

# Statistiques mathématiques

R. Petit

année académique 2016 - 2017

# Table des matières

<b>1</b>	<b>Théorie de l'échantillonnage</b>	<b>2</b>
1.1	Terminologie et définitions . . . . .	2
1.2	Moments . . . . .	3
1.2.1	Indicateurs . . . . .	4
1.3	Quantile . . . . .	4
1.3.1	Lemme de Fisher . . . . .	6
<b>2</b>	<b>Estimation ponctuelle</b>	<b>8</b>
2.1	Introduction . . . . .	8
2.2	Critères d'estimation . . . . .	8
2.2.1	Définitions de convergence . . . . .	8
2.2.2	Résultats élémentaires sur les convergences . . . . .	9
2.2.3	Estimateurs convergents . . . . .	10
2.3	Estimateur exhaustif . . . . .	11
2.3.1	Estimateurs non biaisés . . . . .	13
2.3.2	Estimateurs à dispersion minimale . . . . .	14
2.3.3	Estimateurs efficaces . . . . .	19
2.4	Méthodes d'estimation . . . . .	23
2.4.1	Méthode des moments . . . . .	23
2.4.2	Méthode du maximum de vraisemblance . . . . .	25

# Introduction

En probabilités, une variable aléatoire  $X$  donnée est entièrement définie par sa loi. On peut l'exprimer par la fonction de répartition  $F^X$  ou par la fonction de densité  $f^X = \frac{d}{dx} F^X$ . Ces fonctions permettent de déterminer :

$$\mathbb{P}[a \leq X \leq b] = \int_a^b f^X(x) dx = F^X(b) - F^X(a).$$

Ou encore :

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x f^X(x) dx.$$

Cependant, les fonctions  $f^X$  et  $F^X$  ne sont jamais connues précisément. Elles peuvent être approchées par des modélisations, mais les modèles ne sont jamais exacts. En probabilités, on cherche donc les observations sur base de la loi qui est connue, alors qu'en statistiques, on cherche à retrouver la loi sur base de  $n$  observations  $X_1, \dots, X_n$ .

Nous allons nous intéresser à des *modèles statistiques* sous la forme  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$  où :

$$\mathcal{P}^{(n)} = \left\{ \mathbb{P}^{(n)} \right\} = \left\{ \mathbb{P}_\theta^{(n)} \text{ t.q. } \theta \in \Theta \subset \mathbb{R}^k \right\},$$

et donc les  $\mathbb{P}^{(i)}$  sont chacun une loi possible pour  $(X_1, \dots, X_n)$ .

Ces modèles sont dits *paramétriques* car les différentes lois sont les mêmes au paramètre  $\theta$  près. Nous n'étudierons que des modèles paramétriques où  $\Theta$  est un espace de dimension  $d \in \mathbb{N}$  finie.

*Exemple 0.1.* Soient  $X_1, \dots, X_n$  des variables aléatoires iid (indépendantes et identiquement distribuées).

— Si les  $X_i$  sont de loi normale  $\mathcal{N}(\mu, \sigma^2)$ , alors le paramètre  $\theta$  est donné par :

$$\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \in \Theta = \mathbb{R} \times \mathbb{R}^+ \subset \mathbb{R}^2;$$

— si les  $X_i$  sont de loi uniforme  $\text{Unif}(0, \theta)$ , le paramètre  $\theta$  est donné par  $\theta \in \Theta = \mathbb{R}_0^+ \subset \mathbb{R}$ ;

— si les  $X_i$  sont de loi Bern(p), le paramètre  $\theta$  est donné par  $\theta = p \in \Theta = [0, 1] \subset \mathbb{R}$ .

*Remarque.* Une loi normale  $\mathcal{N}(\mu, \sigma^2)$  est déraisonnable car les valeurs observables ne vont empiriquement pas vers les infinis alors que la distribution le permet théoriquement mais n'est pas **complètement** déraisonnable car ces probabilités sont négligeables grâce à l'exponentielle de  $(-x^2)$  dans la formule de la densité.

# Chapitre 1

## Théorie de l'échantillonnage

### 1.1 Terminologie et définitions

**Définition 1.1.** On appelle *modèle d'échantillonnage* un modèle d'observations iid.

**Définition 1.2.** Soit un modèle statistique  $(E^n, \mathcal{B}(E^n), \mathcal{P}^{(n)})$  où  $\mathcal{P}^{(n)} = \{\mathbb{P}_\theta^{(n)} \text{ t.q. } \theta \in \Theta \subset \mathbb{R}^k\}$ . On note ici  $\mathbb{P}_\theta^{(n)}$  une loi possible pour  $(X_1, \dots, X_n)$  et  $\mathbb{P}_\theta$  une loi possible pour  $X_i$  avec  $i$  fixé. On dit alors que  $\mathbb{P}_\theta^{(n)}$  est déterminé par  $\mathbb{P}_\theta$ .

*Remarque.* Ici, deux visions vont s'opposer et se compléter : la vision *population* qui est associée à  $\mathbb{P}_\theta$  et la version *échantillonnage* (ou *empirique*), qui, elle, est associée à  $\mathbb{P}_\theta^{(n)}$ .

**Définition 1.3.** On définit la fonction indicatrice  $I_{[\cdot]}$  qui vaut 1 quand l'expression entre crochets est vraie et 0 sinon.

**Définition 1.4.** Soit  $X_1, \dots, X_n$  une suite de  $n$  observations. On définit la *ième statistique d'ordre* par  $X_{(i)} = X_k$  t.q.  $|\{X_j \text{ t.q. } X_j < X_k, 1 \leq j \leq n\}| = i$ . On définit également la *statistique d'ordre* par  $(X_{(i)})_i$ .

**Définition 1.5.** On définit les fonctions de répartition comme suit :

— la fonction de répartition population :

$$F_\theta(x) = \mathbb{P}_\theta[X_i \leq x] ;$$

— la fonction de répartition empirique :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[X_i \leq x]}.$$

*Remarque.* La fonction  $F_n$  empirique est une fonction en escaliers. Elle fait des sauts de hauteur  $\frac{1}{n}$ , et est telle que :

$$\lim_{x \rightarrow +\infty} F_n(x) = 1 \quad \text{et} \quad \lim_{x \rightarrow -\infty} F_n(x) = 0.$$

On peut également remarquer que  $F_n(X_{(i)}) = \frac{i}{n}$ . En effet, par définition de  $X_{(i)}$ , il y a exactement  $i$  observations inférieures à  $X_{(i)}$ . Dès lors, la fonction indicatrice donnera  $i$  fois la valeur 1 et  $(n - i)$  fois la valeur 0. La somme donne donc  $i$  et la fonction donne  $\frac{i}{n}$ .

**Définition 1.6.** On appelle *statistique* toute fonction mesurable faisant intervenir **uniquement** des observations.

*Exemple 1.1.* Par exemple  $F_n$  est une statistique car seules les valeurs  $X_i$  sont utilisées, mais  $F_\theta$  n'est pas une statistique car la valeur du paramètre  $\theta$  apparaît et n'est pas une observation.

*Remarque.* Une statistique peut être à valeur scalaire ( $X_{(i)}$  par exemple), à valeur vectorielle ( $(X_{(i)})_{1 \leq i \leq n}$  par exemple), à valeur ensembliste ( $[X_i \pm \bar{X}]$  avec  $i$  fixé par exemple), ou encore à valeur fonctionnelle ( $F_n$  par exemple).

*Remarque.* L'objectif est de pouvoir approximer la loi régissant les populations ( $F_\theta$ ) à l'aide de la loi observée empiriquement. Par la loi des grands nombres, on a :

$$F_n(x) \xrightarrow[n \rightarrow +\infty]{\text{p.s. par } \mathbb{P}_\theta} F_\theta(x).$$

**Théorème 1.7** (Théorème de Glivenko-Cantelli). Si  $F_n$  et  $F_\theta$  sont respectivement une fonction de répartition empirique et de population, alors :

$$\sup_{x \in \mathbb{R}} |F_n(x) - F_\theta(x)| \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 0$$

## 1.2 Moments

**Définition 1.8** (Moments pour populations). On définit  $\mu'_r(\theta)$  le *moment non-centré* d'ordre  $r$  avec  $r \in \mathbb{N}^*$  par :

$$\mu'_r(\theta) := E_\theta[X_1^r].$$

On définit également  $\mu_r(\theta)$ , le *moment centré* d'ordre  $r$  avec  $r \in \mathbb{N}^*$  par :

$$\mu_r(\theta) := E_\theta \left[ (X_1 - \mu'_r(\theta))^r \right].$$

**Définition 1.9** (Moments pour échantillon). On définit  $m'_r$ , le *moment non-centré* d'ordre  $r$  avec  $r \in \mathbb{N}^*$  par :

$$m'_r := \frac{1}{n} \sum_{i=1}^n X_i^r.$$

On définit également le *moment centré* d'ordre  $r$  avec  $r \in \mathbb{N}^*$  par :

$$m_r := \frac{1}{n} \sum_{i=1}^n (X_i - m'_r)^r.$$

*Remarque.* La loi des grands nombres dit que :

$$m'_r \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mu'_r(\theta),$$

mais on ne peut pas dire que :

$$m_r \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mu_r(\theta).$$

Ce n'est donc pas possible car pour  $m'_r$ , il y a une somme de variables iid alors que pour  $m_r$ , les variables sommées ne sont pas iid (mais dépendent toutes de tous les  $X_i$ ).

En réalité, il y a convergence, mais on ne peut pas l'exprimer de manière triviale par la loi des grands nombres.

### 1.2.1 Indicateurs

On peut observer que  $\mu'_1(\theta) = \mathbb{E}_\theta[X_1]$ . Pareil pour  $m'_1 = \bar{X}$ . Le moment d'ordre 1 est donc un indice de position. On a alors  $\mu := \mu_1(\theta) = \mathbb{E}[(X - \mathbb{E}[X_1])] = \mathbb{E}[X_1] - \mathbb{E}[X_1] = 0$ . Cette valeur n'est donc pas intéressante. Par contre :

$$\mu_2(\theta) = \mathbb{E}[(X_1 - \mathbb{E}[X_1])^2] =: \text{Var}(X) \quad \text{si} \quad m_2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 =: s^2.$$

Le moment d'ordre 2 est donc un indice de dispersion.

**Définition 1.10.** On appelle le *coefficient d'asymétrie de Fisher* la quantité :

$$\gamma_1 := \mu_3(\theta) \cdot (\mu_2(\theta))^{-\frac{3}{2}}.$$

*Remarque.* Le dénominateur  $\mu_2(\theta)^{\frac{3}{2}}$  apparait afin de rendre invariant le coefficient d'asymétrie de Fisher aux transformations affines.

**Définition 1.11.** Le coefficient d'asymétrie de Fisher *empirique* est donné par :

$$m_3 \cdot m_2^{-\frac{3}{2}}.$$

**Définition 1.12.** On appelle *coefficient d'aplatissement de Fisher* la quantité :

$$\gamma_2 := \mu_4(\theta) \cdot (\mu_2(\theta))^{-2} - 3.$$

**Définition 1.13.** Le coefficient d'aplatissement de Fisher *empirique* est donné par :

$$m_4 \cdot m_2^{-2} - 3.$$

*Remarque.* Si  $\gamma_2 \geq 0$ , c'est que les événements extrêmes sont de plus haute probabilité et si  $\gamma_2 \leq 0$ , c'est que les événements extrêmes sont de moins haute probabilité.

À nouveau, le dénominateur y a été ajouté afin de rendre le coefficient invariant aux transformations affines. Et le terme  $-3$  sert à annuler le coefficient d'aplatissement de Fisher pour une normale  $\mathcal{N}(\mu, \sigma^2)$ .

## 1.3 Quantile

**Définition 1.14.** Si  $F_\theta$  est inversible, alors on définit  $x_\alpha(\theta) := F_\theta^{-1}(\alpha)$ , et on appelle  $x_\alpha(\theta)$  un *quantile*.

*Remarque.* Il faut cependant faire attention car on peut avoir le cas de  $F_\theta$  discontinue où on choisit  $\alpha = F_\theta^{-1}$  (point de discontinuité) ou alors le cas de  $F_\theta$  admettant un plateau et où on choisit  $\alpha$  sur le plateau.

**Définition 1.15.** On définit alors :

$$x_\alpha(\theta) := \inf \{x \in \mathbb{R} \text{ t.q. } F_\theta(x) \geq \alpha\}.$$

*Remarque.* On donne les noms de *médiane*, *quartile*, *décile*, *percentile* pour  $\alpha$  valant, avec  $k$  entier, respectivement  $\frac{1}{2}$ ,  $\frac{k}{4}$  avec  $k < 4$ ,  $\frac{k}{10}$  avec  $k < 10$ , et  $\frac{k}{100}$  avec  $k < 100$ .

**Définition 1.16.** Pour les échantillons, on définit le *quantile empirique d'ordre  $\alpha$*  par :

$$x_\alpha^{(n)} := \inf \{x \in \mathbb{R} \text{ t.q. } F_n(x) \geq \alpha\}.$$

*Remarque.* On peut également définir des indices de position, dispersion, asymétrie, aplatissement, etc. sur les quantiles plutôt que sur les moments. Ils auront des propriétés différentes et une robustesse différente aux valeurs aberrantes.

