

Première année Master M.A.E.F. 2007 – 2008

Plan du cours

1. Quelques rappels la théorie de la mesure.
2. Quelques rappels sur les applications de la théorie de la mesure aux probabilités.
3. Estimation paramétrique.
4. Tests paramétriques.
5. Introduction à la statistique non-paramétrique.

Bibliographie

• Livres pour revoir les bases....

1. Baillargeon, B. *Probabilités, statistiques et techniques de régression*. SMG.
2. Bercu, B., Pamphile, P. et Azoulay, E. *Probabilités et Applications - Cours Exercices*. Edisciences.
3. Dress, F. *Probabilités et Statistique*. Dunod.
4. Lecoutre, J.-P. *Statistiques et Probabilités*. Dunod.

Théorie de la mesure et applications aux probabilités

- Ansel et Ducel, *Exercices corrigés en théorie de la mesure et de l'intégration*, Ellipses.
- Barbe, P. et Ledoux, M., *Probabilités*, Belin.
- Dacunha-Castelle, D. et Duflo, M., *Probabilités et Statistiques (I)*, Masson
- Jacod, J., *Cours d'intégration*, <http://www.proba.jussieu.fr/pageperso/jacod.html>.
- Jacod, J., *Cours de Probabilités*, <http://www.proba.jussieu.fr/pageperso/jacod.html>.
- Toulouse, P. *Thèmes de probabilités et statistiques*, Masson.

Statistiques inférentielles

- Dacunha-Castelle, D. et Duflo, M., *Probabilités et Statistiques (I)*, Masson.
- Fourdrinier, D., *Statistique inférentielle*, Dunod.
- Lecoutre, J.-M. et Tassi, P., *Statistique non paramétrique et robustesse*, Economica.
- Milhaud, X., *Statistique*, Belin.
- Monfort, A., *Cours de statistique mathématique*, Economica.
- Saporta, G., *Probabilités, analyse des données et statistiques*. Technip.
- Tsybakov, A. *Introduction à la statistique non-paramétrique*. Collection : Mathématiques et Applications, Springer.

Plan détaillé du cours de STATISTIQUES 1

1 Rappels sur la théorie de la mesure

Introduction

Il demeure des choses inconnues à partir des connaissances antérieures en probabilités :

- Qu'est-ce qu'un événement et l'ensemble de tous les événements ?
- Que se passe-t-il pour des probabilités d'événements moins classiques (par exemple l'ensemble des décimaux) ?
- Comment traiter une variable aléatoire qui est continue et discrète à la fois (par exemple le nombre de minutes passées devant la TV) ?

1.1 Mesures

1.1.1 Tribus

Notation. • Ω est un ensemble (fini ou infini).

- $\mathcal{P}(\Omega)$ est l'ensemble de tous les sous-ensembles (parties) de Ω .

Rappel. Soit E un ensemble. E est dit dénombrable s'il existe une bijection entre E et \mathbb{N} ou un sous-ensemble de \mathbb{N} . Par exemple, un ensemble fini, \mathbb{Z} , \mathbb{D} , $\mathbb{Z} \times \mathbb{Z}$, \mathbb{Q} sont dénombrables. En revanche, \mathbb{R} n'est pas dénombrable.

Définition. Soit une famille \mathcal{F} de parties de Ω (donc $\mathcal{F} \subset \mathcal{P}(\Omega)$). On dit que \mathcal{F} est une algèbre si :

- $\Omega \in \mathcal{F}$;
- lorsque $A \in \mathcal{F}$ alors $(\Omega \setminus A) \in \mathcal{F}$;
- pour tout $n \in \mathbb{N}^*$, lorsque $(A_1, \dots, A_n) \in \mathcal{F}^n$ alors $A_1 \cup \dots \cup A_n \in \mathcal{F}$.

Définition. Soit une famille \mathcal{A} de parties de Ω (donc $\mathcal{A} \subset \mathcal{P}(\Omega)$). On dit que \mathcal{A} est une tribu (ou σ -algèbre) sur Ω si :

- $\Omega \in \mathcal{A}$;
- lorsque $A \in \mathcal{A}$ alors $(\Omega \setminus A) \in \mathcal{A}$;
- pour $I \subset \mathbb{N}$, lorsque $(A_i)_{i \in I} \in \mathcal{A}^I$ alors $\bigcup_{i \in I} A_i \in \mathcal{A}$.

Exemple. • Cas du Pile ou Face.

- Cas où Ω est infini : $\Omega = \mathbb{N}$ par exemple.

Propriété. Avec les notations précédentes :

1. $\emptyset \in \mathcal{A}$;
2. si A et B sont dans la tribu \mathcal{A} , alors $A \cap B$ est dans \mathcal{A} ;
3. si \mathcal{A}_1 et \mathcal{A}_2 sont deux tribus sur Ω , alors $\mathcal{A}_1 \cap \mathcal{A}_2$ est une tribu sur Ω . Plus généralement, pour $I \subset \mathbb{N}$, si $(\mathcal{A}_i)_{i \in I}$ ensemble de tribus sur Ω , alors $\bigcap_{i \in I} \mathcal{A}_i$ est une tribu sur Ω ;
4. si \mathcal{A}_1 et \mathcal{A}_2 sont deux tribus sur Ω , alors $\mathcal{A}_1 \cup \mathcal{A}_2$ n'est pas forcément une tribu sur Ω .

Définition. Si \mathcal{E} est une famille de parties de Ω (donc $\mathcal{E} \subset \mathcal{P}(\Omega)$), alors on appelle tribu engendrée par \mathcal{E} , notée $\sigma(\mathcal{E})$, la tribu engendrée par l'intersection de toutes les tribus contenant \mathcal{E} (on peut faire la même chose avec des algèbres).

Remarque. La tribu engendrée est la "plus petite" tribu (au sens de l'inclusion) contenant la famille \mathcal{E} .

Rappel. • Un ensemble ouvert U dans un espace métrique X est telle que pour tout $x \in U$, il existe $r > 0$ tel que $B(x, r) \subset U$.

- On dit qu'un ensemble dans un espace métrique X est fermé si son complémentaire dans X est ouvert.

Définition. Soit Ω un espace métrique. On appelle tribu borélienne sur Ω , notée, $\mathcal{B}(\Omega)$, la tribu engendrée par les ouverts de Ω . Un ensemble de $\mathcal{B}(\Omega)$ est appelé borélien.

Exemple. • Boréliens sur \mathbb{R} , sur $]0, 1[$.

- Boréliens sur \mathbb{R}^2 .

1.1.2 Espace mesurable

Définition. Soit Ω un ensemble et soit \mathcal{A} une tribu sur Ω . On dit que (Ω, \mathcal{A}) est un espace mesurable.

Corollaire. Quand on s'intéressera aux probabilités, on dira que (Ω, \mathcal{A}) est un espace probabilisable.

Propriété. Si $(\Omega_i, \mathcal{A}_i)_i$ sont n espaces mesurables, alors un ensemble élémentaire de $\Omega = \Omega_1 \times \cdots \times \Omega_n$ est une réunion finie d'ensembles $A_1 \times \cdots \times A_n$ où chaque $A_i \in \mathcal{A}_i$. L'ensemble des ensembles élémentaires est une algèbre et on note $\mathcal{A}_1 \otimes \cdots \otimes \mathcal{A}_n$ (on dit \mathcal{A}_1 tensoriel \mathcal{A}_2 ... tensoriel \mathcal{A}_n) la tribu sur Ω engendrée par ces ensembles élémentaires.

Exemple. Pavés de \mathbb{R}^d .

Définition. On appelle espace mesurable produit des $(\Omega_i, \mathcal{A}_i)_i$ l'espace mesurable $\left(\prod_{i=1}^n \Omega_i, \bigotimes_{i=1}^n \mathcal{A}_i \right)$.

Exemple. Pile / Face 2 fois.

1.1.3 Définitions et Propriétés d'une mesure

Définition. Soit (Ω, \mathcal{A}) un espace mesurable. L'application $\mu : \mathcal{A} \rightarrow [0, +\infty]$ est une mesure si :

- $\mu(\emptyset) = 0$.
- Pour tout $I \subset \mathbb{N}$ et pour $(A_i)_{i \in I}$ famille disjointe de \mathcal{A} (telle que $A_i \cap A_j = \emptyset$ pour $i \neq j$), alors
$$\mu\left(\bigcup_{i \in I} A_i\right) = \sum_{i \in I} \mu(A_i)$$
 (propriété dite de σ -additivité).

Définition. Avec les notations précédentes :

- Si $\mu(\Omega) < +\infty$, on dit que μ est finie.
- Si $\mu(\Omega) < M$ avec $M < +\infty$, on dit que μ est bornée.
- Si $\mu(\Omega) = 1$, on dit que μ est une mesure de probabilité.

Exemple. Cas de $\Omega = \mathbb{R}$, de \mathbb{N} , ou \mathbb{R}^2 .

Définition. Si (Ω, \mathcal{A}) est un espace mesurable (resp. probabilisable) alors $(\Omega, \mathcal{A}, \mu)$ est un espace mesuré (resp. probabilisé quand μ est une probabilité).

Remarque. Sur (Ω, \mathcal{A}) , on peut définir une infinité de mesures.

Propriété. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré et $(A_i)_{i \in \mathbb{N}}$, une famille de \mathcal{A} .

1. Si $A_1 \subset A_2$, alors $\mu(A_1) \leq \mu(A_2)$.

2. Si $\mu(A_1) < +\infty$ et $\mu(A_2) < +\infty$, alors $\mu(A_1 \cup A_2) + \mu(A_1 \cap A_2) = \mu(A_1) + \mu(A_2)$.
3. Pour tout $I \subset \mathbb{N}$, on a $\mu\left(\bigcup_{i \in I} A_i\right) \leq \sum_{i \in I} \mu(A_i)$.
4. Si $A_i \subset A_{i+1}$ pour tout $i \in \mathbb{N}$ (suite croissante en sens de l'inclusion), alors $(\mu(A_n))_{n \in \mathbb{N}}$ est une suite croissante convergente telle que
$$\mu\left(\bigcup_{i \in \mathbb{N}} A_i\right) = \lim_{i \rightarrow +\infty} \mu(A_i)$$
 (même si cette limite est $+\infty$).
5. Si $A_{i+1} \subset A_i$ pour tout $i \in \mathbb{N}$ (suite décroissante en sens de l'inclusion) et $\mu(A_0) < +\infty$, alors $(\mu(A_n))_{n \in \mathbb{N}}$ est une suite décroissante convergente telle que
$$\mu\left(\bigcap_{i \in \mathbb{N}} A_i\right) = \lim_{i \rightarrow +\infty} \mu(A_i)$$
.

Exemple. 1. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré. On définit $\nu(A) = \mu(A \cap B)$ où $B \in \mathcal{A}$. ν mesure ?

2. Si μ_1 et μ_2 mesures sur (Ω, \mathcal{A}) , $\mu_1 + \mu_2$ et $\alpha\mu$ sont-elles des mesures ?

Définition. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré et $(A_i)_{i \in \mathbb{N}}$ une famille de \mathcal{A} .

1. On définit $\limsup(A_n)_n = \bigcap_{n \in \mathbb{N}} \bigcup_{m \geq n} A_m$ (intuitivement, $\limsup(A_n)_n$ est l'ensemble des $\omega \in \Omega$ tels que ω appartienne à une infinité de A_n).
2. On définit $\liminf(A_n)_n = \bigcup_{n \in \mathbb{N}} \bigcap_{m \geq n} A_m$ (intuitivement, $\liminf(A_n)_n$ est l'ensemble des $\omega \in \Omega$ tels que ω appartienne à tous les A_n sauf à un nombre fini d'entre eux).

Exemple. Cas des suites croissantes et décroissantes d'ensembles.

