = \mathbb{R}_+^*$.
- 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.

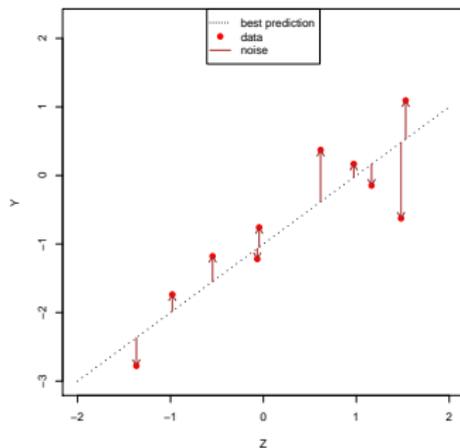
1. Maximum de vraisemblance
2. Méthode des moindres carrés
3. Méthode des moments
4. M- et Z- estimation : cadre général

Cadre de la régression

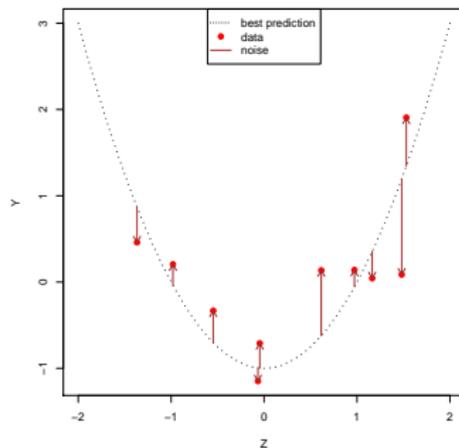
- Observations : $X_i = (Y_i, z_i)$, $Y_i \in \mathbb{R}$ (aléatoire), $z_i \in \mathbb{R}^d$ (donnée du problème, non aléatoire), telles que

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

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



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



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

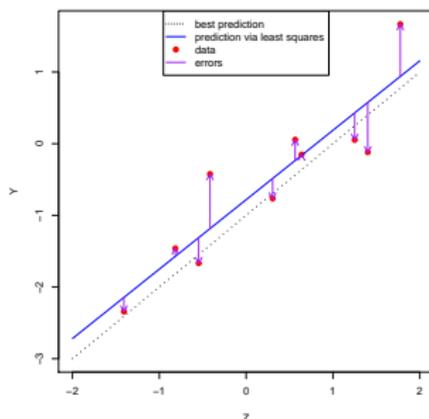
Estimateur des moindres carrés

Méthode très ancienne (Gauss).

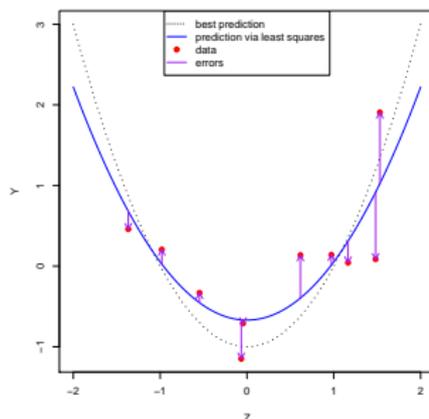
- Contraste : erreur quadratique entre les Y_i et leur meilleure prédiction $f(z_i, t)$:

$$M(X_{1:n}, t) = \sum_{i=1}^n (f(z_i, t) - Y_i)^2$$

$$\hat{\theta}_{MC}(X) = \operatorname{argmin}_{t \in \Theta} M(X_{1:n}, t).$$



$$f(\theta, z) = \theta_0 + \theta_1 z$$
$$\Phi(z) = z$$



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

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

Méthode des moments : principe de substitution

But : estimer θ . Supposons que $\mathbb{E}_\theta(X_1)$ ('moment 'ordre 1') soit une fonction connue de θ , et injective.

principe de substitution

Remplacer $\Phi(\theta) = \mathbb{E}_\theta(X_1)$ (inconnu car θ inconnu) par

$$\Phi_n(X_{1:n}) := \frac{1}{n} \sum_{i=1}^n X_i \quad (\text{moyenne empirique})$$

Définition

L'estimateur de θ par la méthode des moments est l'unique solution $\hat{\theta}_m$ de l'équation (en θ)

$$\Phi(\theta) = \Phi_n(X_{1:n})$$

c'est à dire, $\mathbb{E}_{\hat{\theta}_m}(X) = \frac{1}{n} \sum_{i=1}^n X_i$.

Exemple : estimation du paramètre de localisation d'un loi normale

Supposons $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta, \sigma^2)$, σ^2 connu.

- $\Phi(\theta) = \mathbb{E}_\theta(X_1) = ??$
- Estimateur des moments $\hat{\theta}_m = ??$
- **Problème** : que faire lorsque l'on veut aussi estimer σ^2 ?

→ **idée** : utiliser le moment d'ordre 2, $\frac{1}{n} \sum X_i^2$, en complément de la moyenne empirique

Méthode des moments : généralisation

On se donne 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 θ)
- on peut retrouver θ 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 θ inconnu) par

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

- si $\exists \theta$ t.q. $\Phi(\theta) = \Phi_n(X)$, on pose $\hat{\theta}_m(X_{1:n}) = \Phi^{-1} \circ \Phi_n(X_{1:n})$
- sinon (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

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

définition : estimateur par la méthode des moments

Si $\exists!$ minimiseur de $M(X_{1:n}, \cdot)$, l'estimateur de θ par la méthode des moments est

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

Lemme : Condition suffisante d'existence et unicité

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

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

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

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)

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

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.

Bilan

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

→ $\hat{\theta}$ est un arg min

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

→ $\hat{\theta}$ est une racine.

Question

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

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

- Dans ce cours : pas de réponse absolue.
- On a donné des exemples de construction (max de vraisemblance, moindres carrés, méthode des moments) 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)