**Définition 1.17.** La loi échantillonnée de  $T(X^{(n)})$  est la loi déterminée par :

$$\mathbb{P}_\theta^{(n)} \left[ T(X^{(n)}) \in B \right] = \mathbb{P}_\theta^{(n)} \left[ \left\{ x^{(n)} \in X^{(n)} \text{ t.q. } T(x^{(n)}) \in B \right\} \right], B \in \mathcal{B}(\mathbb{R}^m).$$

*Exemple 1.2 (Bernoulli).*  $X^{(n)} = (X_1, \dots, X_n)$  où les  $X_i$  sont iid Bern( $p$ ). On a alors :  $T(X^{(n)}) = \sum_{i=1}^n X_i$ , sous  $\mathbb{P}_\theta^{(n)}$ , est de loi Bin( $n, p$ ).

*Exemple 1.3 (Normale).*  $X^{(n)} = (X_1, \dots, X_n)$  où les  $X_i$  sont iid  $\mathcal{N}(\mu, \sigma^2)$  et où  $\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \in \Theta = \mathbb{R} \times \mathbb{R}_0^+ \subset \mathbb{R}^2$ .

La statistique  $T_1(X^{(n)}) = \sum_{i=1}^n X_i$ , sous  $\mathbb{P}_\theta^{(n)}$ , est de loi  $\mathcal{N}(n\mu, n\sigma^2)$ .

La statistique  $T_2(X^{(n)}) = \frac{1}{n} \sum_{i=1}^n X_i$ , sous  $\mathbb{P}_\theta^{(n)}$ , est de loi  $\mathcal{N}(\mu, \frac{\sigma^2}{n})$ .

*Exemple 1.4 (Uniforme).*  $X^{(n)} = (X_1, \dots, X_n)$  où les  $X_i$  sont iid Unif( $0, \theta$ ), pour  $\theta \in \Theta = \mathbb{R}_0^+ \subset \mathbb{R}$ . On a donc  $f_\theta^{X_i}(x) = \theta^{-1} I_{[0 \leq x \leq \theta]}$ . Et donc :

$$F_\theta^{X_i}(x) = \begin{cases} 0 & \text{si } x < 0 \\ \frac{x}{\theta} & \text{si } 0 \leq x \leq \theta \\ 1 & \text{sinon} \end{cases}.$$

La statistique  $T(X^{(n)}) = X_{(n)} = \max_{1 \leq k \leq n} \{X_k\}$  a pour fonction de répartition, sous  $\mathbb{P}_\theta^{(n)}$  :

$$F_\theta^{(n)}(x) = \mathbb{P}[X_{(n)} \leq x] = \mathbb{P}[X_1 \leq x, X_2 \leq x, \dots, X_n \leq x].$$

La seconde forme est plus agréable car on a une intersection d'événements indépendants. Donc :

$$F_\theta^{(n)}(x) = \prod_{i=1}^n \mathbb{P}[X_i \leq x] = \prod_{i=1}^n F_\theta^{X_i}(x) = \begin{cases} 0 & \text{si } x < 0 \\ \left(\frac{x}{\theta}\right)^n & \text{si } 0 \leq x \leq \theta \\ 0 & \text{sinon} \end{cases}.$$

On a alors la fonction de densité :

$$\begin{aligned} f_\theta^{X_{(n)}}(x) &= \frac{d}{dx} F_\theta^{X_{(n)}} \Big|_x = \begin{cases} 0 & \text{si } x < 0 \\ \frac{nx^{n-1}}{\theta^n} I_{[0 \leq x \leq \theta]} & \text{si } 0 \leq x \leq \theta \\ 0 & \text{sinon} \end{cases} \\ &= \frac{nx^{n-1}}{\theta^n} I_{[0 \leq x \leq \theta]}. \end{aligned}$$

*Remarque.* La loi échantillonnée n'est pas toujours possible à déterminer exactement analytiquement. Dans ce cas, on donne :

- (i) les/des moments de la loi échantillonnée exacte ;
- (ii) la loi échantillonnée asymptotique.

Et pour de grandes valeurs de  $n$ , la loi asymptotique donne une assez bonne approximation de la loi exacte.

*Remarque.* Ici, les termes *exact* et *asymptotique* s'opposent : on parle d'objet *exact* lorsque l'objet est connu pour  $n$  fixé, et d'objet *asymptotique* lorsque l'objet n'est connu que pour  $n \rightarrow +\infty$ .

*Exemple 1.5.* Voici un cas où on ne peut exprimer de loi exacte mais où il est possible d'exprimer une loi asymptotique. Soit  $X^{(n)} = (X_1, \dots, X_n)$  où les  $X_i$  sont iid  $F$  avec la fonction  $F$  telle que  $\text{Var}_F(X_i) = \sigma^2 < +\infty$  et donc  $E_F(X_i) = \mu < +\infty$ . On peut dès lors appliquer le théorème central limite (TCL) :

$$\sqrt{n}(\bar{X}^{(n)} - \mu) \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \mathcal{N}(0, \sigma^2).$$

Pour  $n \gg$ , on peut alors dire :

$$\bar{X}^{(n)} \approx \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right),$$

où le symbole  $\approx$  se lit *est à peu près de même loi*.

On en conclut donc qu'avec  $n$  suffisamment grand, on peut approximer  $\bar{X}^{(n)}$ , même sans connaître sa loi exacte.

### 1.3.1 Lemme de Fisher

**Définition 1.18.** La variable aléatoire  $Q$  est de loi  $\chi^2$  (chi-carrée) à  $k(\in \mathbb{N}^*)$  degrés de liberté lorsque :

$$Q \stackrel{\mathcal{D}}{=} \sum_{i=1}^k Z_i^2,$$

où les  $Z_i$  sont iid  $\mathcal{N}(0, 1)$  et où «  $\stackrel{\mathcal{D}}{=}$  » veut dire *à la même distribution que*. Cela se note :

$$Q \sim \chi_k^2$$

*Remarque.* Si  $Q \sim \chi_k^2$ , alors :

$$f^Q(x) = \frac{1}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)} x^{\frac{k}{2}-1} \exp\left(-\frac{x}{2}\right) I_{[x>0]},$$

où  $\Gamma$  est la fonction Gamma d'Euler définie par :

$$\Gamma(x) = \int_0^{+\infty} t^{x-1} \exp(-t) dt.$$

De plus,  $\text{Var}(Q) = 2k$ , et  $E(Q) = k$ .

On peut également noter que les  $\chi^2$  sont stables par la somme : si  $Q_1 \sim \chi_{k_1}^2$  et  $Q_2 \sim \chi_{k_2}^2$ , alors :

$$Q_1 + Q_2 \sim \chi_{k_1+k_2}^2.$$

**Lemme 1.19.** Soit  $W = (W_1, \dots, W_k)$  un vecteur de variables aléatoires, où  $f^W : \mathbb{R}^k \rightarrow \mathbb{R}^+$  est la fonction de densité du vecteur  $W$ . Alors :

1.  $\mathbb{P}[W \in B] = \int_B f^W(x) dx$ ;
2. si  $V = AW + b$  où  $A$  est une matrice  $k \times k$  inversible, alors :

$$f^V(v) = \left| \det A^{-1} \right| f^W\left(A^{-1}(v - b)\right).$$

**Théorème 1.20** (Lemme de Fisher). Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$  où  $n \geq 2$ . Alors :

- (i)  $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ ;
- (ii)  $\frac{ns^2}{\sigma^2} \sim \chi_{n-1}^2$ ;
- (iii)  $\bar{X} \perp s^2$ .



*Démonstration.* Posons  $Z_i := \frac{X_i - \mu}{\sigma}$  pour  $i \in \llbracket 1, n \rrbracket$ . Puisque les  $X_i$  sont iid, les  $Z_i$  le sont également (même transformation appliquée à tous les  $X_i$  et chaque  $Z_i$  ne fait intervenir que le  $X_i$  correspondant). Notons que :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n (\sigma Z_i + \mu) = \sigma \bar{Z} + \mu,$$

où  $\bar{Z}$  est la moyenne empirique des  $Z_i$ . Notons également que :

$$ns^2 = \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n \left( (\sigma Z_i + \mu) - (\sigma \bar{Z} + \mu) \right)^2 = \sigma^2 \sum_{i=1}^n (Z_i - \bar{Z})^2 = n\sigma^2 s_Z^2.$$

Il nous faut alors montrer que  $\bar{Z} \sim \mathcal{N}(0, 1)$  et  $ns_Z^2 \sim \chi_{n-1}^2$ , avec  $\bar{Z} \sqcup s_Z^2$ .

Pour cela, on sait que le vecteur  $Z^{(n)} = (Z_1, \dots, Z_n)$  a pour densité :

$$f^{Z^{(n)}}(z^{(n)}) = \prod_{i=1}^n f^{Z_i}(z_i) = \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{z_i^2}{2} \right) \right) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left( -\sum_{i=1}^n \frac{z_i^2}{2} \right) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left( -\frac{1}{2} \|z^{(n)}\|^2 \right).$$

Soit  $O$  une matrice orthogonale de dimension  $n \times n$  telle que  $\forall j \in \llbracket 1, n \rrbracket : O_{1j} = \frac{1}{\sqrt{n}}$ . On pose alors :

$$(Y_1, \dots, Y_n) = Y^{(n)} = OZ^{(n)}.$$

Puisque la matrice  $O$  est orthogonale, on sait que  $O^{-1}$  existe et que  $|\det O| = |\det O^{-1}| = 1$ . Par le lemme 1.19, on peut dire :

$$f^{Y^{(n)}}(y^{(n)}) = |\det O^{-1}| f^{Z^{(n)}}(O^{-1}y^{(n)}) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left( -\frac{1}{2} \|O^{-1}y^{(n)}\|^2 \right) = \left( \frac{1}{\sqrt{2\pi}} \right)^n \exp \left( -\frac{1}{2} \|y^{(n)}\|^2 \right).$$

On a donc  $f^{Y^{(n)}} = f^{Z^{(n)}}$ , ce qui implique que les  $Y_i$  sont iid  $\mathcal{N}(0, 1)$ .

En particulier,  $Y_1 = (Y^{(n)})_1 = (OZ^{(n)})_1 = \sum_{i=1}^n O_{1i}Z_i = \sum_{i=1}^n \frac{Z_i}{\sqrt{n}} = \sqrt{n}\bar{Z} \sim \mathcal{N}(0, 1)$ . On peut alors en déduire que  $\bar{Z} \sim \mathcal{N}(0, n^{-1})$ .

Montrons alors que  $ns_Z^2 \sim \chi_{n-1}^2$  :

$$ns_Z^2 = \sum_{i=1}^n (Z_i - \bar{Z})^2 = \sum_{i=1}^n Z_i^2 - n(\bar{Z})^2 = \|Z^{(n)}\|^2 - (\sqrt{n}\bar{Z})^2 = \|Y^{(n)}\|^2 - Y_1^2 = \sum_{i=2}^n Y_i^2.$$

Or, les  $Y_i$  sont  $\mathcal{N}(0, 1)$ . On a alors bien  $ns_Z^2 \sim \chi_{n-1}^2$  (car la somme sur  $i$  commence à 2, il y a donc  $(n-1)$  variables sommées).

De plus, puisque les  $Y_i$  sont indépendantes deux à deux, que  $\bar{Z}$  ne dépend que de  $Y_1$  et que  $ns_Z^2$  ne dépend pas de  $Y_1$ , on sait que  $\bar{Z} \sqcup ns_Z^2$ .  $\square$

# Chapitre 2

## Estimation ponctuelle

### 2.1 Introduction

Considérons toujours un modèle statistique  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$  avec  $\theta$  le paramètre vectoriel  $\in \Theta \subset \mathbb{R}^k$ .

**Définition 2.1.** Soit  $g : \Theta \rightarrow \mathbb{R}^m$ . Une statistique est appelée *estimateur de  $g(\theta)$*  lorsqu'elle est à valeurs dans  $g(\Theta)$ .

**Définition 2.2.** Soit  $\theta \in \Theta \subset \mathbb{R}^k$ . Si  $g : \Theta \rightarrow \mathbb{R}^m : \theta \mapsto (\theta_{\varphi(1)}, \dots, \theta_{\varphi(m)})$ , on appelle les paramètres  $\theta_{\varphi(i)}$  les paramètres d'intérêt, et on appelle les autres paramètres les paramètres de nuisance.

*Exemple 2.1.* Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$ . On sait  $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_0^+ \subset \mathbb{R}^2$ . Soit  $g : \Theta \rightarrow \mathbb{R} : \theta \mapsto \mu$ .  $\mu$  est le paramètre d'intérêt et  $\sigma^2$  est le paramètre de nuisance.

*Remarque.* Ne pas connaître le paramètre de nuisance induit une *nuisance* pour déterminer le paramètre d'intérêt.

### 2.2 Critères d'estimation

*Remarque.* Afin de définir les estimateurs convergents, il faut définir la notion de convergence, or il n'existe pas une manière canonique de la définir. Il existe donc plusieurs définitions de convergences différentes.

#### 2.2.1 Définitions de convergence

Soient  $Z^{(n)} = (Z_1, \dots, Z_n)$  définis pour  $n \geq 1$  et sur  $(\Omega, \mathcal{F}, \mathbb{P})$ .