Théorème (Théorème d'extension de Hahn - Caratheodory). Si Ω est un ensemble, \mathcal{F} une algèbre sur Ω , et ν une application de \mathcal{F} dans $[0, +\infty]$ additive (telle que $\nu(A \cup B) = \nu(A) + \nu(B)$ pour $A \cup B = \emptyset$), alors si \mathcal{A} est la tribu engendrée par \mathcal{F} , il existe une mesure $\hat{\nu}$ sur la tribu \mathcal{A} qui coïncide avec ν sur \mathcal{F} (c'est-à-dire que pour tout $F \in \mathcal{F}$, $\hat{\nu}(F) = \nu(F)$). On dit que $\hat{\nu}$ prolonge ν sur la tribu \mathcal{A} .

Exemple. Définition de la mesure de Lebesgue sur \mathbb{R} , \mathbb{R}^n, \dots

Définition. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré.

1. Pour $A \in \mathcal{A}$, on dit que A est μ -négligeable si $\mu(A) = 0$.
2. Soit une propriété \mathcal{P} dépendant des éléments ω de Ω . On dit que \mathcal{P} est vraie μ -presque partout (μ -presque sûrement sur un espace probabilisé) si l'ensemble des ω pour laquelle elle n'est pas vérifiée est μ -négligeable.

Exemple. • Mesure de Lebesgue sur \mathbb{N} ou \mathbb{Q} .

- La propriété " la suite de fonction $f_n(x) = x^n$ converge vers la fonction $f(x) = 0$ " est vraie λ -presque partout sur $[0, 1]$.
- Soit $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mu)$ et soit F la fonction définie par $F(x) = \mu(-\infty, x]$ pour $x \in \mathbb{R}$.

1.1.4 Fonctions mesurables

Rappel. Soit $f : E \mapsto F$, où E et F sont 2 espaces métriques.

- Pour $I \subset F$, on appelle ensemble réciproque de I par f , l'ensemble $f^{-1}(I) = \{x \in E, f(x) \in I\}$.
- (f continue) \iff (pour tout ouvert U de F alors $f^{-1}(U)$ est un ouvert de E).

Définition. Soit $f : E \mapsto F$ et soit \mathcal{I} une tribu sur F . On note $f^{-1}(\mathcal{I})$ l'ensemble de sous-ensembles de Ω tel que $f^{-1}(\mathcal{I}) = \{f^{-1}(I), I \in \mathcal{I}\}$.

Propriété. Soit (Ω', \mathcal{A}') un espace mesurable et soit $f : \Omega \mapsto \Omega'$. Alors $f^{-1}(\mathcal{A})$ est une tribu sur Ω appelée tribu engendrée par f .

Définition. Soit (Ω, \mathcal{A}) et (Ω', \mathcal{A}') deux espaces mesurables. Une fonction $f : \Omega \mapsto \Omega'$ est dite mesurable pour les tribus \mathcal{A} et \mathcal{A}' si et seulement si $f^{-1}(\mathcal{A}') \subset \mathcal{A}$ (donc si et seulement si $\forall A' \in \mathcal{A}'$, alors $f^{-1}(A') \in \mathcal{A}$).

Exemple. • Fonction indicatrice.

- Combinaison linéaire de fonctions indicatrices.

Remarque. Dans le cas où (Ω, \mathcal{A}) est un espace probabilisable, et si $f : \Omega \mapsto \mathbb{R}$, alors si f est une fonction mesurable sur \mathcal{A} et $\mathcal{B}(\mathbb{R})$, alors f est une variable aléatoire.

Exemple. Nombre de Piles dans un jeu de Pile/Face.

Remarque. Dans le cas où (Ω, \mathcal{A}) est un espace mesurable, et si $f : \Omega \mapsto (\Omega', \mathcal{B}(\Omega'))$, où Ω' est un espace métrique et $\mathcal{B}(\Omega')$ l'ensemble des boréliens de Ω' , si f est une fonction mesurable sur \mathcal{A} et $\mathcal{B}(\Omega')$, alors f est dite fonction borélienne.

Proposition. Soit (Ω, \mathcal{A}) et (Ω', \mathcal{A}') deux espaces mesurables et $f : \Omega \mapsto \Omega'$. Soit \mathcal{F} une famille de sous-ensembles de Ω' telle que $\sigma(\mathcal{F}) = \mathcal{A}'$. Alors

1. $f^{-1}(\mathcal{F})$ engendre la tribu $f^{-1}(\mathcal{A}')$.
2. $(f \text{ mesurable}) \iff (f^{-1}(\mathcal{F}) \subset \mathcal{A})$

Conséquence. • Si (Ω, \mathcal{A}) et (Ω', \mathcal{A}') sont deux espaces mesurables boréliens, alors toute application continue de $\Omega \mapsto \Omega'$ est mesurable.

- Pour montrer qu'une fonction $f : \Omega \mapsto \mathbb{R}$ est mesurable, il suffit de montrer que la famille d'ensemble $(\{\omega \in \Omega, f(\omega) \leq a\})_{a \in \mathbb{R}} \in \mathcal{A}$.

Propriété. • Soit f mesurable de (Ω, \mathcal{A}) dans (Ω', \mathcal{A}') et g mesurable de (Ω', \mathcal{A}') dans $(\Omega'', \mathcal{A}'')$. Alors $g \circ f$ est mesurable dans \mathcal{A} et \mathcal{A}'' .

- Soit f_1 mesurable de (Ω, \mathcal{A}) dans $(\Omega_1, \mathcal{A}_1)$ et f_2 mesurable de (Ω, \mathcal{A}) dans $(\Omega_2, \mathcal{A}_2)$. Alors $h : \Omega \mapsto \Omega_1 \times \Omega_2$ telle que $h(\omega) = (f_1(\omega), f_2(\omega))$ est mesurable dans \mathcal{A} et $\mathcal{A}_1 \otimes \mathcal{A}_2$.
- Soit $(f_n)_{n \in \mathbb{N}}$ une suite de fonctions mesurables de (Ω, \mathcal{A}) dans $(\Omega', \mathcal{B}(\Omega'))$, où Ω' est un espace métrique, telle qu'il existe une fonction f limite simple de (f_n) (donc $\forall \omega \in \Omega, \lim_{n \rightarrow \infty} f_n(\omega) = f(\omega)$). Alors f est mesurable dans \mathcal{A} et $\mathcal{B}(\Omega')$.

Définition. Soit f mesurable de $(\Omega, \mathcal{A}, \mu)$ dans (Ω', \mathcal{A}') et soit $\mu_f : \mathcal{A}' \mapsto [0, +\infty]$ telle que pour tout $A' \in \mathcal{A}'$, on ait $\mu_f(A') = \mu(f^{-1}(A'))$. Alors μ_f est une mesure sur (Ω', \mathcal{A}') appelée mesure image de μ par f .

Cas particulier. Si μ est une mesure de probabilité et si X est une variable aléatoire alors μ_X est la mesure (loi) de probabilité de la variable aléatoire X .

1.1.5 Cas des fonctions réelles mesurables

Propriété. Soit f et g deux fonctions réelles mesurables (de $(\Omega, \mathcal{A}, \mu)$ dans $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$). Alors $\alpha.f$, $f + g$, $\min(f, g)$ et $\max(f, g)$ sont des fonctions réelles mesurables.

Propriété. Soit $(f_n)_{n \in \mathbb{N}}$ une suite de fonctions réelles mesurables. Alors $\inf(f_n)$ et $\sup(f_n)$ sont des fonctions réelles mesurables.

Définition. Soit $f : \Omega \rightarrow \mathbb{R}$. Alors f est dite étagée s'il existe une famille d'ensembles disjoints $(A_i)_{1 \leq i \leq n}$ de Ω et une famille de réels $(\alpha_i)_{1 \leq i \leq n}$ telles que pour tout $\omega \in \Omega$, on ait $f(\omega) = \sum_{i=1}^n \alpha_i \mathbb{I}_{A_i}(\omega)$.

Remarque. Si les A_i sont tous dans \mathcal{A} tribu sur Ω , alors f est \mathcal{A} -mesurable.

Théorème. Toute fonction réelle mesurable à valeurs dans $[0, +\infty]$ est limite simple d'une suite croissante de fonctions étagées.

Conséquence. Soit f une fonction réelle mesurable. Alors f est limite simple de fonctions étagées.

1.2 Intégration de Lebesgue

Dans toute la suite, on considère $(\Omega, \mathcal{A}, \mu)$ un espace mesuré.

1.2.1 Intégrale de Lebesgue d'une fonction positive

Définition. 1. Soit $f = \mathbb{I}_A$, où $A \in \mathcal{A}$. Alors :

$$\int f d\mu = \int_{\omega} f(\omega) d\mu(\omega) = \mu(A).$$

2. Soit $f = \mathbb{I}_A$, où $A \in \mathcal{A}$ et soit $B \in \mathcal{A}$. Alors :

$$\int_B f d\mu = \int_B f(\omega) d\mu(\omega) = \int \mathbb{I}_B \mu(A)(\omega) f(\omega) d\mu(\omega) = \mu(A \cap B).$$

3. Soit f une fonction étagée positive telle que $f = \sum_{i=1}^n \alpha_i \mathbb{I}_{A_i}$, où les $A_i \in \mathcal{A}$ et $\alpha_i > 0$ et soit $B \in \mathcal{A}$.

Alors :

$$\int_B f d\mu = \int_B f(\omega) d\mu(\omega) = \int \mathbb{I}_B(\omega) f(\omega) d\mu(\omega) = \sum_{i=1}^n \alpha_i \mu(A_i \cap B).$$

Exemple. Fonction $\mathbb{I}_{\mathbb{Q}}$, fonctions en escalier,...

Définition. Soit f une fonction \mathcal{A} -mesurable positive et soit $B \in \mathcal{A}$. Alors l'intégrale de Lebesgue de f par rapport à μ sur B est :

$$\int_B f d\mu = \int \mathbb{I}_B(\omega) f(\omega) d\mu(\omega) = \sup \left\{ \int_B g d\mu, \text{ pour } g \text{ étagée positive telle que } g \leq f \right\}.$$

Propriété. Soit f une fonction \mathcal{A} -mesurable positive et soit A et $B \in \mathcal{A}$. Alors :

1. Pour $c \geq 0$, $\int_B cf d\mu = c \int_B f d\mu$.

2. Si $A \subset B$, alors $\int_A f d\mu \leq \int_B f d\mu$.

3. Si g est une fonction \mathcal{A} -mesurable positive telle que $0 \leq f \leq g$ alors $0 \leq \int_B f d\mu \leq \int_B g d\mu$.

4. Si $\mu(B) = 0$ alors $\int_B f d\mu = 0$.

Théorème (Théorème de convergence monotone (Beppo-Lévi)). Si $(f_n)_n$ est une suite croissante de fonctions mesurables positives convergeant simplement vers f sur Ω , alors :

$$\lim_{n \rightarrow \infty} \left(\int f_n d\mu \right) = \int f d\mu = \int \lim_{n \rightarrow \infty} f_n d\mu.$$

Conséquence. Pour les séries de fonctions mesurables positives, on peut toujours appliquer le Théorème de convergence monotone et donc inverser la somme et l'intégrale.

Lemme (Lemme de Fatou). Soit $(f_n)_n$ est une suite de fonctions mesurables positives alors :

$$\int \left(\liminf_{n \rightarrow \infty} f_n \right) d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu.$$

Exemple. Appliquer Fatou à (f_n) telle que $f_{2n} = \mathbb{I}_A$ et $f_{2n+1} = \mathbb{I}_B$.

1.2.2 Intégrale de Lebesgue d'une fonction réelle et propriétés

Définition. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré, $B \in \mathcal{A}$ et soit f une fonction \mathcal{A} -mesurable à valeurs réelles telle que $f = f^+ - f^-$ avec $f^+ = \max(f, 0)$ et $f^- = \max(-f, 0)$. On dit que f est μ -intégrable sur B si $\int_B |f| d\mu < +\infty$. On a alors

$$\int_B f d\mu = \int_B f^+ d\mu - \int_B f^- d\mu.$$

Notation. Lorsque f est μ -intégrable sur B , soit $\int |f| d\mu < +\infty$, on note $f \in \mathcal{L}^1(\Omega, \mathcal{A}, \mu)$ (on dit que f est \mathcal{L}^1).

Exemple. Intégrale de Riemann et intégrale de Lebesgue.
Cas de la masse de Dirac.