**Définition 2.3.** On dit que  $Z^{(n)}$  converge *presque sûrement* (ou *stochastiquement*) vers  $Z$  lorsque :

$$\mathbb{P} \left[ \left\{ \omega \in \Omega \text{ t.q. } Z^{(n)} \xrightarrow[n \rightarrow +\infty]{} Z(\omega) \right\} \right] = 1.$$

Cela se note :

$$Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} Z.$$

**Définition 2.4.** On dit que  $Z^{(n)}$  converge *en probabilités* vers  $Z$  lorsque :

$$\forall \varepsilon > 0 : \mathbb{P} \left[ \left| Z^{(n)} - Z \right| > \varepsilon \right] \xrightarrow[n \rightarrow +\infty]{} 0.$$

Cela se note :

$$Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} Z.$$

**Définition 2.5.** On dit que  $Z^{(n)}$  converge en  $L_r$  vers  $Z$  lorsque :

$$\mathbb{E} \left[ \left| Z^{(n)} - Z \right| \right] \xrightarrow[n \rightarrow +\infty]{} 0.$$

Cela se note :

$$Z^{(n)} \xrightarrow[n \rightarrow +\infty]{L_r} Z.$$

*Remarque.* Lorsque  $r = 2$ , on parle de convergence en moyenne quadratique.

**Définition 2.6.** On dit que  $Z^{(n)}$  converge en loi (ou en distribution) vers  $Z$  lorsque :

$$\forall z \text{ point de continuité de } F^Z : F^{Z^{(n)}}(z) \xrightarrow[n \rightarrow +\infty]{} F^Z(z).$$

Cela se note :

$$Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} Z.$$

*Remarque.* Ces définitions sont faites pour des variables aléatoires réelles mais peuvent être étendues à  $\mathbb{R}^n$  en appliquant la convergence composante par composante.

## 2.2.2 Résultats élémentaires sur les convergences

**Proposition 2.7.** Les convergences sont induites mutuellement par les assertions suivantes :

1. si  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{p.s.} Z$ , alors  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} Z$ ;
2. si  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} Z$ , alors  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} Z$ ;
3. si  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{L_r} Z$ , alors  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} Z$ .

**Théorème 2.8.** Les convergences presque sûre, en probabilités, et en loi sont stables par transformations continues.

**Théorème 2.9.** Notons  $\rightarrow$  une convergence soit presque sûre, soit en probabilités. Si  $Z^{(n)} \rightarrow Z$ , et  $Y^{(n)} \rightarrow Y$ , alors :

- (i)  $Z^{(n)} + Y^{(n)} \rightarrow Z + Y$ ;
- (ii)  $Z^{(n)} \cdot Y^{(n)} \rightarrow Z \cdot Y$ ;
- (iii) si  $\mathbb{P}[Y^{(n)} = 0] = 0$ , alors  $\frac{Z^{(n)}}{Y^{(n)}} \rightarrow \frac{Z}{Y}$ .

**Lemme 2.10** (Lemme de Slutsky). Si  $Z^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} Z$ , et  $Y^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} c \neq 0$ , alors :

- (i)  $Z^{(n)} + Y^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} Z + c$ ;
- (ii)  $Z^{(n)} \cdot Y^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} Z \cdot c$ ;
- (iii)  $\frac{Z^{(n)}}{Y^{(n)}} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \frac{Z}{c}$ .

**Théorème 2.11** (Loi forte des grands nombres). Soient  $Z_1, Z_2, \dots$  iid avec  $\mathbb{E}[Z_1] < +\infty$ . Alors :

$$\bar{Z}^{(n)} = \frac{1}{n} \sum_{k=1}^n Z_k \xrightarrow[n \rightarrow +\infty]{p.s.} \mu = \mathbb{E}[Z_1].$$

**Théorème 2.12** (Loi faible des grands nombres). Soient  $Z_1, Z_2, \dots$  iid avec  $\mathbb{E}[Z_1] < +\infty$ . Alors :

$$\bar{Z}^{(n)} = \frac{1}{n} \sum_{k=1}^n Z_k \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mu = \mathbb{E}[Z_1].$$

**Théorème 2.13** (Théorème central limite (TCL)). Soient  $Z_1, Z_2, \dots$  iid, avec  $\mathbb{E}[Z_1^2] < +\infty$ . Alors :

$$\sqrt{n} \left( Z^{(n)} - \mu \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} W,$$

où :

$$W \sim \mathcal{N}(0, \sigma^2),$$

avec  $\sigma^2 = \text{Var}(Z_1)$ .

### 2.2.3 Estimateurs convergents

**Définition 2.14.** Un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *faiblement convergent* lorsque :

$$\forall \theta \in \Theta : T^{(n)}(X^{(n)}) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} g(\theta) \quad \text{sur } \mathbb{P}_\theta^{(n)}.$$

**Définition 2.15.** Un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *fortement convergent* lorsque :

$$\forall \theta \in \Theta : T^{(n)}(X^{(n)}) \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} g(\theta) \quad \text{sur } \mathbb{P}_\theta^{(n)}.$$

*Exemple 2.2.* Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$ . Prenons  $g(\theta) = \mu$  et  $T^{(n)}(X^{(n)}) = \bar{X}$ . On a bien :

$$T^{(n)}(X^{(n)}) = \bar{X}^{(n)} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mu = \mathbb{E}[X_1] \quad \text{sur } \mathbb{P}_{\mu, \sigma^2}^{(n)}.$$

$T^{(n)}(X^{(n)})$  est donc un estimateur fortement convergent.

Prenons maintenant  $T_2^{(n)}(X^{(n)}) = s^2$ . On ne peut pas appliquer la loi des grands nombres car les variables aléatoires  $(X_i - \bar{X})^2$  sommées ne sont pas indépendantes. On a alors :

$$\begin{aligned} T_2^{(n)}(X^{(n)}) = s^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n \left( (X_i - \mu)^2 + (\mu - \bar{X})^2 - 2(X_i - \bar{X})(\bar{X} - X_i) \right) \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X} - \mu)^2, \end{aligned}$$

où  $\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mathbb{E}[(X_i - \mu)^2] = \sigma^2$ , par la loi forte des grands nombres, et  $(\bar{X} - \mu) \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 0$ .

Donc, par le théorème 2.9, on a  $T_2^{(n)}(X^{(n)}) \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \sigma^2$

*Remarque.* On a également :

$$\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} s^2 \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 1 \cdot \sigma^2.$$

*Exemple 2.3.* Soient  $X_1, \dots, X_n$  iid  $\text{Unif}(0, \theta)$ , avec  $\theta \in \Theta \subset \mathbb{R}$ . On veut estimer  $g(\theta) = \theta$ . L'estimateur  $T^{(n)}(X^{(n)}) = \bar{X}$  n'est pas un estimateur convergent car :

$$\bar{X} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mathbb{E}_\theta(X_1) = \frac{\theta}{2} \neq \theta \quad \text{sur } \mathbb{P}_\theta^{(n)}.$$

Par contre, si on prend  $T_2^{(n)}(X^{(n)}) = 2\bar{X}$ , on a :

$$2\bar{X} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} 2\mathbb{E}_\theta(X_1) = 2 \cdot \frac{\theta}{2} = \theta \quad \text{sur } \mathbb{P}_\theta^{(n)}.$$

Si on prend  $T_3^{(n)}(X^{(n)}) = X_{(n)}$ , à savoir l'observation maximale, on a :

$$F_{\theta}^{X^{(n)}}(x) = \begin{cases} 0 & \text{si } x < 0 \\ \frac{x^n}{\theta^n} & \text{si } 0 \leq x \leq \theta \\ 1 & \text{sinon} \end{cases}$$

Posons donc  $\varepsilon > 0$ . On calcule :

$$\begin{aligned} \mathbb{P}_{\theta}^{(n)} \left[ \left| X_{(n)} - \theta \right| > \varepsilon \right] &= \mathbb{P}_{\theta}^{(n)} \left[ X_{(n)} \leq \theta - \varepsilon \right] + \mathbb{P}_{\theta}^{(n)} \left[ X_{(n)} \geq \theta + \varepsilon \right] = \mathbb{P}_{\theta}^{(n)} \left[ X_{(n)} \leq \theta - \varepsilon \right] = F_{\theta}^{X^{(n)}}(\theta - \varepsilon) \\ &= \begin{cases} 0 & \text{si } \varepsilon \geq \theta \\ \left( \frac{\theta - \varepsilon}{\theta} \right)^n & \text{si } 0 < \varepsilon < \theta \end{cases} \longrightarrow 0 \end{aligned}$$

on a alors convergence en probabilité de  $X_{(n)}$  vers  $\theta$ . On en déduit que  $T_3^{(n)}(X^{(n)})$  est un estimateur faiblement convergent.

*Remarque.* L'estimateur  $\frac{n+1}{n}X_{(n)}$  est également faiblement convergent.

*Remarque.* Il n'est pas toujours possible de s'en sortir en invoquant le TCL ou la loi des grands nombres pour déterminer la convergence d'un estimateur. Prenons par exemple  $X_1, \dots, X_n$  iid de densité :

$$f^X(x) = \frac{1}{\pi(1 + (x - \theta)^2)}.$$

On a effectivement  $\mathbb{E}[X_1] = +\infty$ . En réalité :

$$\neg \left( \bar{X}^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \theta \quad \text{sur } \mathcal{P}_{\theta}^{(n)} \right),$$

mais bien :

$$\bar{X}^{(n)} \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} X_1$$

## 2.3 Estimateur exhaustif

**Définition 2.16.** Soit  $T^{(n)}(X^{(n)})$ , une statistique. On la dit *exhaustive* lorsque :

$$\forall B \in \mathcal{B}(\mathbb{R}^n) : \forall t \in T^{(n)}(\mathbb{R}^n) : \mathbb{P}_{\theta}^{(n)} \left[ X^{(n)} \in B \mid T^{(n)}(X^{(n)}) = t \right] \text{ ne dépend pas de } \theta.$$

*Remarque.* Puisque l'on travaille sur des modèles paramétriques  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$ , il existe toujours un  $B \in \mathcal{B}(\mathbb{R}^n)$  tel que  $\mathbb{P}_{\theta}[X^{(n)} \in B]$  dépende de  $\theta$ .

*Remarque.* On peut en comprendre qu'une statistique est exhaustive si la valeur prise par  $T^{(n)}(X^{(n)})$  donne toutes les informations contenues par  $X^{(n)}$  sur  $\theta$ .

*Exemple 2.4.* La statistique identité  $x^{(n)} \mapsto x^{(n)}$  est une statistique exhaustive car :

$$\mathbb{P} \left[ X^{(n)} \in B \mid X^{(n)} = x^{(n)} \right] = \begin{cases} 1 & \text{si } x^{(n)} \in B \\ 0 & \text{sinon} \end{cases}.$$

*Exemple 2.5.* Prenons  $X_1, \dots, X_n$  iid Bern(p) avec la statistique  $T^{(n)}(X^{(n)}) = \sum_{i=1}^n X_i$ . Pour évaluer la probabilité :

$$\mathbb{P}_{\mathbf{p}}^{(n)} \left[ X^{(n)} \in \{x^{(n)}\} \mid \sum_{i=1}^n X_i = t \right],$$

on est en présence d'une binomiale. Dès lors, si  $\sum_{i=1}^n X_i \neq t$ , alors la probabilité est nulle. Sinon, la probabilité est  $\frac{1}{\binom{n}{t}}$  car il y a  $\binom{n}{t}$  moyens d'avoir  $n$  observations dont  $t$  valant 1 et  $n - t$  valant 0. Ces probabilités

ne dépendent donc pas de  $\theta$ , la statistique  $T^{(n)}(X^{(n)}) = \sum_{i=1}^n X_i$  est donc une statistique exhaustive.

*Remarque.* Si  $T^{(n)}(X^{(n)})$  est une statistique bijective, alors elle est exhaustive. Cependant, les estimateurs intéressants sont ceux qui « réduisent » l'information de manière à ce qu'elles soient plus facilement analysables.

**Définition 2.17.** Soit  $X^{(n)} = (X_1, \dots, X_n)$ . On appelle la *fonction de vraisemblance* de  $X^{(n)}$  la fonction :

$$L_{\theta}^{(n)} : \mathbb{R}^n \rightarrow \mathbb{R} : x^{(n)} \mapsto \begin{cases} \mathbb{P}[X^{(n)} = x^{(n)}] & \text{si } X^{(n)} \text{ est de loi discrète} \\ f_{\theta}^{X^{(n)}}(x^{(n)}) & \text{sinon} \end{cases}.$$

*Remarque.* Dans le cas de variables  $X_1, \dots, X_n$  iid, la fonction de vraisemblance correspond toujours à un produit :

$$L_{\theta}^{(n)}(X^{(n)}) = \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] \stackrel{\text{indépendance}}{=} \prod_{i=1}^n \mathbb{P}[X_i = x_i].$$

**Théorème 2.18** (Critère de factorisation de Neymann-Fisher). Dans un modèle paramétrique  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$ , une statistique  $T^{(n)}(X^{(n)})$  est exhaustive **si et seulement si** pour tout  $\theta \in \Theta$ , la fonction de vraisemblance  $L_{\theta}^{(n)}(X^{(n)})$  est factorisable sous la forme :

$$L_{\theta}^{(n)}(X^{(n)}) \left( = g_{\theta} \circ T^{(n)} \right) h(x^{(n)}),$$

et ce  $\mathbb{P}_{\theta}^{(n)}$ -sûrement.

*Remarque.* Dans cette factorisation, la fonction  $h$  ne peut dépendre de  $\theta$ , et seule la fonction  $g_{\theta}$  peut en dépendre, mais uniquement par l'intermédiaire de  $T^{(n)}$ .

*Exemple 2.6.* En reprenant l'exemple d'au-dessus :  $X_1, \dots, X_n$  iid Bern( $p$ ) et  $T^{(n)}(X^{(n)}) = \sum_{i=1}^n X_i$ , on a :

$$L_{\theta}^{(n)}(x^{(n)}) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}.$$

Dès lors, en posant  $g_{\theta}(x) = p^x (1-p)^{n-x}$  et  $h(x^{(n)}) = 1$ , on a bien une factorisation de Neymann-Fisher, ce qui implique que la statistique est exhaustive.

*Remarque.* Pour chaque statistique exhaustive, il en existe une infinité définies à bijection près. En effet, si  $T(X^{(n)})$  est une statistique exhaustive et si  $H$  est une fonction bijective quelconque, alors :

$$L_{\theta}^{(n)}(x^{(n)}) = g_{\theta}(T(x^{(n)})) h(x^{(n)}) = (g_{\theta} \circ H^{-1} \circ H \circ T)(x^{(n)}) h(x^{(n)}).$$

La fonction  $H \circ T$  est donc également une statistique exhaustive.

*Remarque.* Le critère précédent peut également donner une manière de deviner des statistiques exhaustives. Prenons par exemple  $X^{(n)} = (X_1, \dots, X_n)$  iid  $\mathcal{N}(\mu, \sigma^2)$ . Si  $\theta = (\mu, \sigma^2) \in \Theta \subset \mathbb{R}^2$ , on peut écrire :

$$\begin{aligned} L_{\theta}^{(n)}(x^{(n)}) &= \prod_{i=1}^n f_{\theta}^{X_i}(x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{2}{2\sigma^2} \sum_{i=1}^n x_i^2 - \frac{n\mu}{\sigma^2} + \frac{\mu}{\sigma^2} \sum_{i=1}^n x_i\right). \end{aligned}$$

Dès lors, en prenant  $T(X^{(n)}) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$ , on a bien une statistique exhaustive. Alors, de même, on peut dire que  $(\frac{1}{n} \sum_{i=1}^n x_i, \frac{1}{n} \sum_{i=1}^n x_i^2)$  est un estimateur exhaustif (composition avec une bijection).

On peut également dire que  $\left(\frac{1}{n} \sum_{i=1}^n x_i, \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2\right)$  est un estimateur exhaustif (même argument). Or ce dernier vecteur correspond à  $(\bar{X}, s^2)$ .

En prenant cette fois  $X^{(n)} = (X_1, \dots, X_n)$  iid  $\text{Unif}(0, \theta)$ , on peut à nouveau construire des statistiques exhaustives :

$$L_{\theta}^{(n)}(x^{(n)}) = \prod_{i=1}^n f_{\theta}^{X_i}(x_i) = \prod_{i=1}^n \frac{1}{\theta} I_{[0 \leq x_i \leq \theta]} = \frac{1}{\theta^n} I_{[0 \leq x_1, \dots, x_n \leq \theta]}.$$

Cela garantit bien que  $x^{(n)}$  est une statistique exhaustive. Mais de plus, par commutativité du produit :

$$L_{\theta}^{(n)}(x^{(n)}) = \prod_{i=1}^n f_{\theta}^{X_{(i)}}(x_i) = \frac{1}{\theta^n} I_{[0 \leq x_{(1)}, \dots, x_{(n)} \leq \theta]}.$$

On a donc que la statistique d'ordre est une statistique exhaustive. Or, la condition  $0 \leq x_{(1)}, \dots, x_{(n)} \leq \theta$  revient à la condition  $0 \leq x_{(1)}, x_{(n)} \leq \theta$ . La statistique  $(x_{(1)}, x_{(n)})$  est donc également une statistique exhaustive. Pour aller plus loin, décomposant la fonction caractéristique  $I_{[0 \leq x_{(1)}, x_{(n)} \leq \theta]}$  en  $I_{[0 \leq x_{(1)}]} I_{[x_{(n)} \leq \theta]}$ , on peut poser  $h(x^{(n)}) = I_{[0 \leq x_{(1)}]}$ , ce qui amène à une nouvelle statistique exhaustive :  $x_{(n)}$ .

À chaque étape du raisonnement, la statistique exhaustive contient *de moins en moins d'information* générale mais conserve l'information sur  $\theta$  qui est donc en quelque sorte contenue dans  $x_{(n)}$ .

### 2.3.1 Estimateurs non biaisés

En se situant toujours dans un modèle statistique  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$ , on veut estimer  $g(\theta)$ , avec  $g : \Theta \rightarrow \mathbb{R}^m$ .

**Définition 2.19.** Un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *non biaisé* lorsque :

$$\forall \theta \in \Theta : \mathbb{E}_{\theta}[T^{(n)}(X^{(n)})] = g(\theta).$$

**Définition 2.20.** Un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *asymptotiquement non biaisé* lorsque :

$$\forall \theta \in \Theta : \mathbb{E}_{\theta}[T^{(n)}(X^{(n)})] \xrightarrow{n \rightarrow +\infty} g(\theta).$$

**Définition 2.21.** Le *biais* d'un estimateur  $T^{(n)}(X^{(n)})$  est la quantité :

$$b_{\theta}^{(n)} = b^{(n)}(\theta) = \mathbb{E}_{\theta}[T^{(n)}(X^{(n)})] - g(\theta).$$

*Remarque.* On remarque donc qu'un estimateur non biaisé a un biais de 0 et qu'un estimateur asymptotiquement non biaisé a un biais qui tend vers 0 pour  $n$  tendant vers  $+\infty$ .

*Exemple 2.7.*  $X_1, \dots, X_n$  iid  $\text{Unif}(0, \theta)$ , avec  $\theta \in \Theta = \mathbb{R}^+_0 \subset \mathbb{R}$ .

— Si  $T^{(n)}(X^{(n)}) = 2\bar{X}$ , on a :

$$\mathbb{E}_{\theta}^{(n)}[2\bar{X}] = \frac{2}{n} \mathbb{E}_{\theta}^{(n)}\left[\sum_{i=1}^n X_i\right] = \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{\theta}[X_i] = \frac{2}{n} \frac{n\theta}{2} = \theta.$$

— Si  $T^{(n)}(X^{(n)}) = X_{(n)}$ , on a :

$$\mathbb{E}_{\theta}^{(n)}[X_{(n)}] = \int_{\mathbb{R}} x f_{\theta}^{X_{(n)}}(x) dx = \int_0^{\theta} x \frac{n x^{n-1}}{\theta^n} dx = \frac{n}{\theta^n} \int_0^{\theta} x^n dx = \frac{n}{\theta^n} \left[ \frac{x^{n+1}}{n+1} \right]_0^{\theta} = \frac{n}{n+1} \theta.$$

On en déduit que  $T^{(n)}(X^{(n)})$  est biaisé car il *voise en moyenne trop à gauche* et a un biais de  $-\frac{\theta}{n+1}$ . Il est cependant asymptotiquement non biaisé car  $\frac{n}{n+1} \xrightarrow{n \rightarrow +\infty} 1$ .

*Remarque.* Si le biais d'un estimateur est négatif, alors c'est que l'estimateur sous-estime en moyenne, alors que si le biais est positif, c'est que l'estimateur surestime.

*Remarque.* On peut tout de même dire que l'estimateur  $\frac{n+1}{n}X_{(n)}$  est non biaisé car :

$$\mathbb{E}_{\theta}^{(n)} \left[ \frac{n+1}{n} X_{(n)} \right] = \frac{n+1}{n} \mathbb{E}_{\theta}^{(n)} [X_{(n)}] = \frac{n+1}{n} \frac{n}{n+1} \theta = \theta.$$

*Exemple 2.8.* Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$ , avec  $\theta = (\mu, \sigma^2)$ .

- pour tout  $\theta \in \Theta$ , on a :  $\mathbb{E}_{\theta}[\bar{X}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \frac{n\mu}{n} = \mu$ .  $\bar{X}$  est donc un estimateur sans biais de  $\mu$  ;
- $\mathbb{E}_{\theta}[s^2] = \mathbb{E}_{\theta}[X_1^2] - \mathbb{E}_{\theta}[\bar{X}^2]$ . On remarque que pour toute variable aléatoire  $Z$ , on a :

$$\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2 \iff \mathbb{E}[Z^2] = \text{Var}(Z) + \mathbb{E}[Z]^2.$$

On peut donc remplacer dans la formule de l'espérance de  $s^2$ , et on obtient :

$$\mathbb{E}_{\theta}[s^2] = \text{Var}_{\theta}(X_1) + \mathbb{E}(X_1)^2 - \text{Var}_{\theta}(\bar{X}) - \mathbb{E}(\bar{X})^2.$$

Or on sait  $\mathbb{E}_{\theta}(X_1) = \mathbb{E}_{\theta}(\bar{X}) = \mu$ . On a donc :

$$\mathbb{E}_{\theta}[s^2] = \text{Var}_{\theta}(X_1) - \text{Var}_{\theta}(\bar{X}) = \sigma^2 - \frac{1}{n^2} \text{Var}_{\theta} \left( \sum_{i=1}^n X_i \right) = \sigma^2 - \frac{1}{n^2} \sum_{i=1}^n \text{Var}_{\theta}(X_i) = \sigma^2 - \frac{n\sigma^2}{n^2} = \frac{n-1}{n} \sigma^2.$$

On en conclut que l'estimateur  $s^2$  est biaisé pour  $\sigma^2$  et de biais  $\frac{-\sigma^2}{n} \xrightarrow[n \rightarrow +\infty]{} 0$ .

Notons alors  $S^2 := \frac{n}{n-1} s^2$ . On a que  $\mathbb{E}_{\theta}[S^2] = \sigma^2$ , et donc  $S^2$  est un estimateur sans biais.

*Remarque.* Le non biais est une propriété fragile. Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$  où l'on veut estimer  $g(\theta) = \sigma$  (et pas  $\sigma^2$ ).

Que vaut  $\mathbb{E}_{\theta}[S]$  ? On sait  $\text{Var}_{\theta}(S) = \mathbb{E}_{\theta}[S^2] - \mathbb{E}_{\theta}[S]^2 = \sigma^2 - \mathbb{E}_{\theta}[S]^2$ . Or  $\text{Var}_{\theta}(S) \geq 0$ , et donc  $\sigma^2 \geq \mathbb{E}_{\theta}[S]^2$ . On en déduit  $\mathbb{E}_{\theta}[S] \leq \sigma$ .

*Remarque.* De plus, il n'existe pas toujours d'estimateur sans biais. Soit  $X_1 \sim \text{Bern}(p)$ . On veut estimer  $g(p) = p^2 \in [0, 1]$ . L'estimateur  $T(X^{(n)})$  est entièrement déterminé par  $T(0)$  et  $T(1)$ . Imposons donc pour tout  $p \in [0, 1] : \mathbb{E}_{\theta}[T(X_1)] = p^2$ . On a donc :

$$p^2 = \mathbb{E}_{\theta}[T(X_1)] = T(0)\mathbb{P}[X_1 = 0] + T(1)\mathbb{P}[X_1 = 1] = T(0)(1-p) + T(1)p.$$

En réarrangeant cette équation du second degré, on obtient :

$$\forall p \in [0, 1] : p^2 + (T(0) - T(1))p - T(0) = 0.$$

Or une telle équation ne peut avoir que 2 racines tout au plus, et ici, une infinité non-dénombrable est requise. Il n'existe donc pas de telle fonction  $T$ , et donc par extension, il n'existe pas d'estimateur sans biais de  $p^2$ .

*Exemple 2.9.* Prenons  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, 1)$ ,  $g(\theta) = g(\mu) = \mu$ . Prenons  $T^{(n)}(X^{(n)}) = \overline{X^{(n)}} + Y$ , où  $Y = \pm 10^9$ , avec probabilité  $\frac{1}{2}$  pour chaque. On a alors :

$$\mathbb{E}_{\mu}[T^{(n)}(X^{(n)})] = \mathbb{E}_{\mu}[\overline{X^{(n)}}] + \mathbb{E}_{\mu}[Y] = \mu + 0 = \mu.$$

L'estimateur  $T^{(n)}(X^{(n)})$  est donc sans biais, mais est très mauvais : sur-/sous-estime toujours à  $\simeq 10^9$  près.