Propriété. On suppose que f et $g \in \mathcal{L}^1(\Omega, \mathcal{A}, \mu)$. Alors :

1. $\int (\alpha f + \beta g) d\mu = \alpha \int f d\mu + \beta \int g d\mu$ pour $(\alpha, \beta) \in \mathbb{R}^2$.
2. Si $f \leq g$ alors $\int f d\mu \leq \int g d\mu$.

Théorème (Théorème de convergence dominée de Lebesgue). Soit $(f_n)_n$ est une suite de fonctions de $\mathcal{L}^1(\Omega, \mathcal{A}, \mu)$ telles que pour tout $n \in \mathbb{N}$, $|f_n| \leq g$ avec $g \in \mathcal{L}^1(\Omega, \mathcal{A}, \mu)$. Si on suppose que (f_n) converge simplement vers f sur Ω alors :

$$\lim_{n \rightarrow \infty} \int f_n d\mu = \int f d\mu.$$

Extension. Le Théorème de Lebesgue s'applique également dans le cas où $(f_n)_n$ converge presque partout vers f .

Exemple. Convergence d'intégrale dépendant d'un paramètre : par exemple $\int_0^\infty \frac{f(x)}{1+x^n} dx$.

Théorème (Inégalité de Jensen). Soit $(\Omega, \mathcal{A}, \mathbb{P})$ un espace probabilisé, soit $\phi : \mathbb{R} \mapsto \mathbb{R}$ une fonction convexe et soit $f : \Omega \mapsto \mathbb{R}$ mesurable telle que $\phi(f)$ soit une fonction intégrable par rapport à P . Alors :

$$\phi \left(\int f d\mathbb{P} \right) \leq \int \phi(f) d\mathbb{P}.$$

Exemple. Soit X une v.a. sur $(\Omega, \mathcal{A}, \mathbb{P})$. Alors $\phi(\mathbb{E}X) \leq \mathbb{E}(\phi(X))$.

1.2.3 Mesures induites et densités

Théorème (Théorème du Transport). Soit f une fonction mesurable de $(\Omega, \mathcal{A}, \mu)$ dans (Ω', \mathcal{A}') telle que μ_f soit la mesure induite par f (donc $\mu_f(A') = \mu(f^{-1}(A'))$) pour $A' \in \mathcal{A}'$ et soit ϕ une fonction mesurable de (Ω', \mathcal{A}') dans $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Alors, si $\phi \circ f \in \mathcal{L}^1(\Omega, \mathcal{A}, \mu)$,

$$\int_{\Omega'} \phi d\mu_f = \int_{\Omega} \phi \circ f d\mu.$$

Définition. Soit μ et ν deux mesures sur (Ω, \mathcal{A}) . On dit que μ domine ν (ou ν est dominée par μ) et que ν est absolument continue par rapport à μ lorsque pour tout $A \in \mathcal{A}$, $\mu(A) = 0 \implies \nu(A) = 0$.

Propriété. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré et f une fonction définie sur (Ω, \mathcal{A}) mesurable et positive. On suppose que pour $A \in \mathcal{A}$, $\nu(A) = \int_A f d\mu$. Alors, ν est une mesure sur (Ω, \mathcal{A}) , dominée par μ . De plus, pour toute fonction g définie sur (Ω, \mathcal{A}) mesurable et positive,

$$\int g d\nu = \int g \cdot f d\mu.$$

Enfin, g est ν intégrable si et seulement si $g \cdot f$ est μ intégrable.

Définition. On dit que μ mesure sur (Ω, \mathcal{A}) est σ -finie lorsqu'il existe une famille $(A_i)_{i \in I}$, avec I dénombrable, d'ensembles de \mathcal{A} telle que $\bigcup A_i = \Omega$ et $\mu(A_i) < +\infty$ pour tout $i \in I$.

Théorème (Théorème de Radon-Nikodym). On suppose que μ et ν sont deux mesures σ -finies sur (Ω, \mathcal{A}) telles que μ domine ν . Alors il existe une fonction f définie sur (Ω, \mathcal{A}) mesurable et positive, appelée densité de ν par rapport à μ , telle que pour tout $A \in \mathcal{A}$, $\nu(A) = \int_A f d\mu$.

Théorème (Théorème de Fubini). Soit $\Omega = \Omega_1 \times \Omega_2$, $\mathcal{A} = \mathcal{A}_1 \otimes \mathcal{A}_2$ et $\mu = \mu_1 \otimes \mu_2$ (mesures σ finies), où $(\Omega_1, \mathcal{A}_1, \mu_1)$ et $(\Omega_2, \mathcal{A}_2, \mu_2)$ sont des espaces mesurés. Soit une fonction $f : \Omega \mapsto \mathbb{R}$, \mathcal{A} -mesurable et μ -intégrable. alors :

$$\int_{\Omega} f d\mu = \int_{\Omega_1} \left(\int_{\Omega_2} f(\omega_1, \omega_2) d\mu_2(\omega_2) \right) d\mu_1(\omega_1) = \int_{\Omega_2} \left(\int_{\Omega_1} f(\omega_1, \omega_2) d\mu_1(\omega_1) \right) d\mu_2(\omega_2).$$

1.2.4 Espaces \mathcal{L}^p

Définition. Soit $(\Omega, \mathcal{A}, \mu)$ un espace mesuré. On appelle espace $\mathcal{L}^p(\Omega, \mathcal{A}, \mu)$, où $p > 0$, l'ensemble des fonctions $f : \Omega \mapsto \mathbb{R}$, mesurables et telles que $\int |f|^p d\mu < +\infty$.

Définition. Pour $f \in \mathcal{L}^p(\Omega, \mathcal{A}, \mu)$, où $p > 0$, on note $\|f\|_p = \left(\int |f|^p d\mu \right)^{1/p}$.

Propriété (Inégalité de Hölder). Soit $p > 1$ et $q > 1$ tels que $\frac{1}{p} + \frac{1}{q} = 1$, et soit $f \in \mathcal{L}^p(\Omega, \mathcal{A}, \mu)$ et $g \in \mathcal{L}^q(\Omega, \mathcal{A}, \mu)$. Alors, $fg \in \mathcal{L}^1(\Omega, \mathcal{A}, \mu)$ et

$$\|fg\|_1 \leq \|f\|_p \cdot \|g\|_q.$$

Propriété (Inégalité de Minkowski). Soit $p > 1$ et soit f et $g \in \mathcal{L}^p(\Omega, \mathcal{A}, \mu)$. Alors, $f + g \in \mathcal{L}^p(\Omega, \mathcal{A}, \mu)$ et

$$\|f + g\|_p \leq \|f\|_p + \|g\|_p.$$

Remarque. Pour $p > 1$, $\|\cdot\|_p$ définie ainsi sur une semi-norme sur $\mathcal{L}^p(\Omega, \mathcal{A}, \mu)$. Pour obtenir une norme, il faut se placer dans l'espace $\mathbb{L}^p(\Omega, \mathcal{A}, \mu)$ obtenu en "quotientant" $\mathcal{L}^p(\Omega, \mathcal{A}, \mu)$ par la relation d'équivalence $f = g$ μ -presque partout (c'est-à-dire que dans $\mathbb{L}^p(\Omega, \mathcal{A}, \mu)$ on dira que $f = g$ lorsque $f = g$ μ -presque partout).

Définition. Pour f et $g \in \mathbb{L}^2(\Omega, \mathcal{A}, \mu)$, on définit le produit scalaire $\langle f, g \rangle = \int f \cdot g d\mu$. On muni ainsi $\mathbb{L}^2(\Omega, \mathcal{A}, \mu)$ d'une structure d'espace de Hilbert. On dira que f est orthogonale à g lorsque $\langle f, g \rangle = 0$.

Conséquence. Si A est un sous-espace vectoriel fermé de $\mathbb{L}^2(\Omega, \mathcal{A}, \mu)$ (par exemple un sous-espace de dimension finie), alors pour tout $f \in \mathbb{L}^2(\Omega, \mathcal{A}, \mu)$, il existe un unique projeté orthogonal de f sur A , noté f_A , qui vérifie $f_A = \operatorname{Arg} \inf_{g \in A} \|g - f\|_2$.

2 Applications de la théorie de la mesure et de l'intégration en Probabilités

2.1 Espérance de variables aléatoires

Définition. Soit X une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ un espace probabilisé. Alors si $X \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, on définit l'espérance de X par le nombre $\mathbb{E}X = \int X d\mathbb{P}$. Plus généralement, si $\phi : \mathbb{R} \mapsto \mathbb{R}$ est borélienne et si $\phi(X) \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, on définit l'espérance de $\phi(X)$ par $\mathbb{E}\phi(X) = \int \phi(X) d\mathbb{P}$.

Propriété. Si X est une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, si $\phi : \mathbb{R} \mapsto \mathbb{R}$ est borélienne telle que $\phi(X) \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, et si \mathbb{P}_X est la mesure de probabilité de X alors :

$$\mathbb{E}\phi(X) = \int_{\mathbb{R}} \phi(x) d\mathbb{P}_X(x).$$

Conséquence. • Si \mathbb{P}_X est absolument continue par rapport à la mesure de Lebesgue (donc X est une v.a. dite absolument continue), de densité f_X , alors $\mathbb{E}\phi(X) = \int_{\mathbb{R}} \phi(x)f_X(x)dx$.

• Si \mathbb{P}_X est absolument continue par rapport à la mesure de comptage sur \mathbb{N} (donc X est une v.a. dite discrète), de densité p_X , alors $\mathbb{E}\phi(X) = \sum_{k=0}^{\infty} p_X(k)\phi(k)$.

Propriété. 1. Soit X et Y des variables aléatoires telles que X et $Y \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$. Alors pour tout $(a, b) \in \mathbb{R}^2$, $aX + bY \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$ et

$$\mathbb{E}(aX + bY) = a\mathbb{E}X + b\mathbb{E}Y.$$

2. Soit X une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, et soit $A \in \mathcal{A}$. Alors $\mathbb{E}(\mathbb{I}_A(X)) = \mathbb{P}(X \in A)$.

3. Soit X et Y des variables aléatoires telles que $X \in \mathbb{L}^p(\Omega, \mathcal{A}, \mathbb{P})$ et $Y \in \mathbb{L}^q(\Omega, \mathcal{A}, \mathbb{P})$ avec $\frac{1}{p} + \frac{1}{q} = 1$ et $p > 1$, $q > 1$. Alors $X.Y \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$ et

$$\mathbb{E}|X.Y| \leq (\mathbb{E}|X|^p)^{1/p}(\mathbb{E}|Y|^q)^{1/q}.$$

4. Soit X et Y des variables aléatoires telles que X et $Y \in \mathbb{L}^p(\Omega, \mathcal{A}, \mathbb{P})$, avec $p \geq 1$. Alors $X + Y \in \mathbb{L}^p(\Omega, \mathcal{A}, \mathbb{P})$ et

$$(\mathbb{E}|X + Y|^p)^{1/p} \leq (\mathbb{E}|X|^p)^{1/p} + (\mathbb{E}|Y|^p)^{1/p}.$$

5. Soit X une variable aléatoire telle que $X \in \mathbb{L}^p(\Omega, \mathcal{A}, \mathbb{P})$ pour $p > 0$. Alors pour tout $0 < r \leq p$, $X \in \mathbb{L}^r(\Omega, \mathcal{A}, \mathbb{P})$ et

$$(\mathbb{E}|X|^r)^{1/r} \leq (\mathbb{E}|X|^p)^{1/p}.$$

6. Si X est une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, si $\phi : \mathbb{R} \mapsto \mathbb{R}$ est une fonction borélienne convexe telle que X et $\phi(X) \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, alors

$$\mathbb{E}(\phi(X)) \geq \phi(\mathbb{E}X).$$

Définition. Pour X et Y des variables aléatoires telles que X et $Y \in \mathbb{L}^2(\Omega, \mathcal{A}, \mathbb{P})$, on définit la covariance de X et Y par

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)];$$

On appelle variance de X , $\text{var}(X) = \text{cov}(X, X) = \mathbb{E}[(X - \mathbb{E}X)^2] = \mathbb{E}(X^2) - (\mathbb{E}X)^2$.