## 2.3.2 Estimateurs à dispersion minimale

**Définition 2.22.** L'erreur quadratique moyenne d'un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  ( $\in \mathbb{R}$ ) est la quantité :

$$\text{MSE}_{\theta}[T^{(n)}(X^{(n)})] := \mathbb{E}_{\theta} \left[ \left| T^{(n)}(X^{(n)}) - g(\theta) \right|^2 \right].$$



**Définition 2.23.** On appelle l'erreur absolue moyenne la quantité :

$$\text{MAE}_\theta[T^{(n)}(X^{(n)})] = \mathbb{E}_\theta \left[ \left| T^{(n)}(X^{(n)}) - g(\theta) \right| \right]$$

*Remarque.*

$$\begin{aligned} \text{MSE}_\theta[T^{(n)}(X^{(n)})] &= \mathbb{E}_\theta \left[ \left| T^{(n)}(X^{(n)}) - g(\theta) \right|^2 \right] = \text{Var}_\theta \left[ T^{(n)}(X^{(n)}) - g(\theta) \right] - \mathbb{E}_\theta \left[ T^{(n)}(X^{(n)}) - g(\theta) \right]^2 \\ &= \text{Var}_\theta[T^{(n)}(X^{(n)})] - b_\theta^{(n)}(T^{(n)}(X^{(n)}))^2. \end{aligned}$$

Puisque  $\text{Var}_\theta(T^{(n)}(X^{(n)})) \geq 0$  et  $b_\theta^{(n)}(T^{(n)}(X^{(n)})) \geq 0$ , pour avoir  $\text{MSE}[T^{(n)}(X^{(n)})] \xrightarrow[n \rightarrow +\infty]{} 0$ , il faut :

$$\begin{cases} \text{Var}_\theta(T^{(n)}(X^{(n)})) \xrightarrow[n \rightarrow +\infty]{} 0, \\ b_\theta^{(n)}(T^{(n)}(X^{(n)})) \xrightarrow[n \rightarrow +\infty]{} 0. \end{cases}$$

Cela dit que sous  $\mathbb{P}_\theta^{(n)}$ , on a  $T^{(n)}(X^{(n)}) \xrightarrow{L^2} g(\theta)$  (et également  $T^{(n)}(X^{(n)}) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} g(\theta)$ ).

On peut également dire que si  $\text{MSE}_\theta[T^{(n)}(X^{(n)})] \rightarrow 0$ , alors  $T^{(n)}(X^{(n)})$  est un estimateur faiblement convergent.

*Remarque.* Le MSE sert d'« arbitrage » entre la variance et le biais. Cet arbitrage est bien souvent nécessaire car il arrive fréquemment que baisser le biais (respectivement la variance) augmente la variance (respectivement le biais).

*Exemple 2.10.*  $X_1, \dots, X_n$  iid  $\text{Unif}(0, \theta)$ . Prenons  $\theta \in \Theta = \mathbb{R}_0^+ \subset \mathbb{R}$ . Prenons  $T_1^{(n)}(X^{(n)}) = X_{(n)}$ . On sait que :

$$\mathbb{E}_\theta[T_1^{(n)}(X^{(n)})] = \frac{n}{n+1}\theta,$$

ce qui nous permet de calculer la variance :

$$\text{Var}_\theta(X_{(n)}) = \mathbb{E}_\theta(X_{(n)}^2) - \mathbb{E}_\theta(X_{(n)})^2,$$

où :

$$\mathbb{E}_\theta(X_{(n)}^2) = \int_{\mathbb{R}} x^2 f^{X_{(n)}}(x) dx = \int_0^\theta x^2 \frac{nx^{n-1}}{\theta^n} dx = \frac{n}{\theta^n} \int_0^\theta x^{n+1} dx = \frac{n}{\theta^n} \left[ \frac{x^{n+2}}{n+2} \right]_0^\theta = \frac{n\theta^2}{n+2}.$$

On trouve finalement la variance donnée par :

$$\frac{n\theta^2}{n+2} - \frac{n^2\theta^2}{(n+1)^2} = n\theta^2 \left( \frac{(n+1)^2 - n(n+2)}{(n+2)(n+1)^2} \right) = \frac{n\theta^2}{(n+2)(n+1)^2}.$$

On trouve finalement une erreur quadratique moyenne de :

$$\text{MSE}_\theta[T_1^{(n)}(X^{(n)})] = \frac{2\theta^2}{(n+1)(n+2)} \xrightarrow[n \rightarrow +\infty]{} 0.$$

En prenant  $T_2^{(n)}(X^{(n)}) = \frac{n+1}{n}X_{(n)}$ , on annule le biais, mais on augmente la variance :

$$\text{Var}_\theta[T_2^{(n)}(X^{(n)})] = \left( \frac{n+1}{n} \right)^2 \text{Var}_\theta[X_{(n)}] = \frac{\theta^2}{n(n+2)}.$$

On trouve alors :

$$\forall \theta \in \Theta : \text{MSE}_\theta[T_2^{(n)}(X^{(n)})] = \frac{\theta^2}{n(n+2)} \leq \text{MSE}_\theta[T_1^{(n)}(X^{(n)})].$$

Dans le cas présent, annuler le biais donne une erreur quadratique moyenne plus faible. Donc ici, la variance a « moins augmenté que le biais n'a diminué ».

*Exemple 2.11.* Prenons maintenant  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma^2)$  et  $g(\theta) = \sigma^2$ .

Pour  $T_1^{(n)}(X^{(n)}) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ , on trouve :

$$\text{MSE}_\theta[T_1^{(n)}(X^{(n)})] = \frac{2(n-1)\sigma^4}{n^2}.$$

À nouveau, il existe un estimateur sans-biais :

$$T_2^{(n)}(X^{(n)}) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n-1}{n} T_1^{(n)}(X^{(n)}).$$

Dans ce cas, on trouve :

$$\forall \theta \in \Theta : \text{MSE}_\theta[T_2^{(n)}(X^{(n)})] = \frac{2\sigma^4}{n-1} \geq \text{MSE}_\theta[T_1^{(n)}(X^{(n)})].$$

*Remarque.* Dans les deux exemples ci-dessus, l'inégalité est stricte pour  $n \geq 2$ .

De plus, l'inégalité tient pour tout  $\theta \in \Theta$ . Cela permet de définir qu'un estimateur est meilleur que l'autre. Avoir deux estimateurs tels que l'erreur quadratique moyenne de l'un est meilleure que l'autre pour certaines valeurs de  $g(\theta)$ , et inversement pour d'autres ne permet pas de dire que l'un est meilleur que l'autre car la véritable valeur de  $g(\theta)$  n'est pas connue. On ne peut donc pas savoir quel estimateur choisir.

**Définition 2.24.** Soit  $\mathcal{C}$ , une classe d'estimateurs de  $g(\theta)$ . On dit que  $T_*^{(n)}$  est à erreur quadratique moyenne minimale dans  $\mathcal{C}$  lorsque :

- $T_*^{(n)} \in \mathcal{C}$  ;
- $\forall T^{(n)} \in \mathcal{C} : \forall \theta \in \Theta : \text{MSE}_\theta[T_*^{(n)}] \leq \text{MSE}_\theta[T^{(n)}]$ .

*Remarque.* On peut prendre  $\mathcal{C} = \{T^{(n)} \text{ t.q. } \mathbb{E}_\theta[T^{(n)2}] < +\infty\}$ , mais on n'en prend qu'un sous-ensemble strict bien souvent car dans un tel  $\mathcal{C}$ , on peut prendre  $T_{\theta_0}^{(n)}$ , pour un certain  $\theta_0 \in \Theta$  défini par  $T_{\theta_0}^{(n)}(X^{(n)}) = g(\theta_0)$ . On trouve donc :

$$\text{MSE}_\theta[T_{\theta_0}^{(n)}(X^{(n)})] = \mathbb{E}_\theta[(g(\theta_0) - g(\theta))^2] = (g(\theta_0) - g(\theta))^2.$$

La fonction d'erreur quadratique moyenne est donc une parabole qui s'annule en  $\theta_0 = \theta$ . Pour avoir un tel  $T_*^{(n)}$ , il faut que  $\text{MSE}_{\theta_0}[T_*^{(n)}] \leq \text{MSE}_{\theta_0}[T_{\theta_0}^{(n)}] = 0$ . Et ce, pour tout  $\theta_0 \in \Theta$ . Il faut donc avoir :

$$\forall \theta \in \Theta : T_*^{(n)} = g(\theta) \quad \text{sous } \mathbb{P}_\theta^{(n)}.$$

Pour résoudre ce problème, il faut donc réduire  $\mathcal{C}$  afin de ne plus avoir de tels estimateurs *pathologiques*.

**Définition 2.25.** Soit  $Z = (Z_1, \dots, Z_m)^T$ , un vecteur aléatoire (vecteur de variables aléatoires). On définit :

$$\mathbb{E}_\theta(Z) := (\mathbb{E}_\theta[Z_1], \dots, \mathbb{E}_\theta[Z_m])^T.$$

On définit également :

$$\text{Var}_\theta[Z] = \mathbb{E} \left[ (Z - \mathbb{E}_\theta[Z])(Z - \mathbb{E}_\theta[Z])^T \right] \in \mathbb{R}^{m \times m}.$$

*Remarque.* On observe que :

$$\text{Var}_\theta[Z]_{ij} = \mathbb{E}_\theta \left[ (Z - \mathbb{E}_\theta[Z])_i (Z - \mathbb{E}_\theta[Z])_j^T \right] = \mathbb{E}_\theta \left[ (Z_i - \mathbb{E}_\theta[Z_i]) (Z_j - \mathbb{E}_\theta[Z_j]) \right] = \text{Cov}[Z_i, Z_j].$$

Cela veut dire :

$$\text{Var}_\theta[Z] = \begin{pmatrix} \text{Var}_\theta[Z_1] & \text{Cov}[Z_1, Z_2] & \dots & \text{Cov}[Z_1, Z_n] \\ \text{Cov}[Z_2, Z_1] & \text{Var}_\theta[Z_2] & \ddots & \text{Cov}[Z_2, Z_n] \\ \vdots & & \ddots & \vdots \\ \text{Cov}[Z_n, Z_1] & \dots & \dots & \text{Var}_\theta[Z_n] \end{pmatrix}$$

On appelle cette matrice la *matrice de variance/covariance* (pour des raisons évidentes). On voit donc que cette matrice est symétrique car  $\text{Cov}[Z_i, Z_j] = \text{Cov}[Z_j, Z_i]$ .

**Définition 2.26.** Soit  $A$  une matrice. On dit que  $A \geq 0$  si la forme bilinéaire associée est semi-définie positive.

**Proposition 2.27.** Soit  $Z$  un vecteur aléatoire, et soient  $A \in \mathbb{R}^{k \times m}$ ,  $b \in \mathbb{R}^k$ . Alors :

- $\mathbb{E}_\theta[AZ + b] = A\mathbb{E}_\theta[Z] + b$  ;
- $\text{Var}_\theta[AZ + b] = A \text{Var}_\theta[AZ + b] A^T$  ;
- $\text{Var}_\theta[Z] \geq 0$ .

*Démonstration.*

- On observe que  $\mathbb{E}_\theta[AZ + b] = (\mathbb{E}_\theta[(AZ)_i + b_i])_i$ , où :

$$\begin{aligned} \mathbb{E}_\theta[(AZ)_i + b_i] &= \mathbb{E}_\theta[(AZ)_i] + b_i = \mathbb{E}_\theta \left[ \sum_{j=1}^k A_{ij} Z_j \right] + b_i = \sum_{j=1}^k \mathbb{E}_\theta[A_{ij} Z_j] + b_i \\ &= \sum_{j=1}^k A_{ij} \mathbb{E}_\theta[Z_j] + b_i = \sum_{j=1}^k A_{ij} (\mathbb{E}_\theta[Z])_j + b_i = (A\mathbb{E}_\theta[Z] + b)_i. \end{aligned}$$

- On remarque premièrement que  $\text{Var}_\theta(AZ+b) = \mathbb{E}_\theta \left[ (AZ + b - A\mathbb{E}_\theta(Z) - b)(AZ + b - A\mathbb{E}_\theta(Z) - b)^T \right] = \mathbb{E}_\theta \left[ (AZ - A\mathbb{E}_\theta(Z))(AZ - A\mathbb{E}_\theta(Z))^T \right] = \text{Var}_\theta[AZ]$ . De là, on trouve :

$$\begin{aligned} \text{Var}_\theta[AZ]_{ij} &= \mathbb{E}_\theta \left[ (AZ - A\mathbb{E}_\theta(Z))_i (AZ - A\mathbb{E}_\theta(Z))_j \right] = \mathbb{E}_\theta \left[ ((AZ)_i - \mathbb{E}_\theta((AZ)_i))((AZ)_j - \mathbb{E}_\theta((AZ)_j)) \right] \\ &= \mathbb{E}_\theta[(AZ)_i (AZ)_j] - \mathbb{E}_\theta[(AZ)_i] \mathbb{E}_\theta[(AZ)_j] \\ &= \mathbb{E}_\theta \left[ \left( \sum_{\delta=1}^m A_{i\delta} Z_\delta \right) \left( \sum_{\gamma=1}^m A_{j\gamma} Z_\gamma \right) \right] - \mathbb{E}_\theta \left[ \sum_{\delta=1}^m A_{i\delta} Z_\delta \right] \mathbb{E}_\theta \left[ \sum_{\gamma=1}^m A_{j\gamma} Z_\gamma \right] \\ &= \sum_{\delta=1}^m \sum_{\gamma=1}^m A_{i\delta} A_{j\gamma} \mathbb{E}_\theta[Z_\delta Z_\gamma] - \sum_{\delta=1}^m \sum_{\gamma=1}^m A_{i\delta} A_{j\gamma} \mathbb{E}_\theta[Z_\delta] \mathbb{E}_\theta[Z_\gamma] \\ &= \sum_{\delta=1}^m \sum_{\gamma=1}^m A_{i\delta} A_{j\gamma} (\mathbb{E}_\theta[Z_\delta Z_\gamma] - \mathbb{E}_\theta[Z_\delta] \mathbb{E}_\theta[Z_\gamma]) \\ &= \sum_{\delta=1}^m \sum_{\gamma=1}^m A_{i\delta} A_{j\gamma} \text{Var}_\theta[Z]_{\delta\gamma} \\ &= \sum_{\delta=1}^m \sum_{\gamma=1}^m A_{i\delta} \text{Var}_\theta[Z]_{\delta\gamma} (A^T)_{\gamma j} \\ &= \left( A \text{Var}_\theta[Z] A^T \right)_{ij} \end{aligned}$$

— soit  $v \in \mathbb{R}^m$ . On remarque que  $v^T \text{Var}_\theta[Z]v = \text{Var}_\theta[v^T Z] \geq 0$ .

□

**Définition 2.28.** Soient  $Z = (Z_1, \dots, Z_m)^T$ , et  $Y = (Y_1, \dots, Y_\ell)^T$ . On définit la covariance de  $Y$  et  $Z$  par :

$$\text{Cov}[Y, Z] := \mathbb{E}_\theta [(Y - \mathbb{E}_\theta[Y])(Z - \mathbb{E}_\theta[Z])].$$

**Proposition 2.29.** Soient  $Z = (Z_1, \dots, Z_m)^T$ , et  $Y = (Y_1, \dots, Y_\ell)^T$ . Alors :

- $\text{Cov}[Y, Z]_{ij} = \text{Cov}[Y_i, Z_j]$ ;
- $\text{Cov}[Z, Y] = \text{Cov}[Y, Z]^T$ ;
- $\text{Cov}[Z, Z] = \text{Var}_\theta[Z]$ ;
- $\text{Cov}[AY + b, CZ + d] = A \text{Cov}[Y, Z] C^T$

**Définition 2.30.** Soient  $g(\theta) \in \mathbb{R}^m$  et  $T^{(n)}(X^{(n)})$ , un estimateur de  $g(\theta)$ . On définit :

$$\text{MSE}_\theta[T^{(n)}(X^{(n)})] := \mathbb{E}_\theta \left[ (T^{(n)}(X^{(n)}) - g(\theta))(T^{(n)}(X^{(n)}) - g(\theta))^T \right] \in \text{Mat}(m, m).$$

*Remarque.*  $\text{MSE}_\theta[T^{(n)}(X^{(n)})] = \mathbb{E}_\theta \left[ \left( T_i^{(n)}(X^{(n)}) - g(\theta) \right) \left( T_j^{(n)}(X^{(n)}) - g(\theta) \right)^T \right]$ . En particulier :

$$\text{MSE}_\theta[T^{(n)}(X^{(n)})]_{ii} = \text{MSE}_\theta[T_i^{(n)}].$$

**Définition 2.31.** Soient  $A, B \in \text{Mat}(m, m)$ . On dit que  $A \geq B$  lorsque  $A - B \geq 0$ , c-à-d lorsque la forme bilinéaire définie par  $(A - B)$  est semi définie positive.

**Définition 2.32.** Soit  $\mathcal{C}$ , une classe d'estimateurs de  $g(\theta)$ . Alors  $T_*^{(n)}$  est à erreur quadratique moyenne minimale dans  $\mathcal{C}$  lorsque :

- (i)  $T_*^{(n)} \in \mathcal{C}$ ;
- (ii)  $\forall T^{(n)} \in \mathcal{C} : \forall \theta \in \Theta : \text{MSE}_\theta[T^{(n)}] \geq \text{MSE}_\theta[T_*^{(n)}]$ .

**Proposition 2.33.** Pour tout  $\theta \in \Theta$ ,  $T^{(n)} \in \mathcal{C}$ , on a :

$$\text{MSE}_\theta[T^{(n)}] \geq 0.$$

*Démonstration.* Soit  $v \in \mathbb{R}^m$ . On trouve :

$$\begin{aligned} v^T \text{MSE}_\theta[T^{(n)}]v &= v^T \mathbb{E}_\theta \left[ (T^{(n)} - g(\theta)) (T^{(n)} - g(\theta))^T \right] v = \mathbb{E}_\theta \left[ v^T (T^{(n)} - g(\theta)) (T^{(n)} - g(\theta))^T v \right] \\ &= \mathbb{E}_\theta \left[ \left( v^T (T^{(n)} - g(\theta)) \right) \left( v^T (T^{(n)} - g(\theta)) \right)^T \right] \\ &= \mathbb{E}_\theta \left[ \left( v^T (T^{(n)} - g(\theta)) \right)^2 \right] \geq 0. \end{aligned}$$

□

**Proposition 2.34.** Soient  $T_1^{(n)}, T_2^{(n)}$ , deux estimateurs de  $g(\theta)$  tels que :

$$\text{MSE}_\theta[T_1^{(n)}] \geq \text{MSE}_\theta[T_2^{(n)}].$$

Alors :

$$\forall v \in \mathbb{R}^m : \text{MSE}_\theta[v^T T_1^{(n)}] \geq \text{MSE}_\theta[v^T T_2^{(n)}].$$

*Démonstration.*

$$\begin{aligned}
\text{MSE}_\theta[\mathbf{v}^T \mathbf{T}_1^{(n)}] - \text{MSE}_\theta[\mathbf{v}^T \mathbf{T}_2^{(n)}] &= \mathbb{E}_\theta \left[ \left( \mathbf{v}^T \mathbf{T}_1^{(n)} - \mathbf{v}^T \mathbf{g}(\theta) \right)^2 \right] - \mathbb{E}_\theta \left[ \left( \mathbf{v}^T \mathbf{T}_2^{(n)} - \mathbf{v}^T \mathbf{g}(\theta) \right)^2 \right] \\
&= \mathbb{E}_\theta \left[ \mathbf{v}^T \left( \mathbf{T}_1^{(n)} - \mathbf{g}(\theta) \right) \right] - \mathbb{E}_\theta \left[ \mathbf{v}^T \left( \mathbf{T}_2^{(n)} - \mathbf{g}(\theta) \right) \right] \\
&= \mathbf{v}^T \text{MSE}_\theta[\mathbf{T}_1^{(n)}] \mathbf{v} - \mathbf{v}^T \text{MSE}_\theta[\mathbf{T}_2^{(n)}] \mathbf{v} = \mathbf{v}^T \left( \text{MSE}_\theta[\mathbf{T}_1^{(n)}] - \text{MSE}_\theta[\mathbf{T}_2^{(n)}] \right) \mathbf{v} \geq 0,
\end{aligned}$$

par définition de semi-positivité.  $\square$

### 2.3.3 Estimateurs efficaces

**Définition 2.35.** On note l'ensemble des valeurs  $\mathbf{x}^{(n)}$  de  $\mathbb{R}^n$  prises par  $\mathbf{X}^{(n)}$  avec probabilité non-nulle par l'ensemble :

$$\mathcal{X}^{(n)} := \{\mathbf{x}^{(n)} \in \mathbb{R}^n \text{ t.q. } L_\theta^{(n)}(\mathbf{x}^{(n)}) > 0\}.$$

**Définition 2.36.** Soit un modèle statistique. On définit la *matrice d'information de Fisher* par :

$$I^{(n)}(\theta) := \int_{\mathcal{X}^{(n)}} \left[ \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right] \left[ \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right]^T L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)}.$$

**Définition 2.37.** Soit  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$ , un modèle statistique. On le dit *régulier* lorsque :

- H1  $\Theta$  est ouvert ;
- H2  $\mathcal{X}^{(n)}$  ne dépend pas de  $\theta$  ;
- H3  $\forall \mathbf{x}^{(n)} \in \mathcal{X}^{(n)} : \theta \mapsto L_\theta^{(n)}(\mathbf{x}^{(n)})$  est différentiable sur  $\theta$  ;
- H4 l'expression  $\int_{\mathcal{X}^{(n)}} L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)}$  est différentiable sous le signe ;
- H5  $\forall \theta \in \Theta : I^{(n)}(\theta)$  existe, est finie, et est inversible.

*Remarque.* Les intégrales doivent être comprises comme des sommes (voire des séries dans le cas infini) dans les cas discrets.

*Remarque.* Un artifice fréquemment utilisé ici est le suivant :

$$\frac{\partial}{\partial \theta_i} L_\theta^{(n)}(\mathbf{x}^{(n)}) = \frac{\partial}{\partial \theta_i} L_\theta^{(n)}(\mathbf{x}^{(n)}) \frac{L_\theta^{(n)}(\mathbf{x}^{(n)})}{L_\theta^{(n)}(\mathbf{x}^{(n)})} = \frac{\partial}{\partial \theta_i} \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) L_\theta^{(n)}(\mathbf{x}^{(n)}).$$

**Proposition 2.38.** Dans un modèle statistique régulier, on a :

$$\mathbb{E}_\theta \left[ \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right] = 0.$$

*Démonstration.* Par définition de l'espérance :

$$\mathbb{E}_\theta \left[ \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right] = \int_{\mathcal{X}^{(n)}} \left( \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right) L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)} = \int_{\mathcal{X}^{(n)}} \nabla_\theta L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)}.$$

Par hypothèse de régularité, on sait :

$$\int_{\mathcal{X}^{(n)}} \nabla_\theta L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)} = \nabla_\theta \int_{\mathcal{X}^{(n)}} L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)}.$$

On trouve finalement :

$$\mathbb{E}_\theta \left[ \nabla_\theta \ln L_\theta^{(n)}(\mathbf{x}^{(n)}) \right] = \nabla_\theta \int_{\mathcal{X}^{(n)}} L_\theta^{(n)}(\mathbf{x}^{(n)}) d\mathbf{x}^{(n)} = \nabla_\theta 1 = 0.$$

$\square$

*Remarque.* Par la proposition précédente, on peut dire :

$$I^{(n)}(\theta) = \mathbb{E}_\theta \left[ \left[ \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right] \left[ \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right]^T \right] = \text{Var}_\theta \left[ \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right].$$

**Proposition 2.39.** Dans un modèle d'échantillonnage, on a :

$$I^{(n)}(\theta) = nI^{(1)}(\theta).$$

*Démonstration.* Les variables  $X_1, \dots, X_n$  étant iid (modèle d'échantillonnage), on peut exprimer :

$$\begin{aligned} I^{(n)}(\theta) &= \text{Var}_\theta \left[ \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right] = \text{Var}_\theta \left[ \nabla_\theta \ln \left( \prod_{i=1}^n L_\theta^{(1)}(X_i) \right) \right] \\ &= \text{Var}_\theta \left[ \sum_{i=1}^n \nabla_\theta \ln L_\theta^{(1)}(X_i) \right] = \sum_{i=1}^n \text{Var}_\theta \left[ \nabla_\theta \ln L_\theta^{(1)}(X_i) \right] \\ &= n \text{Var}_\theta \left[ \nabla_\theta \ln L_\theta^{(1)}(X_1) \right] = nI^{(1)}(\theta). \end{aligned}$$

□

*Remarque.* On en déduit que  $n$  variables aléatoires iid contiennent  $n$  fois plus d'information qu'une seule d'entre elles.

*Remarque.* pour comprendre l'importance de l'hypothèse H5, supposons  $k = 1$  et  $\forall \theta \in \Theta : I^{(n)}(\theta) = 0$ . 0 n'est pas inversible, ce qui viole H5. Alors par l'égalité d'une remarque précédente :

$$\text{Var}_\theta \left[ \frac{d}{d\theta} \ln L_\theta^{(n)}(X^{(n)}) \right] = 0.$$

On en déduit que  $\frac{d}{d\theta} \ln L_\theta^{(n)}(X^{(n)}) = 0$ , et donc que  $\ln L_\theta^{(n)}(X^{(n)})$  est constante sur son espérance, c-à-d 0. On en déduit :

$$\frac{d}{d\theta} L_\theta^{(n)}(X^{(n)}) = 0 \quad \forall \theta \in \Theta \quad \mathbb{P}_\theta - \text{p.s.}.$$

Et donc, pour tout  $\theta \in \Theta$ , on a :

$$\mathbb{P}[X^{(n)} \in B] = \int_B L_\theta^{(n)} dx^{(n)},$$

et en dérivant les deux membres de l'égalité :

$$\frac{d}{d\theta} \mathbb{P}[X^{(n)} \in B] = \int_B \frac{d}{d\theta} L_\theta^{(n)}(x^{(n)}) dx^{(n)}.$$

Cela veut dire que la probabilité  $\mathbb{P}[X^{(n)} \in B]$  est constante pour toutes valeurs de  $\theta$ . Les  $\mathbb{P}_\theta$  ne sont donc plus différenciables. On a donc perdu toute information sur  $\theta$ .