Propriété. Sur $\mathbb{L}^2(\Omega, \mathcal{A}, \mathbb{P})$, $\text{cov}(\cdot, \cdot)$ définit un produit scalaire. De plus

$$|\text{cov}(X, Y)|^2 \leq \text{var}(X).\text{var}(Y).$$

2.2 Fonction de répartition et quantiles d'une loi de probabilité

Il y a une correspondance bijective entre la connaissance de \mathbb{P}_X et celle de $F_X = \mathbb{P}_X(\cdot - \infty, x]$. La fonction de répartition permet également de définir les quantiles qui sont essentiels à la construction d'intervalles de confiance et de test.

Soit $\alpha \in [0, 1]$. Des propriétés de la fonction de répartition, on en déduit qu'il existe $x_\alpha \in \mathbb{R}$, tel que :

$$\lim_{x \rightarrow x_\alpha} F_X(x) \leq \alpha \leq F_X(x_\alpha). \quad (1)$$

Soit $I_\alpha = \{x_\alpha \in \mathbb{R} \text{ tel que } x_\alpha \text{ vérifie (1)}\}$. On appelle **quantile** (ou fractile, ou percentile en anglais) d'ordre α de la loi \mathbb{P}_X , noté q_α , le milieu de l'intervalle I_α . Evidemment, lorsque X admet une distribution absolument continue par rapport à la mesure de Lebesgue, $q_\alpha = F_X^{-1}(\alpha)$, où F_X^{-1} désigne la fonction réciproque de F_X .

Deux cas particuliers sont à connaître :

- 1/ pour $\alpha = 0.5$, $q_{0.5}$ est appelé la **médiane** de \mathbb{P}_X ;
- 2/ pour $\alpha = 0.25$ et $\alpha = 0.75$ (respectivement), $q_{0.25}$ et $q_{0.75}$ sont appelés premier et troisième **quartile** (respectivement) de \mathbb{P}_X .
- 3/ pour $\alpha = 0.1, \dots, 0.9$, on parlera de **décile** de \mathbb{P}_X .

2.3 Principales lois de probabilités

Loi uniforme discrète :

C'est la loi de probabilité discrète à valeurs dans $\{x_1, \dots, x_n\}$ telle que

$$\mathbb{P}(X = x_i) = \frac{1}{n}.$$

On alors : $\mathbb{E}X = \frac{1}{n}(x_1 + \dots + x_n)$ et $\text{var}(X) = \frac{1}{n}(x_1^2 + \dots + x_n^2) - (\mathbb{E}X)^2$.

Loi de Bernoulli :

C'est la loi de probabilité discrète notée $\mathcal{B}(p)$ à valeurs dans $\{0, 1\}$ telle que

$$\mathbb{P}(X = 1) = p \quad \text{et} \quad \mathbb{P}(X = 0) = 1 - p.$$

On alors : $\mathbb{E}X = p$ et $\text{var}(X) = p(1 - p)$.

Loi binomiale :

C'est la loi de probabilité discrète notée $\mathcal{B}(n, p)$ à valeurs dans $\{0, 1, \dots, n\}$ telle que

$$\mathbb{P}(X = k) = C_n^k \cdot p^k \cdot (1 - p)^{n-k} \quad \text{pour } k \in \{0, 1, \dots, n\}.$$

On alors : $X = X_1 + \dots + X_n$, où (X_i) est une suite de v.a.i.i.d. de loi $\mathcal{B}(p)$, d'où $\mathbb{E}X = n \cdot p$ et $\text{var}(X) = n \cdot p(1 - p)$.

Loi de Poisson :

C'est la loi de probabilité discrète notée $\mathcal{P}(\theta)$ à valeurs dans \mathbb{N} telle que

$$\mathbb{P}(X = k) = \frac{\theta^k}{k!} \cdot e^{-\theta} \quad \text{pour } k \in \mathbb{N}.$$

On alors $\mathbb{E}X = \theta$ et $\text{var}(X) = \theta$.

Loi uniforme sur $[a, b]$:

Cette loi est généralement notée $\mathcal{U}([a, b])$, où $-\infty < a < b < \infty$. C'est la loi de probabilité à valeurs dans $[a, b]$ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{b - a} \cdot \mathbb{I}_{x \in [a, b]}.$$

On a alors $\mathbb{E}X = \frac{b + a}{2}$ et $\text{var}(X) = \frac{(b - a)^2}{12}$.

Loi Gamma :

Cette loi est généralement notée $\gamma(p, \theta)$, où $p > 0$ et $\theta > 0$. C'est la loi de probabilité à valeurs dans \mathbb{R}_+ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{\theta^p}{\Gamma(p)} \cdot e^{-\theta \cdot x} \cdot x^{p-1} \cdot \mathbb{I}_{x \in \mathbb{R}_+}.$$

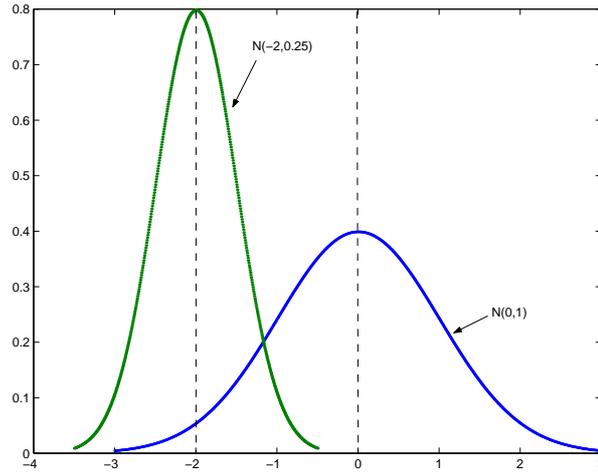


Figure 1: Représentation de la densité de la loi $\mathcal{N}(0, 1)$ et de la loi $\mathcal{N}(-2, 0.5^2)$

On a alors $\mathbb{E}X = \frac{p}{\theta}$ et $\text{var}(X) = \frac{p}{\theta^2}$.

Si $X \sim \gamma(p, \theta)$ et $Y \sim \gamma(q, \theta)$ avec X et Y indépendantes et $p > 0$ et $q > 0$, alors $X + Y \sim \gamma(p + q, \theta)$.

Pour $p = 1$, la loi $\gamma(p, \theta)$ est la loi exponentielle $\mathcal{E}(\theta)$.

Loi Béta :

Cette loi est généralement notée $\beta(p, \theta)$, où $p > 0$ et $q > 0$. C'est la loi de probabilité à valeurs dans $[0, 1]$ de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{x^p(1-x)^{q-1}}{B(p, q)} \cdot x^{p-1} \cdot \mathbb{I}_{x \in [0, 1]}, \quad \text{où } B(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}.$$

On a alors $\mathbb{E}X = \frac{B(p+1, q)}{B(p, q)}$ et $\text{var}(X) = \frac{p \cdot q}{(p+q)^2(p+q+1)}$.

Si $X \sim \gamma(p, \theta)$ et $Y \sim \gamma(q, \theta)$ avec X et Y indépendantes et $p > 0$ et $q > 0$, alors $\frac{X}{X+Y} \sim \beta(p, q)$.

Pour $p = 1$, la loi $\gamma(p, \theta)$ est la loi exponentielle $\mathcal{E}(\theta)$.

Loi normale (ou gaussienne) centrée réduite :

Cette loi est généralement notée $\mathcal{N}(0, 1)$. C'est la loi de probabilité à valeurs dans \mathbb{R} de densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right).$$

On a :

$$\mathbb{E}(X) = 0 \quad \text{et} \quad \text{var}(X) = 1.$$

Loi normale (ou gaussienne) de moyenne m et de variance σ^2 :

Si Z suit la loi $\mathcal{N}(0, 1)$, $X = m + \sigma Z$ suit par définition la loi $\mathcal{N}(m, \sigma^2)$, loi normale d'espérance m et de variance σ^2 . La densité de X est donnée par :

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right).$$

La figure A.1. représente la densité de la loi normale centrée réduite et celle d'une loi normale non centrée et non réduite.

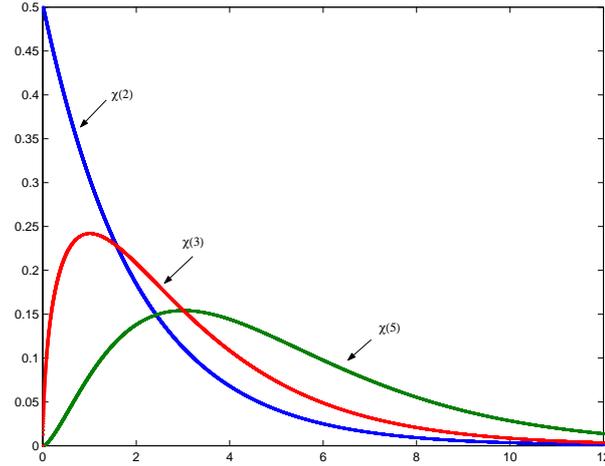


Figure 2: Représentation de la densité des lois $\chi^2(2)$, $\chi^2(3)$ et $\chi^2(5)$

A partir de la loi gaussienne, on peut en déduire les lois suivantes.

Loi du χ^2 à n degrés de libertés :

Soit X_1, \dots, X_n , n variables aléatoires indépendantes de loi $\mathcal{N}(0, 1)$, alors

$$S = X_1^2 + \dots + X_n^2$$

suit une loi du χ^2 à n degrés de libertés, loi notée $\chi^2(n)$. Cette loi est à valeurs dans \mathbb{R}_+ , d'espérance n et de variance $2n$. C'est aussi la loi Gamma $\gamma(n/2, 1/2)$, c'est-à-dire que $X \sim \chi^2(n)$ admet pour densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{2^{n/2} \cdot \Gamma(n/2)} x^{n/2-1} \exp\left(-\frac{x}{2}\right) \cdot \mathbb{I}_{\{x \geq 0\}},$$

où la fonction Gamma est telle que $\Gamma(a) = \int_0^\infty x^{a-1} \cdot e^{-x}$ pour $a \geq 0$. Enfin, si X suit une loi $\chi^2(n)$, par définition on dira que $Y = \sigma^2 \cdot X$ suit une loi $\sigma^2 \cdot \chi^2(n)$. La figure A.2. exhibe trois tracés différents de densité de loi du χ^2 .

Loi de Student à n degrés de libertés :

La loi de Student à n degrés de liberté, notée $T(n)$, est la loi du quotient

$$T = \frac{N}{\sqrt{S/n}}$$

où N suit une loi $\mathcal{N}(0, 1)$ et S suit une loi $\chi^2(n)$, N et S étant deux variables aléatoires indépendantes. Il est également possible de déterminer la densité d'une telle loi par rapport à la mesure de Lebesgue, à savoir,

$$f_X(x) = \frac{1}{\sqrt{n} \cdot B(1/2, n/2)} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2},$$

où la fonction Beta est telle que $B(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)}$ pour $a > 0$ et $b > 0$. La figure A.3. illustre deux exemples de cette densité, que l'on compare également avec la densité de la loi normale centrée réduite.

Remarque : Par la loi des grands nombres, plus n est grand, plus S est proche de son espérance qui vaut n . Le dénominateur est donc proche de 1. Il s'ensuit que la loi $T(n)$ est d'autant plus proche d'une loi normale

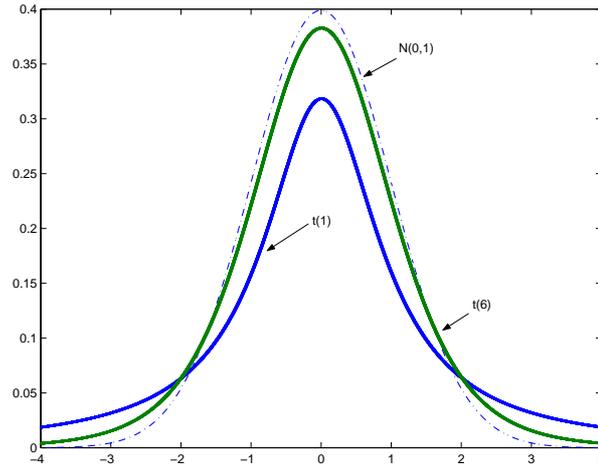


Figure 3: Représentation de la densité des lois $T(2)$, $T(6)$, à comparer avec celle de la loi $\mathcal{N}(0, 1)$

que n est grand.