**Définition 2.40.** Dans un modèle statistique régulier, un estimateur  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *régulier* lorsque :

- (i)  $\forall \theta \in \Theta : \mathbb{E}_\theta \left[ \left\| T^{(n)}(X^{(n)}) \right\|^2 \right] < +\infty$  ;
- (ii)  $\psi_\theta := \mathbb{E}_\theta [T^{(n)}(X^{(n)})] = \int_{\mathcal{X}^{(n)}} T^{(n)}(x^{(n)}) L_\theta^{(n)}(x^{(n)}) dx^{(n)}$  est différentiable sous le signe pour tout  $\theta \in \Theta$ .

**Théorème 2.41.** Soient un modèle statistique régulier, et  $T^{(n)}(X^{(n)})$  un estimateur régulier de  $g(\theta)$ . Alors :

$$\forall \theta \in \Theta : \text{MSE}_\theta(T^{(n)}(X^{(n)})) = \mathbb{E}_\theta \left[ \left[ T^{(n)}(X^{(n)}) - g(\theta) \right] \left[ T^{(n)}(X^{(n)}) - g(\theta) \right]^T \right] \geq \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T,$$

où  $\Delta_\theta = \text{Jac } \psi = \left[ \frac{\partial \psi_i}{\partial \theta_j}(\theta) \right] \in \text{Mat}(m, k)$ .

*Remarque.* Dans le cas où  $g(\theta) = \theta = \psi(\theta)$ , on a un estimateur non-biaisé, et donc :

$$\text{Jac } \psi = \text{Jac Id} = \text{Id}.$$

Et donc, on a :

$$\forall \theta \in \Theta : \text{MSE}_\theta(T^{(n)}(X^{(n)})) \geq I^{(n)}(\theta)^{-1}.$$

*Remarque.* Dans le cas d'un modèle d'échantillonnage, on a  $(I^{(n)})^{-1} = n^{-1} I^{(n)}(\theta)^{-1}$ , et donc :

$$\forall \theta \in \Theta : \text{MSE}_\theta(T^{(n)}(X^{(n)})) \geq n^{-1} \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T \xrightarrow{n \rightarrow +\infty} 0.$$

*Preuve du Théorème 2.41.* Notons d'abord :

$$\begin{aligned} (\Delta_\theta)_{ij} &= \frac{\partial \psi_i}{\partial \theta_j}(\theta) = \frac{\partial}{\partial \theta_j} \mathbb{E}_\theta [T_i^{(n)}(X^{(n)})] = \frac{\partial}{\partial \theta_j} \int_{\mathcal{X}^{(n)}} T_i^{(n)}(x^{(n)}) L_\theta^{(n)} dx^{(n)} \\ &= \int_{\mathcal{X}^{(n)}} T_i(x^{(n)}) \frac{\partial}{\partial \theta_j} L_\theta^{(n)}(x^{(n)}) dx^{(n)} = \int_{\mathcal{X}^{(n)}} T_i(x^{(n)}) \frac{\partial}{\partial \theta_j} L_\theta^{(n)}(x^{(n)}) \ln L_\theta^{(n)}(x^{(n)}) dx^{(n)}. \end{aligned}$$

On a donc une espérance donnée par :

$$(\Delta_\theta)_{ij} = \mathbb{E}_\theta \left[ T_i(X^{(n)}) \frac{\partial}{\partial \theta_j} \ln L_\theta^{(n)}(X^{(n)}) \right] = \mathbb{E}_\theta \left[ T^{(n)}(X^{(n)}) \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right]_{ij}.$$

On a alors l'égalité suivante :

$$\Delta_\theta = \mathbb{E}_\theta \left[ T^{(n)}(X^{(n)}) \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right]_{ij}.$$

On veut ensuite montrer :

$$\begin{aligned} \text{MSE}_\theta[T^{(n)}] &= \mathbb{E}_\theta \left[ \left( T^{(n)} - g(\theta) \right) \left( T^{(n)} - g(\theta) \right)^T \right] \stackrel{(1)}{=} \text{Var}[T^{(n)}] + (\psi(\theta) - g(\theta)) (\psi(\theta) - g(\theta))^T \\ &\stackrel{(2)}{\geq} \text{Var}_\theta[T^{(n)}] \\ &\stackrel{(3)}{\geq} \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T \end{aligned}$$

Montrons d'abord (1). Pour cela :

$$\begin{aligned} \text{MSE}_\theta(T^{(n)}) &= \mathbb{E}_\theta \left[ \left( (T^{(n)} - \psi(\theta)) + (\psi(\theta) - g(\theta)) \right) \left( (T^{(n)} - \psi(\theta)) + (\psi(\theta) - g(\theta)) \right)^T \right] \\ &= \mathbb{E}_\theta [(T^{(n)} - \psi(\theta))(T^{(n)} - \psi(\theta))^T] + \mathbb{E}_\theta [(T^{(n)} - \psi(\theta))(\psi(\theta) - g(\theta))^T] \\ &\quad - (\psi(\theta) - g(\theta)) \mathbb{E}_\theta [(T^{(n)} - \psi(\theta))^T] + (\psi(\theta) - g(\theta))(\psi(\theta) - g(\theta))^T \\ &= \text{Var}_\theta[T^{(n)}] + 0 + 0 + (\psi(\theta) - g(\theta))(\psi(\theta) - g(\theta))^T. \end{aligned}$$

Pour montrer (2), on voit bien que  $(\psi(\theta) - g(\theta))(\psi(\theta) - g(\theta))^T$  est semi-définie positive par construction ( $uu^T$  est toujours semi-définie positive).

Montrons alors (3). Posons pour cela :

$$S_\theta := T^{(n)}(X^{(n)}) - \Delta_\theta I^{(n)}(\theta)^{-1} \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \in \mathbb{R}^m.$$

Calculons ensuite :

$$\mathbb{E}_\theta[S_\theta] = \mathbb{E}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \mathbb{E}_\theta[\nabla_\theta \ln L_\theta^{(n)}(X^{(n)})] = \mathbb{E}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \cdot 0.$$

On sait également que :

$$\begin{aligned} \text{Var}_\theta[S_\theta] &= \mathbb{E}_\theta[S_\theta S_\theta^T] - \mathbb{E}_\theta[S_\theta] \mathbb{E}_\theta[S_\theta]^T \\ &= \mathbb{E}_\theta \left[ \left( T^{(n)} - \Delta_\theta I^{(n)}(\theta)^{-1} \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right) \left( T^{(n)} - \Delta_\theta I^{(n)}(\theta)^{-1} \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right)^T \right] - \psi(\theta) \psi(\theta)^T \\ &= \text{Var}_\theta[T^{(n)}] - \mathbb{E}_\theta \left[ T^{(n)}(X^{(n)}) \left( \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right) \right] \left( I^{(n)}(\theta)^{-1} \right)^T \Delta_\theta^T \\ &\quad - \Delta_\theta I^{(n)}(\theta)^{-1} \mathbb{E}_\theta \left[ \left( \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right) T^{(n)}(X^{(n)}) \right] \\ &\quad + \Delta_\theta I^{(n)}(\theta)^{-1} \mathbb{E}_\theta \left[ \left( \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right) \left( \nabla_\theta \ln L_\theta^{(n)}(X^{(n)}) \right)^T \right] \left( I^{(n)}(\theta)^{-1} \right)^T \Delta_\theta^T \\ &= \text{Var}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta \left( I^{(n)}(\theta)^{-1} \right)^T \Delta_\theta^T - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T + \Delta_\theta I^{(n)}(\theta)^{-1} I^{(n)}(\theta) \left( I^{(n)}(\theta)^{-1} \right)^T \Delta_\theta^T, \end{aligned}$$

par définition de  $I^{(n)}(\theta)$ . Puisque  $I^{(n)}(\theta)$  est symétrique et inversible, on sait que  $I^{(n)}(\theta)^{-1}$  est également symétrique. On trouve alors :

$$\text{Var}_\theta[S_\theta] = \text{Var}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T + \Delta_\theta I^{(n)}(\theta)^{-1} I^{(n)}(\theta) I^{(n)}(\theta)^{-1} \Delta_\theta^T.$$

En simplifiant les produits de matrice avec leur inverse, on obtient :

$$\begin{aligned} \text{Var}_\theta[S_\theta] &= \text{Var}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T + \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T \\ &= \text{Var}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T \end{aligned}$$

Or  $\text{Var}_\theta[S_\theta]$  est semi-définie positive par construction. Dès lors, on sait  $\text{Var}_\theta[T^{(n)}(X^{(n)})] - \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T \geq 0$ , ou encore  $\text{Var}_\theta[T^{(n)}(X^{(n)})] \geq \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T$ .  $\square$

**Définition 2.42.** Dans un modèle statistique régulier, l'estimateur régulier  $T^{(n)}(X^{(n)})$  de  $g(\theta)$  est dit *efficace* lorsque :

$$\forall \theta \in \Theta : \text{MSE}_\theta[T^{(n)}(X^{(n)})] = \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T.$$

*Remarque.* Pour avoir  $T^{(n)}(X^{(n)})$  efficace, il faut  $g(\theta) = \psi(\theta)$  et  $\text{Var}_\theta[T^{(n)}(X^{(n)})] = \Delta_\theta I^{(n)}(\theta)^{-1} \Delta_\theta^T$ . L'estimateur doit donc être sans biais. De plus, la variance est définie par la borne de Cramer.

*Exemple 2.12.* Soient  $X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma_0^2)$ ,  $g(\theta) = \mu$ . On a donc  $k = m = 1$ . On sait que  $T^{(n)}(X^{(n)}) = \bar{X}^{(n)}$  est un estimateur sans biais. Comparons alors la variance de  $T^{(n)}(X^{(n)})$  avec la borne de Cramer.

$$\Delta_\mu = \frac{d}{d\mu} \psi(\mu) = \frac{d}{d\mu} \mathbb{E}_\theta[\bar{X}^{(n)}] = \frac{d}{d\mu} \mu = 1.$$



Calculons alors  $I^{(n)}(\theta)$ , qui est donné par :

$$L_{\mu}^{(n)}(X^{(n)}) = \left( \frac{1}{\sqrt{2\pi}\sigma_0} \right)^n \exp \left( -\frac{1}{2\sigma_0^2} \sum_{i=1}^n (X_i - \mu)^2 \right),$$

et donc :

$$\begin{aligned} I^{(n)}(\theta) &= \text{Var}_{\mu} \left[ \frac{d}{d\mu} \ln L_{\mu}^{(n)}(X^{(n)}) \right] = \text{Var}_{\mu} \left[ \frac{d}{d\mu} \left( n \ln \frac{1}{\sqrt{2\pi}\sigma_0} - \frac{1}{2\sigma_0^2} \sum_{i=1}^n (X_i - \mu)^2 \right) \right] \\ &= \frac{1}{(2\sigma_0^2)^2} \text{Var}_{\mu} \left[ -2 \sum_{i=1}^n (X_i - \mu) \right] = \frac{1}{\sigma_0^4} \sum_{i=1}^n \text{Var}_{\mu}(X_i) = \frac{\sigma_0^2}{\sigma_0^4} = \frac{n}{\sigma_0^2} = \text{Var}[T^{(n)}(X^{(n)})]^{-1} = \text{MSE}_{\theta}[T^{(n)}(X^{(n)})]. \end{aligned}$$

On en déduit que l'estimateur  $\bar{X}^{(n)}$  est efficace pour  $g(\theta) = \mu$ .

*Remarque.* On observe que  $I^{(n)}(\theta) \propto \text{Var}_{\theta}[T^{(n)}(X^{(n)})]^{-1}$ , ce qui est cohérent avec l'intuition de la notion d'information : au plus la variance est faible, au plus les observations sont groupées, et au plus il y a d'information à déduire sur  $\theta$ , alors qu'au plus la variance est élevée, au plus les observations sont dispersées, et au moins il y a d'information contenue sur  $\theta$ .

*Exemple 2.13.*  $X_1, \dots, X_n$  iid Bern( $p$ ), avec  $p \in (0, 1)$ . On sait que  $\mathbb{E}_p[X_i] = p$  et  $\text{Var}_p[X_i] = p(1-p)$ . On pose  $g(p) = p$ . On calcule ensuite :

$$\begin{aligned} I^{(n)}(\theta) &= \text{Var}_p \left[ \frac{d}{dp} \ln L_p^{(n)}(X^{(n)}) \right] = \text{Var}_p \left[ \frac{d}{dp} \ln \left( p^{\sum_{i=1}^n X_i} (1-p)^{n-\sum_{i=1}^n X_i} \right) \right] \\ &= \text{Var}_p \left[ \frac{d}{dp} \left( \sum_{i=1}^n X_i \ln p + \left( n - \sum_{i=1}^n X_i \right) \ln(1-p) \right) \right] \\ &= \text{Var}_p \left[ \frac{1}{p} \sum_{i=1}^n X_i - \frac{1}{1-p} \left( n - \sum_{i=1}^n X_i \right) \right] = \text{Var}_p \left[ \left( \frac{1}{p} + \frac{1}{1-p} \right) \sum_{i=1}^n X_i - \frac{n}{1-p} \right] \\ &= \frac{1}{(p(1-p))^2} \sum_{i=1}^n \text{Var}_p[X_i] = \frac{1}{(p(1-p))^2} n(p(1-p)) = \frac{n}{p(1-p)}. \end{aligned}$$

En prenant  $T^{(n)}(X^{(n)}) = \bar{X}^{(n)}$ , on a  $\mathbb{E}_p[T^{(n)}(X^{(n)})] = p$ , et  $\text{Var}_p[T^{(n)}(X^{(n)})] = \frac{\text{Var}[X_1]}{n} = \frac{p(1-p)}{n} = I^{(n)}(\theta)^{-1}$ .