Un des principaux intérêt de la loi de Student réside dans le fait que si X_1, \dots, X_n sont n variables aléatoires indépendantes de loi $\mathcal{N}(m, \sigma^2)$, si on considère la moyenne et la variance empiriques :

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n) \quad \text{et} \quad \bar{\sigma}_n^2 = \frac{1}{n-1} ((X_1 - \bar{X}_n)^2 + \dots + (X_n - \bar{X}_n)^2),$$

alors

$$T = \frac{\sqrt{n} \cdot (\bar{X}_n - m)}{\sqrt{\bar{\sigma}_n^2}}$$

suit une loi de Student à $(n - 1)$ degrés de liberté.

Loi de Fisher à n_1 et n_2 degrés de liberté :

Soit S_1 et S_2 deux variables aléatoires indépendantes de loi respectives $\chi^2(n_1)$ et $\chi^2(n_2)$. Alors par définition :

$$F = \frac{S_1/n_1}{S_2/n_2}$$

suit une loi de Fisher à n_1 et n_2 degrés de liberté, notée $F(n_1, n_2)$.

Remarque : Par les mêmes considérations que précédemment, la loi F est d'autant plus proche de 1 que les degrés de liberté n_1 et n_2 sont grands.

On a également les propriétés suivantes :

- Si F suit une loi $F(n_1, n_2)$, alors la loi de $\frac{n_1}{n_2} F$ est une loi beta de seconde espèce de paramètres $(n_1/2, n_2/2)$, c'est-à-dire que F est à valeurs dans \mathbb{R}_+ et admet la densité par rapport à la mesure de Lebesgue :

$$f_X(x) = \frac{1}{B(n_1/2, n_2/2)} n_1^{n_1/2} \cdot n_2^{n_2/2} \frac{x^{n_1/2-1}}{(n_2 + n_1 \cdot x)^{(n_1+n_2)/2}} \mathbb{I}_{\{x \geq 0\}},$$

la notation B désignant encore la fonction Beta.

- Si $F \sim F(n_1, n_2)$, alors $\mathbb{E}(F) = \frac{n_2}{n_2 - 2}$ lorsque $n_2 > 2$ et $\text{var}(F) = \frac{2n_2^2(n_1 + n_2 - 2)}{n_1(n_2 - 4)(n_2 - 2)^2}$ lorsque $n_2 > 4$.
- Si T suit une loi de Student $T(n)$, alors T^2 suit une loi de Fisher $F(1, n)$.

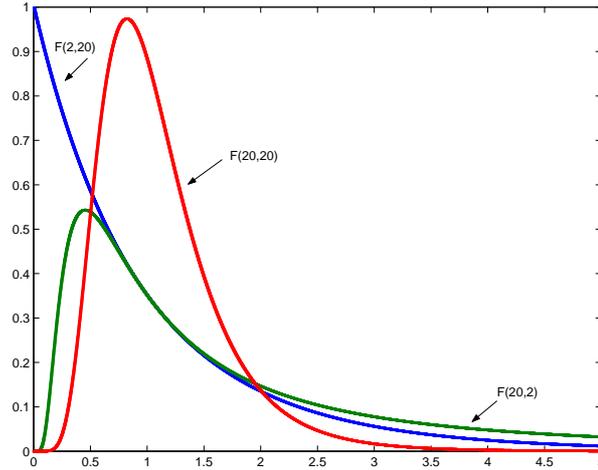


Figure 4: Représentation de la densité des lois $F(2, 20)$, $F(20, 2)$ et $F(20, 20)$

La figure A.4. donne une idée de la distribution d'une loi de Fisher pour différents choix des paramètres.

2.4 Indépendance

Définition. Soit $(\Omega, \mathcal{A}, \mathbb{P})$ un espace probabilisé.

- Soit $(A_i)_{i \in I}$ une famille dénombrable d'événements de \mathcal{A} . On dit que les événements $(A_i)_{i \in I}$ sont indépendants si et seulement si pour tous les sous-ensembles finis $K \subset I$,

$$\mathbb{P} \left(\bigcap_{i \in K} A_i \right) = \prod_{i \in K} \mathbb{P}(A_i).$$

- Soit $(\mathcal{A}_i)_{i \in I}$ une famille de sous-tribus de \mathcal{A} (donc pour tout $i \in I$, $\mathcal{A}_i \subset \mathcal{A}$). On dit que les tribus $(\mathcal{A}_i)_{i \in I}$ sont indépendantes si et seulement si pour tous les sous-ensembles finis $K \subset I$, et pour tous les événements $A_k \in \mathcal{A}_k$ avec $k \in K$, les A_k sont indépendants.
- Soit $(X_i)_{i \in I}$ des variables aléatoires sur (Ω, \mathcal{A}) à valeurs dans $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$. On dit que les v.a. $(X_i)_{i \in I}$ sont indépendantes si et seulement si les tribus engendrées $(X_i^{-1}(\mathcal{B}(\mathbb{R})))_{i \in I}$ sont indépendantes.

Proposition. Si (X_1, \dots, X_n) sont des variables aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$. Alors les (X_i) sont indépendantes si et seulement si $\mathbb{P}_{(X_1, \dots, X_n)} = \bigotimes_{i=1}^n \mathbb{P}_{X_i}$.

Proposition. Si $(X_i)_{i \in I}$ sont des variables aléatoires indépendantes sur $(\Omega, \mathcal{A}, \mathbb{P})$. Alors les (X_i) sont indépendantes si et seulement si pour tout $J \subset I$, J fini, pour toutes fonctions boréliennes $(g_j)_{j \in J}$ telles que $g_j(X_j)$ soit intégrable, alors

$$\mathbb{E} \left(\prod_{j \in J} g_j(X_j) \right) = \prod_{j \in J} \mathbb{E}(g_j(X_j)).$$

Corollaire. (X_1, \dots, X_n) sont des variables aléatoires indépendantes si et seulement si pour tout $(t_1, \dots, t_n) \in \mathbb{R}^n$,

$$\phi_{(X_1, \dots, X_n)}(t_1, \dots, t_n) = \prod_{j=1}^n \phi_{X_j}(t_j).$$

Lemme (Lemme de Borel-Cantelli). Soit $(A_n)_{n \in \mathbb{N}}$ une suite d'événements sur $(\Omega, \mathcal{A}, \mathbb{P})$.

1. Si $\sum \mathbb{P}(A_n) < +\infty$ alors $\mathbb{P}(\limsup A_n) = 0$.
2. Si les (A_n) sont indépendants, $\sum \mathbb{P}(A_n) = +\infty$ implique que $\mathbb{P}(\limsup A_n) = 1$.

2.5 Vecteurs aléatoires

Définition. On dit que X est un vecteur aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, un espace probabilisé, si X est une fonction mesurable de (Ω, \mathcal{A}) dans $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$.

Définition. Soit X un vecteur aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeurs dans \mathbb{R}^d . Alors la loi (ou mesure) de probabilité de X , \mathbb{P}_X , est définie de façon univoque à partir de la fonction de répartition de X , telle que pour $x = (x_1, \dots, x_d)$,

$$F_X(x) = \mathbb{P}_X\left(\prod_{i=1}^d]-\infty, x_i]\right) = \mathbb{P}(X \in \prod_{i=1}^d]-\infty, x_i]).$$

Propriété. Soit X un vecteur aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeurs dans \mathbb{R}^d . On suppose que $X = (X_1, \dots, X_d)$. Alors les X_i sont des variables aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$, de fonction de répartition

$$F_{X_i}(x_i) = \lim_{\substack{x_j \rightarrow +\infty \\ j \neq i}} F_X(x_1, \dots, x_i, \dots, x_d).$$

Les mesures de probabilités P_{X_i} déterminées de façon univoque à partir des F_{X_i} sont appelées lois marginales de X .

On se place maintenant dans la base canonique orthonormale de \mathbb{R}^d . Si Z est un vecteur aléatoire à valeurs sur \mathbb{R}^d , on définit $\mathbb{E}(Z)$, le vecteur dont les coordonnées sont les espérances des coordonnées de Z . Ainsi, si dans la base canonique de \mathbb{R}^d , $Z = (Z_1, \dots, Z_d)'$,

$$\mathbb{E}(Z) = \mathbb{E} \begin{pmatrix} Z_1 \\ \vdots \\ Z_d \end{pmatrix} = \begin{pmatrix} \mathbb{E}(Z_1) \\ \vdots \\ \mathbb{E}(Z_d) \end{pmatrix}.$$

De la même manière, on définira l'espérance d'une matrice dont les coordonnées sont des variables aléatoires par la matrice dont les coordonnées sont les espérances de chacune de ces variables aléatoires.

Ceci nous permet de définir la matrice de variance-covariance de Z de la manière suivante :

$$\text{var}(Z) = \mathbb{E}[(Z - \mathbb{E}(Z))(Z - \mathbb{E}(Z))']$$

donc si $Z = (Z_1, \dots, Z_d)'$,

$$\text{var} \begin{pmatrix} Z_1 \\ \vdots \\ Z_d \end{pmatrix} = \begin{pmatrix} \text{var}(Z_1) & \text{Cov}(Z_1, Z_2) & \cdots & \text{Cov}(Z_1, Z_d) \\ \text{Cov}(Z_1, Z_2) & \text{var}(Z_2) & \cdots & \text{Cov}(Z_2, Z_d) \\ \vdots & \vdots & \cdots & \vdots \\ \text{Cov}(Z_1, Z_d) & \text{Cov}(Z_2, Z_d) & \cdots & \text{var}(Z_d) \end{pmatrix}$$

matrice (d, d) dont les éléments diagonaux sont les variances et les éléments non diagonaux sont les covariances des coordonnées de Z (remarquons que la variance de Z_1 est aussi la covariance de Z_1 et de Z_1).

On vérifie également le résultat suivant : si C est une matrice (p, d) à coordonnées constituées de réels constants et si Z est un vecteur aléatoire à valeurs dans \mathbb{R}^d , alors $C \cdot Z$ est un vecteur de taille p de matrice de variance-covariance

$$\text{var}(C \cdot Z) = C \cdot \text{var}(Z) \cdot C'.$$

En particulier, si p vaut 1, alors $C = h'$ où h est un vecteur de taille d , et :

$$\text{var}(h' \cdot Z) = h' \cdot \text{var}(Z) \cdot h.$$

Notez que cette dernière quantité est un scalaire. Soit Y_1, \dots, Y_d des variables aléatoires indépendantes de même loi $\mathcal{N}(0, \sigma^2)$, indépendantes (ce qui, dans le cas gaussien, est équivalent à $\text{cov}(Y_i, Y_j) = 0$ pour $i \neq j$). On considère le vecteur $Y = (Y_1, \dots, Y_d)'$. En raison de l'indépendance, Y est un vecteur gaussien admettant

une densité f_Y (par rapport à la mesure de Lebesgue sur \mathbb{R}^d) qui est le produit des densités de chacune des coordonnées, soit :

$$\begin{aligned} f_Y(y_1, \dots, y_d) &= f_{Y_1}(y_1) \times f_{Y_2}(y_2) \times \dots \times f_{Y_d}(y_d) \\ &= (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{1}{2\sigma^2}(y_1^2 + \dots + y_d^2)\right) \\ &= (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{\|y\|^2}{2\sigma^2}\right), \end{aligned}$$

avec $y = (y_1, \dots, y_d)$. On voit donc que la densité de Y ne dépend que de la norme $\|Y\|$: elle est constante sur toutes les sphères centrées en zéro. Cela implique qu'elle est invariante par rotation ou symétrie orthogonale d'axe passant par 0 : elle est invariante par toutes les isométries de \mathbb{R}^d : on dira que Y suit une loi gaussienne isotrope. Rappelons que les isométries correspondent à des changements de bases orthonormées (BON). En conséquence, on a la première propriété importante :

Propriété. Soit Y un vecteur aléatoire de \mathbb{R}^d de loi normale isotrope variance σ^2 , c'est-à-dire que dans une BON les coordonnées de Y vérifient $\mathbb{E}(Y) = 0$ et $\text{var}(Y) = \sigma^2 \cdot \text{Id}$. Alors les coordonnées de Y dans toute BON sont encore des lois $\mathcal{N}(0, \sigma^2)$ indépendantes.