À nouveau, on en déduit que  $\bar{X}^{(n)}$  est un estimateur efficace.

## 2.4 Méthodes d'estimation

### 2.4.1 Méthode des moments

Plaçons-nous dans un modèle paramétrique  $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^{(n)})$ , avec  $\theta \in \Theta \subseteq \mathbb{R}^k$ . On suppose que pour tout  $1 \leq r \leq k$ , on a :  $\mathbb{E}_{\theta}[|X_i|^r] < +\infty$ .

**Définition 2.43.** Soit  $M : \Theta \rightarrow M(\Theta) : \theta \mapsto [\mu'_i(\theta)]_{1 \leq i \leq k}$ . Supposons  $M$  inversible. Soit  $X^{(n)}$  un vecteur d'observations de loi  $\mathbb{P}_\theta^{(n)}$ . On appelle *estimateur de  $\theta$  par la méthode des moments* la statistique :

$$T^{(n)}(X^{(n)}) = M^{-1}([m'_i]_{1 \leq i \leq k}),$$

où  $m'_i$  est le moment empirique d'ordre  $i$  de  $\theta$

*Exemple 2.14.* Soient  $X_1, \dots, X_n$  iid  $N(\mu, \sigma^2)$ , avec  $\theta \in \Theta \subset \mathbb{R}^2$ . On connaît les moments d'ordre 1 et 2 :

$$\begin{aligned}\mu'_1(\theta) &= \mathbb{E}_\theta[X_1] = \mu \\ \mu'_2(\theta) &= \mathbb{E}_\theta[X_1^2] = \text{Var}_\theta[X_1] + \mathbb{E}_\theta[X_1]^2 = \mu^2 + \sigma^2\end{aligned}$$

On a alors la fonction  $M$  donnée par :

$$M : \Theta \rightarrow M(\Theta) = \Theta : \theta = \begin{pmatrix} \mu \\ \sigma \end{pmatrix} \mapsto \begin{pmatrix} \mu \\ \mu^2 + \sigma^2 \end{pmatrix},$$

qui est bien inversible. Il faut alors résoudre le système :

$$\begin{cases} \mu'_1(\theta) = m'_1 \\ \mu'_2(\theta) = m'_2, \end{cases}$$

ce qui donne :

$$\begin{cases} \mu = \mu'_1(\theta) = m'_1 = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X} \\ \sigma^2 + \mu^2 = \mu'_2(\theta) = m'_2 = \frac{1}{n} \sum_{i=1}^n X_i^2. \end{cases}$$

On prend alors pour estimateur :

$$T^{(n)}(X^{(n)}) = \begin{pmatrix} \bar{X} \\ \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \end{pmatrix} = \begin{pmatrix} \bar{X} \\ s^2 \end{pmatrix}.$$

Cet estimateur est convergent (car convergent en chaque composante), est exhaustif, est biaisé (car  $s^2$  est un estimateur biaisé de  $\sigma^2$ ), et est asymptotiquement efficace.

*Exemple 2.15.* Soient  $X_1, \dots, X_n$  iid  $\text{Bern}(p)$ ,  $p \in [0, 1] \subset \mathbb{R}$ . On obtient :

$$p = \mu = \mu'_1(\theta) = m'_1 = \bar{X},$$

et donc on prend l'estimateur  $T^{(n)}(X^{(n)}) = \bar{X}$ .

*Exemple 2.16.* Soient  $X_1, \dots, X_n$  iid  $\text{Unif}(0, \theta)$ ,  $\theta \in \mathbb{R}^+_0 \subset \mathbb{R}$ . Le système nous donne donc :

$$\frac{\theta}{2} = \mathbb{E}_\theta[X_1] = \mu'_1(\theta) = m'_1 = \bar{X}.$$

L'estimateur de  $\theta$  par la méthode des moments est alors donné par :

$$T^{(n)}(X^{(n)}) = 2\bar{X}.$$

*Remarque.*

$$\text{MSE}_\theta[2\bar{X}] = 0^2 + \text{Var}_\theta[2\bar{X}] = \frac{4}{n} \text{Var}[X_1] \propto \frac{c}{n},$$

alors qu'un estimateur efficace a un  $\text{MSE}_\theta \propto \frac{c}{n^2}$ .

## 2.4.2 Méthode du maximum de vraisemblance

*Remarque.* Une notation usuelle pour noter un estimateur de  $\theta$  est  $\hat{\theta}$ . Cette notation remplace l'estimateur  $T^{(n)}(X^{(n)})$  pour lequel il faut explicitement préciser la variable estimée.  $\hat{\theta}$  est une statistique, et ne dépend aucunement de  $\theta$  !

L'idée d'un tel estimateur est de se baser sur la maximisation (selon  $\theta$ ) de la vraisemblance. On est donc tenté d'écrire :

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L_{\theta}^{(n)}(X^{(n)}).$$

Or il peut arriver qu'un tel  $\theta$  de soit pas unique, et donc qu'il existe plusieurs estimateurs de  $\theta$  par le maximum de vraisemblance.

**Définition 2.44.** L'estimateur  $\hat{\theta}$  est un estimateur de  $\theta$  par la méthode du maximum de vraisemblance lorsque :

$$\forall \theta \in \Theta : L_{\hat{\theta}}^{(n)}(X^{(n)}) \geq L_{\theta}^{(n)}(X^{(n)}).$$

*Exemple 2.17.* Soient  $X_1, \dots, X_n$  iid Unif  $\left(\left[\theta - \frac{1}{2}, \theta + \frac{1}{2}\right]\right)$ , avec  $\theta \in \mathbb{R}$ . On peut écrire la vraisemblance :

$$L_{\theta}^{(n)}(X^{(n)}) = \prod_{i=1}^n f_{\theta}^{X_i}(X_i) = \prod_{i=1}^n I_{[\theta - \frac{1}{2} \leq X_i \leq \theta + \frac{1}{2}]} = I_{[\theta - \frac{1}{2} \leq X_{(1)}]} I_{[X_{(n)} \leq \theta + \frac{1}{2}]} = I_{[X_{(n)} - \frac{1}{2} \leq \theta \leq X_{(1)} + \frac{1}{2}]}.$$

La fonction de vraisemblance valant ici soit 1 soit 0, toute valeur de  $\theta$  telle que  $\theta \in \left[X_{(n)} - \frac{1}{2}, X_{(1)} + \frac{1}{2}\right]$  maximise la vraisemblance.

*Exemple 2.18.* Soient  $X_1, \dots, X_n$  iid Unif  $(0, \theta)$ , avec  $\theta \in \mathbb{R}^+_0$ . Ici, la fonction de vraisemblance ne vaut plus uniquement 0 ou 1 :

$$L_{\theta}^{(n)}(X^{(n)}) = \frac{1}{\theta^n} I_{[0 \leq X_{(1)}]} I_{[X_{(n)} \leq \theta]} = \frac{1}{\theta^n} I_{[0 \leq \theta \leq X_{(n)}]}.$$

Pour  $\theta \in [0, X_{(n)}]$ , la fonction de vraisemblance est nulle, et puis vaut  $\theta^{-n}$  pour  $\theta \geq X_{(n)}$ . Le maximum est donc en  $\theta = X_{(n)}$ .

De plus, la fonction  $\theta \mapsto \theta^{-n}$  est strictement décroissante, donc le maximum de vraisemblance est unique.

*Remarque.* Par le critère de factorisation de Neymann-Fisher, si  $S(X^{(n)})$  est une statistique exhaustive, on sait :

$$L_{\theta}^{(n)}(X^{(n)}) = (g_{\theta} \circ S)(X^{(n)}) \cdot h(X^{(n)}).$$

Donc (en considérant  $\arg \max$  comme l'ensemble des arguments maximisant la valeur), on a :

$$\hat{\theta} \in \arg \max_{\theta \in \Theta} L_{\theta}^{(n)}(X^{(n)}) = \arg \max_{\theta \in \Theta} (g_{\theta} \circ S)(X^{(n)}) = (\alpha \circ S)(X^{(n)}).$$

Cela implique que dans le cas du maximum de vraisemblance, l'estimateur dépend toujours des observations à travers une statistique exhaustive.

**Proposition 2.45.** Si  $\theta \mapsto L_{\theta}^{(n)}(X^{(n)})$  est différentiable sur  $\text{int}(\Theta)$ , alors tout estimateur du maximum de vraisemblance à valeurs dans  $\text{int}(\Theta)$  vérifie :

$$\nabla_{\theta} L_{\theta}^{(n)}(X^{(n)}) \Big|_{\theta=\hat{\theta}} = 0.$$

*Remarque.* Par continuité et croissance stricte de  $\ln(\cdot)$ , on sait également que :

$$\nabla_{\theta} \ln \left( L_{\theta}^{(n)}(X^{(n)}) \right) \Big|_{\theta=\hat{\theta}} = 0.$$

L'intérêt de passer par un logarithme est de transformer les produits en somme, ce qui est plus agréable à dériver.

*Exemple 2.19.* Soient  $X_1, \dots, X_n$  iid  $\text{Bern}(p)$ . On peut calculer la fonction de vraisemblance :

$$L_p^{(n)}(X^{(n)}) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i} = p^{\sum_{i=1}^n X_i} (1-p)^{n-\sum_{i=1}^n X_i} = p^{n\bar{X}} (1-p)^{n(1-\bar{X})}.$$

Annulons ensuite le gradient de  $\ln \circ L_p^{(n)}$  :

$$0 = \frac{\partial}{\partial p} = n\bar{X} \frac{\partial}{\partial p} \ln p + n(1-\bar{X}) \frac{\partial}{\partial p} \ln(1-p) = n\bar{X} \frac{1}{p} + n(\bar{X}-1) \frac{1}{1-p}.$$

En multipliant par  $p(1-p)$  de part et d'autre ,on obtient :

$$0 = n\bar{X}(1-p) + n(\bar{X}-1)p = n\bar{X} - n\bar{X}p + n\bar{X}p - np = n\bar{X} - np,$$

ou encore  $p = \bar{X}$ . Dès lors,  $\hat{p} = \bar{X}$  est un estimateur de  $p$  par le maximum de vraisemblance.

**Proposition 2.46.** Soit  $g : \Theta \rightarrow g(\Theta)$ . Alors :

$$\forall \theta \in \Theta : \widehat{g(\theta)} = g(\hat{\theta}).$$

*Démonstration.* Distinguons les deux cas où  $g$  est bijective, et où  $g$  n'est pas bijective.

Si  $g$  est bijective, alors c'est évident car  $g(\theta)$  transforme juste le paramètre de manière univoque.

Si  $g$  n'est pas bijective, on observe que  $g$  est surjective par construction. Et donc  $g$  n'est pas injective. On pose alors :

$$G : \Theta \rightarrow g(\Theta) \times \Theta : \theta \mapsto \begin{pmatrix} g(\theta) \\ \theta \end{pmatrix}.$$

$G$  est injective par construction, et donc est bijective (car surjective en chaque composante). On sait donc que :

$$\widehat{G(\theta)} = G(\hat{\theta}),$$

ce qui implique l'égalité composante par composante, et donc :

$$\widehat{g(\theta)} = g(\hat{\theta}).$$

□

**Définition 2.47.** Soit  $\hat{\theta}$  un estimateur. On dit qu'il est *asymptotiquement efficace* lorsque :

$$\text{Var}_{\theta}[\hat{\theta}] = \frac{1}{nI^{(1)}(\theta)}.$$

*Exemple 2.20.* Prenons  $X_1, \dots, X_n$  iid  $\mathbb{P}_{\theta}^{(1)}$ , avec  $\theta \in \Theta \subseteq \mathbb{R}$ . Par la loi des grands nombres, on a :

$$\forall \theta \in \Theta : \hat{\theta} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \theta \quad \text{sous } \mathbb{P}_{\theta}^{(1)}.$$

On a également :

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \mathcal{N}\left(0, \frac{1}{I^{(1)}(\theta)}\right).$$

Pour  $n \gg$ , on a :

$$\hat{\theta} \approx \mathcal{N}\left(\theta, \frac{1}{nI^{(1)}(\theta)}\right),$$

et donc :

$$\begin{aligned}\mathbb{E}_{\theta}[\hat{\theta}] &\approx \theta \\ \text{Var}_{\theta}[\hat{\theta}] &\approx \frac{1}{nI^{(1)}(\theta)}.\end{aligned}$$

On sait dès lors que pour  $n$  grand, on a :

$$\text{MSE}_{\theta}[\hat{\theta}] \approx b_{\theta}(\hat{\theta})^2 + \text{Var}_{\theta}[\hat{\theta}] = \text{Var}_{\theta}[\hat{\theta}] = \frac{1}{nI^{(1)}(\theta)},$$

qui est la borne de Cramer-Rao.  $\hat{\theta}$  est donc asymptotiquement efficace.