Voici maintenant l'un des résultats (encore appelé Théorème de Cochran) que nous utilisons le plus et nous en donnons donc une démonstration.

Théorème (Théorème de Cochran). Soit E_1 et E_2 , deux sous-espaces vectoriels orthogonaux de $E = \mathbb{R}^d$ de dimensions respectives k_1 et k_2 et soit Y un vecteur aléatoire de \mathbb{R}^d de loi normale centrée isotrope de variance σ^2 . Alors $P_{E_1}(Y)$ et $P_{E_2}(Y)$ sont deux variables aléatoires gaussiennes centrées indépendantes et $\|P_{E_1}(Y)\|^2$ (resp. $\|P_{E_2}(Y)\|^2$) est une loi $\sigma^2 \cdot \chi^2(k_1)$ (resp. $\sigma^2 \cdot \chi^2(k_2)$). Ce théorème se généralise naturellement pour $2 < m \leq d$ sous-espaces vectoriels orthogonaux $(E_i)_{1 \leq i \leq m}$ de $E = \mathbb{R}^d$.

Démonstration : Soit (e_1, \dots, e_{k_1}) et $(e_{k_1+1}, \dots, e_{k_1+k_2})$ deux BON de E_1 et E_2 (respectivement). L'ensemble de ces deux bases peut être complété en

$$(e_1, \dots, e_{k_1}, e_{k_1+1}, \dots, e_{k_1+k_2}, e_{k_1+k_2+1}, \dots, e_d)$$

pour former une BON de \mathbb{R}^d (du fait que E_1 et E_2 sont orthogonaux).

Soit (Y_1, \dots, Y_d) , les coordonnées de Y dans cette base; elles sont indépendantes de loi $\mathcal{N}(0, \sigma^2)$ car le changement de base est orthonormal et nous avons vu que la distribution de Y était conservé par transformation isométrique. Comme

$$\begin{aligned} P_{E_1}(Y) = Y_1 e_1 + \dots + Y_{k_1} e_{k_1} &\implies \|P_{E_1}(Y)\|^2 = \sigma^2 \left(\left(\frac{Y_1}{\sigma}\right)^2 + \dots + \left(\frac{Y_{k_1}}{\sigma}\right)^2 \right) \\ P_{E_2}(Y) = Y_{k_1+1} e_{k_1+1} + \dots + Y_{k_1+k_2} e_{k_1+k_2} &\implies \|P_{E_2}(Y)\|^2 = \sigma^2 \left(\left(\frac{Y_{k_1+1}}{\sigma}\right)^2 + \dots + \left(\frac{Y_{k_1+k_2}}{\sigma}\right)^2 \right). \end{aligned}$$

On voit bien ainsi l'indépendance entre les deux projections et le fait que la loi de $\|P_{E_1}(Y)\|^2$ (resp. $\|P_{E_2}(Y)\|^2$) est une loi $\sigma^2 \cdot \chi^2(k_1)$ (resp. $\sigma^2 \cdot \chi^2(k_2)$). ■

On peut définir plus généralement un vecteur gaussien Y à valeurs dans \mathbb{R}^d (non dégénéré), d'espérance $\mu \in \mathbb{R}^d$ et de matrice de variance-covariance Σ quelconques (du moment que Σ soit une matrice de Toeplitz définie positive). Cela équivaut à définir un vecteur aléatoire de densité par rapport à la mesure de Lebesgue sur \mathbb{R}^d ,

$$f_Y(y) = \frac{(2\pi)^{-n/2}}{|\Sigma|} \exp\left(-\frac{1}{2}(y - \mu)' \cdot \Sigma^{-1} \cdot (y - \mu)\right),$$

pour $y \in \mathbb{R}^d$, et avec $|\Sigma|$ le déterminant de la matrice Σ . Remarquons une nouvelle fois que l'espérance et la variance définissent complètement la loi de probabilité d'un vecteur gaussien.

A partir des propriétés générales sur les vecteurs aléatoires, on obtient le fait que :

Propriété. Soit Y un vecteur gaussien à valeurs dans \mathbb{R}^d (non dégénéré), d'espérance $\mu \in \mathbb{R}^d$ et de matrice de variance-covariance Σ . Soit C une matrice réelle de taille (p, d) où $p \in \mathbb{N}^*$. Alors $C \cdot Y$ est un vecteur gaussien tel que :

$$C \cdot Y \sim \mathcal{N}(C \cdot \mu, C \cdot \Sigma \cdot C')$$

On en déduit les conséquences suivantes :

- si Y est un vecteur gaussien isotrope de \mathbb{R}^d de variance σ^2 et h un vecteur de \mathbb{R}^d , alors $h' \cdot Y$ est une combinaison linéaire des coordonnées de Y tel que :

$$h' \cdot Y \text{ suit la loi } \mathcal{N}(0, \sigma^2 \cdot h' \cdot h) = \mathcal{N}(0, \sigma^2 \cdot \|h\|^2)$$

- si Y est un vecteur gaussien d'espérance μ et de matrice de variance Σ et si h un vecteur de \mathbb{R}^d , alors $h' \cdot Y$ est une combinaison linéaire des coordonnées de Y et :

$$h' \cdot Y \text{ suit la loi unidimensionnelle } \mathcal{N}(h' \cdot \mu, h' \cdot \Sigma \cdot h)$$

(Pour une présentation plus détaillée des notions sur les vecteurs gaussiens on peut consulter le livre P. Toulouse, 1999, chap.2)

2.6 Fonctions caractéristiques et génératrices

Définition. Soit X un vecteur aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeurs dans \mathbb{R}^d . La fonction caractéristique de X est la fonction $\phi_X : \mathbb{R}^d \mapsto \mathbb{C}$ telle que

$$\phi_X(t) = \mathbb{E}[\exp(i \langle t, X \rangle)] = \int_{\mathbb{R}^d} e^{i \langle t, x \rangle} d\mathbb{P}_X(x),$$

où $\langle \cdot, \cdot \rangle$ désigne le produit scalaire euclidien sur \mathbb{R}^d tel que $\langle t, x \rangle = \sum_{i=1}^d t_i x_i$ pour $t = (t_1, \dots, t_d)$ et $x = (x_1, \dots, x_d)$.

Remarque. La fonction génératrice existe sur \mathbb{R} et $\phi_X(0) = 1$. ϕ_X est aussi la transformée de Fourier de la mesure \mathbb{P}_X .

Théorème. Soit X et Y des vecteurs aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeurs dans \mathbb{R}^d , de lois \mathbb{P}_X et \mathbb{P}_Y . Alors $\mathbb{P}_X = \mathbb{P}_Y$ si et seulement si $\phi_X = \phi_Y$.

Théorème (Théorème d'inversion). Si X est un vecteur aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeurs dans \mathbb{R}^d et si ϕ_X est une fonction intégrable par rapport à la mesure de Lebesgue λ_d sur \mathbb{R}^d , alors X admet une densité f_X par rapport à λ_d telle que pour $x \in \mathbb{R}^d$,

$$f_X(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{-i \langle t, x \rangle} \phi_X(t) dt.$$

Proposition. Si X est une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$ de fonction génératrice ϕ_X . Alors si $\mathbb{E}(|X|^n) < +\infty$ (ou $X \in \mathbb{L}^n(\Omega, \mathcal{A}, \mathbb{P})$), ϕ_X est n fois dérivable et $\phi_X^{(n)}(t) = i^n \mathbb{E}(X^n e^{itX})$.

Remarque. Lorsque ces moments existent, on a $i^n \mathbb{E}(X^n) = \phi_X^{(n)}(0)$.

2.7 Convergence de suites de variables aléatoires

Définition. Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$. On dit que

- (X_n) converge en probabilité vers X , noté $X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} X$, lorsque pour tout $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0.$$

- (X_n) converge dans $\mathbb{L}^p(\Omega, \mathcal{A}, \mathbb{P})$ vers X , noté $X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{L}^p} X$, avec $p > 0$, lorsque

$$\lim_{n \rightarrow \infty} \mathbb{E}|X_n - X|^p = 0.$$

- (X_n) converge en loi vers X , noté $X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X$, lorsque,

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x) \text{ pour tout } x \in \mathbb{R} \text{ tel que } F_X \text{ continue en } x.$$

- (X_n) converge presque sûrement vers X , noté $X_n \xrightarrow[n \rightarrow +\infty]{p.s.} X$, lorsque pour tout $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\sup_{m \geq n} |X_m - X| > \varepsilon) = 0.$$

Propriété. 1. $p.s.$ et $\mathbb{L}^p \longrightarrow \mathcal{P} \longrightarrow \mathcal{L}$.

2. pour $q \geq p$, $\mathbb{L}^q \longrightarrow \mathbb{L}^p$.

3. La convergence en loi n'entraîne pas la convergence en probabilité. Mais $(X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} C) \iff (X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} C)$ pour C une constante.

4. Si g est une fonction borélienne continue alors $(X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} X) \implies (g(X_n) \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} g(X))$.

Propriété. 1. Si pour tout $\varepsilon > 0$, $\sum_{n=0}^{\infty} \mathbb{P}(|X_n - X| > \varepsilon) < +\infty$ alors $X_n \xrightarrow[n \rightarrow +\infty]{p.s.} X$ (application du Lemme de Borel-Cantelli).

2. Si il existe $r > 0$ tel que $\mathbb{E}(|X_n|^r) < +\infty$ et $\sum_{n=0}^{\infty} \mathbb{E}(|X_n - X|^r) < +\infty$ alors $X_n \xrightarrow[n \rightarrow +\infty]{p.s.} X$.

Théorème (Loi faible des Grands Nombres). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si $\mathbb{E}(|X_i|) < +\infty$,

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} m = \mathbb{E}X_i.$$

Théorème (Loi forte des Grands Nombres). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si $\mathbb{E}(|X_i|) < +\infty$,

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow[n \rightarrow +\infty]{p.s.} m = \mathbb{E}X_i.$$

Théorème (Théorème de la limite centrale). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires indépendantes et identiquement distribuées. Alors si $\sigma^2 = \mathbb{E}X_i^2 < +\infty$, et $m = \mathbb{E}X_i$,

$$\sqrt{n} \frac{\bar{X}_n - m}{\sigma} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

Théorème (Loi forte des Grands Nombres multidimensionnelle). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de vecteurs aléatoires à valeurs dans \mathbb{R}^d , indépendants et identiquement distribués. Alors si $\mathbb{E}(\|X_i\|) < +\infty$ (pour $\|\cdot\|$ une norme sur \mathbb{R}^d),

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n} \xrightarrow[n \rightarrow +\infty]{p.s.} m = \mathbb{E}X_i.$$

Théorème (Théorème de la limite centrale multidimensionnel). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de vecteurs aléatoires à valeurs dans \mathbb{R}^d , indépendants et identiquement distribués. Alors si Σ matrice de covariance de chaque X_i existe, et $m = \mathbb{E}X_i$,

$$\sqrt{n}(\bar{X}_n - m) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, \Sigma).$$

Théorème (Delta-method). Soit $(X_n)_{n \in \mathbb{N}}$ une suite de vecteurs aléatoires à valeurs dans \mathbb{R}^d , indépendants et identiquement distribués, telle que Σ matrice de covariance de chaque X_i existe, et $m = \mathbb{E}X_i$. Soit $g: \mathbb{R}^d \rightarrow \mathbb{R}^p$ une fonction de classe \mathcal{C}^1 sur un voisinage autour de m , de matrice Jacobienne $J_g(m)$ en m . Alors,

$$\sqrt{n}(g(\bar{X}_n) - g(m)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, J_g(m) \cdot \Sigma \cdot J_g'(m)).$$

2.8 Espérance conditionnelle

Définition. Soit Y une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$. Si \mathcal{B} est une sous-tribu de \mathcal{A} et si $Y \in \mathbb{L}^2(\Omega, \mathcal{A}, \mathbb{P})$. Alors on note $\mathbb{E}(Y | \mathcal{B})$ la projection orthogonale de Y sur $\mathbb{L}^2(\Omega, \mathcal{B}, \mathbb{P})$, appelée espérance conditionnelle de Y sachant \mathcal{B} . Ainsi :

$$\mathbb{E}|Y - \mathbb{E}(Y | \mathcal{B})|^2 = \inf_{Z \in \mathbb{L}^2(\Omega, \mathcal{B}, \mathbb{P})} \left\{ \mathbb{E}|Y - Z|^2 \right\}.$$

Par extension, si $Y \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, on définit l'espérance conditionnelle par rapport à \mathcal{B} , comme l'unique (p.s.) variable aléatoire, \mathcal{B} -mesurable vérifiant p.s. :

$$\int_B \mathbb{E}(Y | \mathcal{B}) d\mathbb{P} = \int_B Y d\mathbb{P}, \quad \text{pour tout } B \in \mathcal{B}.$$

Définition. Par convention, si X un vecteur aléatoire à valeurs dans \mathbb{R}^n sur $(\Omega, \mathcal{A}, \mathbb{P})$ et si Y une variable aléatoire sur $(\Omega, \mathcal{A}, \mathbb{P})$, on note $\mathbb{E}(Y | X) = \mathbb{E}(Y | X^{-1}(\mathcal{B}(\mathbb{R})))$.

Propriété. 1. Lemme de Doob : Pour $Y \in \mathbb{L}^1(\Omega, \mathcal{A}, \mathbb{P})$, et X une v.a. de $(\Omega, \mathcal{A}, \mathbb{P})$, alors p.s. $\mathbb{E}(Y | X) = h(X)$, avec h une fonction borélienne.

2. Pour Y_1 et Y_2 deux variables aléatoires sur $(\Omega, \mathcal{A}, \mathbb{P})$, et $(a, b, c) \in \mathbb{R}^3$, alors

$$\mathbb{E}(aY_1 + bY_2 + c | \mathcal{B}) = a\mathbb{E}(Y_1 | \mathcal{B}) + b\mathbb{E}(Y_2 | \mathcal{B}) + c.$$

3. Si $Y_1 \leq Y_2$, alors $\mathbb{E}(Y_1 | \mathcal{B}) \leq \mathbb{E}(Y_2 | \mathcal{B})$.

4. Le Lemme de Fatou, les théorèmes de Beppo-Levi, Lebesgue et Jensen s'appliquent avec l'espérance conditionnelle.

5. Si $Y \in \mathbb{L}^2(\Omega, \mathcal{B}, \mathbb{P})$, alors $\mathbb{E}(Y | \mathcal{B}) = Y$; ainsi $\mathbb{E}(g(X) | X) = g(X)$ pour g une fonction mesurable réelle.

6. On a $\mathbb{E}(\mathbb{E}(Y | \mathcal{B})) = \mathbb{E}Y$.

7. Si $Y^{-1}(\mathcal{B}(\mathbb{R}))$ et \mathcal{B} sont indépendantes alors $\mathbb{E}(Y | \mathcal{B}) = \mathbb{E}Y$; ainsi, si X et Y sont indépendantes, $\mathbb{E}(Y | X) = \mathbb{E}Y$.

8. Si (X, Y) est un couple de v.a. à valeurs dans \mathbb{R}^2 possédant une densité $f_{(X,Y)}$ par rapport à la mesure de Lebesgue, alors si X est intégrable ,

$$\mathbb{E}(Y | X = x) = \frac{\int_{\mathbb{R}} y \cdot f_{(X,Y)}(x, y) dy}{\int_{\mathbb{R}} f_{(X,Y)}(x, y) dy}, \quad \text{pour tout } x \text{ tel que } \int_{\mathbb{R}} f_{(X,Y)}(x, y) dy > 0.$$

Proposition. Si (Y, X_1, \dots, X_n) est un vecteur gaussien, alors $\mathbb{E}(Y | (X_1, \dots, X_n)) = a_0 + a_1X_1 + \dots + a_nX_n$, où les a_i sont des réels.

3 Estimation paramétrique

3.1 Définitions

Dans toute la suite, on se place sur $(\Omega, \mathcal{A}, \mathbb{P})$ un espace de probabilité. On considère $(X_n)_{n \in \mathbb{N}}$ une suite de variable aléatoire, où chaque X_i est définie sur $(\Omega, \mathcal{A}, \mathbb{P})$ et est à valeur dans $\Omega' \subset \mathbb{R}$.

Définition. • On appelle modèle statistique de dimension n un espace $((\Omega')^n, \mathcal{A}'_n, \mu)$, où \mathcal{A}'_n est une tribu sur $(\Omega')^n$ et μ une mesure de probabilité sur $((\Omega')^n, \mathcal{A}'_n)$.

- On appelle échantillon de taille n du modèle statistique $((\Omega')^n, \mathcal{A}'_n, \mu)$ le vecteur aléatoire (X_1, \dots, X_n) distribuée selon la loi μ . Pour $\omega \in \Omega$, $(X_1(\omega), \dots, X_n(\omega))$ vecteur de \mathbb{R}^n est appelé échantillon observé. C'est à partir et sur ce vecteur que le travail statistique s'effectue (en général).

Définition. On appelle :

- *Modèle statistique paramétrique, une famille de modèle de la forme : $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$.*
- *Modèle statistique semi-paramétrique, une famille de modèle de la forme : $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_{(\theta, f)}, \theta \in \Theta, f \in \mathcal{F})$, où $\Theta \subset \mathbb{R}^p$ et \mathcal{F} n'est pas de dimension finie.*
- *Modèle statistique non-paramétrique, une famille de modèle de la forme : $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_f, f \in \mathcal{F})$, où \mathcal{F} n'est pas de dimension finie.*

Définition. • *On dit que le modèle paramétrique : $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, est dominé par une mesure μ lorsque \mathbb{P}_θ est absolument continue par rapport à μ pour tout $\theta \in \Theta$.*

- *On se place dans le cadre d'un modèle paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ . Pour $(x_1, \dots, x_n) \in (\Omega')^n$, la fonction $\theta \in \Theta \mapsto L_\theta(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mu}(x_1, \dots, x_n)$ est appelée une vraisemblance du modèle statistique.*

Exemple. • *Dans le cas où μ est la mesure de Lebesgue sur \mathbb{R}^n , la vraisemblance sera la densité (classique) en (x_1, \dots, x_n) .*

- *Dans le cas où μ est comptage sur \mathbb{N}^n , la vraisemblance sera la probabilité en (x_1, \dots, x_n) .*
- *Attention ! si le support de \mathbb{P}_θ dépend de θ , la mesure qui domine (ainsi que Ω' et \mathcal{A}'_n) ne peut dépendre de θ : il ne faut pas oublier de le préciser dans l'expression de la vraisemblance.*

Définition. *Lorsque l'on dispose d'un échantillon (X_1, \dots, X_n) du modèle statistique $((\Omega')^n, \mathcal{A}'_n, \mu)$, une statistique \widehat{T}_n est une application mesurable de $((\Omega')^n, \mathcal{A}'_n)$ dans $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, donc un vecteur aléatoire défini sur $(\Omega, \mathcal{A}, \mathbb{P})$ à valeur dans \mathbb{R}^d , et telle que :*

$$\widehat{T}_n = h(X_1, \dots, X_n), \quad \text{où } h : (\Omega')^n \mapsto \mathbb{R}^d \text{ est mesurable.}$$

Exemple. *Estimateur du paramètre d'une loi de Bernoulli.*

Estimateur de l'espérance et de la variance par la moyenne et la variance empirique.

Estimateurs du paramètre θ d'un n-échantillon (X_1, \dots, X_n) de loi uniforme sur $[0, \theta]$.

Test sur la moyenne.

3.2 Statistiques exhaustives

On se place désormais dans le cadre d'une modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ .

Exemple. 1. *Soit le modèle statistique paramétrique $([0, \infty[^n, \mathcal{B}([0, \infty[^n, \mathcal{U}([0, \theta])^{\otimes n}, \theta \in]0, +\infty[)$. On dispose donc d'un n-échantillon (X_1, \dots, X_n) de v.a.i.i.d. suivant une loi uniforme sur $[0, \theta]$. Si on considère $\max\{X_1, \dots, X_n\}$ cela semble suffire pour posséder toute l'information sur θ que contenait (X_1, \dots, X_n) : on a donc résumé l'"information" sur θ contenait (X_1, \dots, X_n) , un vecteur de taille n, par une statistique de taille 1.*

2. *De même, si on considère le modèle statistique paramétrique $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), \mathcal{B}(p)^{\otimes n}, p \in [0, 1])$ (on dispose donc d'un n-échantillon (X_1, \dots, X_n) de v.a.i.i.d. suivant une loi de Bernoulli de paramètre p) alors la statistique $X_1 + \dots + X_n$ contient toute l'"information" sur p contenue dans l'échantillon (X_1, \dots, X_n) .*

Comment exprimer formellement ce fait qu'une statistique puisse résumer à elle seule toute l'information sur le paramètre ?

Définition. *Soit \widehat{T} une statistique du modèle statistique paramétrique dominé à valeurs dans \mathbb{R}^d . On dit que \widehat{T} est une statistique exhaustive si pour toute statistique S intégrable (donc dans $\mathbb{L}^1((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta)$) alors $\mathbb{E}_\theta(S \mid \widehat{T})$ ne dépend (\mathbb{P}_θ -presque sûrement) pas de θ .*

Théorème (Théorème de factorisation de Neyman). *Soit (X_1, \dots, X_n) un n-échantillon et soit \widehat{T} une statistique du modèle statistique paramétrique dominé avec \widehat{T} à valeurs dans \mathbb{R}^d , où $d \in \mathbb{N}^*$. La statistique \widehat{T} est exhaustive si et seulement s'il existe une fonction $h : \mathbb{R}^n \rightarrow \mathbb{R}_+$ et une fonction $g_\theta(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}_+$, telle que l'on puisse écrire pour tout $(x_1, \dots, x_n) \in (\Omega')^n$:*

$$L_\theta(x_1, \dots, x_n) = g_\theta(\widehat{T}(x_1, \dots, x_n)) \cdot h(x_1, \dots, x_n) \quad \text{pour tout } \theta \in \Theta.$$

Lemme. Soit le modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$. Alors ce modèle est dominé si et seulement si il existe une sous-famille dénombrable $(\mathbb{P}_{\theta_i})_{i \in \mathbb{N}}$ telle que pour tout $A \in \mathcal{A}$, $\forall i \in \mathbb{N}$, $\mathbb{P}_{\theta_i}(A) = 0$ entraîne $\forall \theta$, $\mathbb{P}_\theta(A) = 0$. Toute mesure de probabilité de la forme $\mathbb{P}^* = \sum_{i \in \mathbb{N}} a_i \cdot \mathbb{P}_{\theta_i}$ avec $c_i > 0$ pour tout $i \in \mathbb{N}$ et $\sum_{i \in \mathbb{N}} c_i = 1$ domine le modèle.

Démonstration du lemme : \Leftarrow Il est bien clair que si une telle mesure P^* existe, le modèle est dominé.
 \Rightarrow Montrons maintenant que si le modèle est dominé par une mesure μ alors la famille $(P_{\theta_i})_{i \in \mathbb{N}}$ existe. En premier lieu, si μ est une mesure non finie mais σ -finie (par exemple la mesure de Lebesgue), alors μ_P définie par $\mu_P(A) = \sum_{i=1}^{\infty} \frac{1}{2^i} \frac{\mu(A \cap A_i)}{\mu(A_i)}$ pour tout $A \in \mathcal{A}$, est une mesure de probabilité équivalente à μ (avec $(A_i)_{i \in \mathbb{N}^*}$ une partition de $(\Omega')^n$ telle que $0 < \mu(A_i) < \infty$ pour tout $i \in \mathbb{N}^*$). On travaille donc désormais avec μ_P .

Pour $\theta \in \Theta$, soit B_θ le sous-ensemble de $(\Omega')^n \subset \mathbb{R}^n$ qui est le support de la densité de \mathbb{P}_θ par rapport à μ . Soit

$$\mathcal{C} = \left\{ \bigcup_{i \in I} B_{\theta_i}, I \subset \mathbb{N}, \theta_i \in \Theta \right\},$$

l'ensemble de toutes les unions dénombrables d'ensembles B_θ . On note $M = \sup_{C \in \mathcal{C}} \mu_P(C)$. Soit $(C_n)_{n \in \mathbb{N}}$ une suite d'ensembles de \mathcal{C} telle que la suite $(\mu_P(C_n))_n$ converge vers M (une telle suite existe forcément sinon M ne serait pas le supremum). Remarquons que chaque C_i étant une union dénombrable de B_{θ_k} , alors une suite (θ_n) de θ suffit pour engendrer la suite $(C_n)_{n \in \mathbb{N}}$. Si on pose :

$$D = \bigcup_{n \in \mathbb{N}} C_n = \bigcup_{k \in \mathbb{N}} B_{\theta_k},$$

alors $M = \mu_P(D)$ et pour tout $\theta \in \Theta$, $B_\theta \cup D \in \mathcal{C}$ et :

$$\mu_P(B_\theta \cup D) \leq M \leq \mu_P(B_\theta \cup D) = \mu_P(B_\theta \cap D^c) + \mu_P(D)$$

Donc pour tout $\theta \in \Theta$, $\mu_P(B_\theta \cap D^c) = 0$ soit $\forall \theta \in \Theta$, $\mathbb{P}_\theta(B_\theta \cap D^c) = 0$ puisque $\mathbb{P}_\theta \ll \mu_P$. En conséquence, pour tout $A \in \mathcal{A}'_n$, $A \subset B_\theta \cup B_\theta^c = (\Omega')^n$, soit :

$$\mathbb{P}_\theta(A \cap D^c) = 0, \quad \text{car par définition des } B_\theta, \mathbb{P}_\theta(B_\theta^c) = 0.$$

Si on suppose maintenant que $A \in \mathcal{A}'_n$ est tel que $\mathbb{P}_{\theta_k}(A) = 0$, avec la suite (θ_k) précédemment définie, alors $\mu_P(A \cap B_{\theta_k}) = 0$ par définition des B_θ et donc $\mu_P(A \cap D) = 0$ (par la propriété de σ -additivité d'une mesure). Comme $\mathbb{P}_\theta \ll \mu_P$, on en déduit que $\forall \theta \in \Theta$, $\mathbb{P}_\theta(A \cap D) = 0$ et donc $\mathbb{P}_\theta(A) = \mathbb{P}_\theta(A \cap D) + \mathbb{P}_\theta(A \cap D^c) = 0$. Ainsi, \mathbb{P}^* domine bien \mathbb{P}_θ pour tout $\theta \in \Theta$. \blacksquare

Démonstration du Théorème de factorisation de Neyman : Soit $\mathbb{P}^* = \sum_{i \in \mathbb{N}} a_i \cdot \mathbb{P}_{\theta_i}$ une mesure de probabilité dominante construite comme dans le lemme.

\Leftarrow Si $g_\theta(\widehat{T}(x)) \cdot h(x)$ avec $x \in (\Omega')^n$ est la densité de \mathbb{P}_θ par rapport à μ , alors $\sum_{i \in \mathbb{N}} a_i \cdot g_{\theta_i}(\widehat{T}(x)) \cdot h(x) = g_*(\widehat{T}(x)) \cdot h(x)$ est une densité de P^* par rapport à μ . Alors, comme $g_*(\widehat{T}(x)) \cdot h(x) > 0$ P^* -p.s., donc \mathbb{P}_θ -p.s., pour toute variable aléatoire S intégrable, pour tout $\theta \in \Theta$:

$$\begin{aligned} \mathbb{E}_\theta(S \cdot \mathbb{I}_B) &= \int_B S d\mathbb{P}_\theta, \quad \text{pour tout } B \in \sigma(\widehat{T}), \text{ tribu engendrée par } \widehat{T} \\ &= \int_B S(x) \cdot g_\theta(\widehat{T}(x)) \cdot h(x) d\mu(x) \\ &= \int_B S(x) \cdot \frac{g_\theta(\widehat{T}(x)) \cdot h(x)}{g_*(\widehat{T}(x)) \cdot h(x)} d\mathbb{P}^*(x) \\ &= \mathbb{E}_* \left(\mathbb{I}_B \cdot \frac{g_\theta(\widehat{T})}{g_*(\widehat{T})} \cdot S \right) \\ &= \mathbb{E}_* \left(\mathbb{I}_B \cdot \frac{g_\theta(\widehat{T})}{g_*(\widehat{T})} \cdot \mathbb{E}_*(S \mid \widehat{T}) \right) \quad (\text{d'après la définition de l'espérance conditionnelle}) \\ &= \mathbb{E}_\theta \left(\mathbb{I}_B \cdot \mathbb{E}_*(S \mid \widehat{T}) \right). \end{aligned}$$

En conséquence, d'après la définition de l'espérance conditionnelle dans $\mathbb{L}^1((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta)$, on a \mathbb{P}_θ -p.s., $\mathbb{E}_*(S | \hat{T}) = \mathbb{E}_\theta(S | \hat{T})$: la statistique \hat{T} est bien exhaustive.

\implies On suppose que \hat{T} est une statistique exhaustive pour le modèle. Donc pour toute statistique intégrable S , $\forall \theta$, $\mathbb{E}_*(S | \hat{T}) = \mathbb{E}_\theta(S | \hat{T})$. En conséquence, si on note $\phi(x, \theta) = \frac{d\mathbb{P}_\theta}{d\mathbb{P}_*}(x)$ la densité de \mathbb{P}_θ par rapport à \mathbb{P}_* ,

$$\begin{aligned} \mathbb{E}_\theta(S) &= \mathbb{E}_\theta(\mathbb{E}_*(S | \hat{T})), \quad (\text{car } \hat{T} \text{ est exhaustive et d'après les propriétés de l'espérance conditionnelle}) \\ &= \mathbb{E}_*(\phi(X, \theta) \cdot \mathbb{E}_*(S | \hat{T})), \quad \text{où } X \sim \mathbb{P}_* \\ &= \mathbb{E}_* \left[\mathbb{E}_*(\phi(X, \theta) \cdot \mathbb{E}_*(S | \hat{T}) | \hat{T}) \right], \quad (\text{d'après les propriétés de l'espérance conditionnelle}) \\ &= \mathbb{E}_* \left[\mathbb{E}_*(\phi(X, \theta) | \hat{T}) \cdot \mathbb{E}_*(S | \hat{T}) \right], \quad (\text{car } \mathbb{E}_*(S | \hat{T}) \text{ est une fonction de } \hat{T}) \\ &= \mathbb{E}_* \left[\mathbb{E}_*(S \cdot \mathbb{E}_*(\phi(X, \theta) | \hat{T}) | \hat{T}) \right] \\ &= \mathbb{E}_* \left[S \cdot \mathbb{E}_*(\phi(X, \theta) | \hat{T}) \right] \end{aligned}$$

Ainsi, la variable aléatoire $\mathbb{E}_*(\phi(X, \theta) | \hat{T})$, qui est une fonction de \hat{T} (qui est elle-même une fonction sur $(\Omega')^n$), est la densité de \mathbb{P}_θ par rapport à \mathbb{P}_* . Par suite, la vraisemblance, qui est la densité de \mathbb{P}_θ par rapport à μ , s'écrit :

$$L_\theta(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mu}(x_1, \dots, x_n) = \frac{d\mathbb{P}_\theta}{d\mathbb{P}_*}(x_1, \dots, x_n) \cdot \frac{d\mathbb{P}_*}{d\mu}(x_1, \dots, x_n) = \mathbb{E}_*(\phi(X, \theta) | \hat{T}) \cdot h(x_1, \dots, x_n),$$

avec h une fonction mesurable. ■

Exemple. Différentes statistiques exhaustives pour les modèles paramétriques de loi uniforme, de loi de Bernoulli, de loi gaussienne...

Propriété. On se place dans le cadre d'un modèle paramétrique dominé.

1. La statistique $\hat{T} = (X_1, \dots, X_n)$ est exhaustive.
2. Si \hat{T} est une statistique exhaustive et s'il existe une fonction borélienne h telle qu'une autre statistique \hat{U} vérifie $\hat{T} = h(\hat{U})$, alors \hat{U} est également exhaustive.

On vient de voir que l'on peut toujours trouver une statistique exhaustive (l'échantillon lui-même par exemple). Comme on aurait plutôt tendance à vouloir le "maximum d'information" dans une statistique exhaustive, lorsque le paramètre est dans \mathbb{R}^d , on aimerait savoir quelle dimension minimale peut avoir cette statistique. En particulier, si $d = 1$, peut-on toujours trouver une statistique exhaustive de taille 1 ? L'exemple suivant montre que ce n'est pas toujours le cas :

Exemple. Soit le modèle statistique $([0, \infty[^n, \mathcal{B}([0, \infty[^n), (\mathbb{P}_\theta)^{\otimes n}, \theta \in \mathbb{R}_+)$, où la densité de \mathbb{P}_θ par rapport à la mesure de Lebesgue est : $f_\theta(x) = \theta(e^{\theta^2} - 1) \cdot e^{-\theta \cdot x} \cdot \mathbb{1}_{x \in [0, \theta]}$. Alors les statistiques $\hat{T}_1 = \max(X_1, \dots, X_n)$ et $\hat{T}_2 = X_1 + \dots + X_n$ ne sont pas chacune exhaustive alors que $\hat{T} = (\hat{T}_1, \hat{T}_2)$ est exhaustive. On pourra même montrer que cette statistique est de taille minimale...

Définition. Une statistique exhaustive \hat{T} du modèle statistique paramétrique dominé avec \hat{T} est dite minimale si pour toute autre statistique exhaustive \hat{U} est telle qu'il existe une fonction borélienne h vérifiant :

$$\hat{T} = h(\hat{U}).$$

Proposition. Soit un modèle statistique paramétrique dominé et soit $L_\theta(x_1, \dots, x_n)$ sa vraisemblance. Alors \hat{T} est une statistique exhaustive minimale pour ce modèle lorsque $\forall (x_1, \dots, x_n) \in (\Omega')^n$ et $\forall (y_1, \dots, y_n) \in (\Omega')^n$,

$$\left(\theta \mapsto \frac{L_\theta(x_1, \dots, x_n)}{L_\theta(y_1, \dots, y_n)} \quad \text{ne dépend pas de } \theta \right) \iff \left(\hat{T}(x_1, \dots, x_n) = \hat{T}(y_1, \dots, y_n) \right) \quad (2)$$

Démonstration de la proposition : On suppose que (2) est vraie et on suppose (sans perte de généralité) que la vraisemblance est strictement positive. Soit $t \in \text{Im}(\widehat{T}((\Omega')^n))$. Notons $x^{(t)} \in \widehat{T}^{-1}(\{t\}) \subset (\Omega')^n$. Alors $\forall x \in \widehat{T}^{-1}(\{t\})$, $\widehat{T}(x) = \widehat{T}(x^{(t)})$ et donc d'après (2),

$$h(x) = \frac{L_\theta(x)}{L_\theta(x^{(t)})} \text{ est indépendant de } \theta.$$

Posons $g_\theta(t) = L_\theta(x^{(t)})$. Alors $L_\theta(x) = g_\theta(\widehat{T}(x)) \cdot h(x)$. Comme ceci est vrai pour tout $x \in (\Omega')^n$, la statistique \widehat{T} est bien exhaustive.

Supposons maintenant que \widehat{S} est une autre statistique exhaustive. Alors par le théorème de factorisation de Neyman, il existe deux fonctions $g_\theta^{(s)}$ et $h^{(s)}$ (ne dépendant pas de θ) telles que pour tout $x \in (\Omega')^n$, $L_\theta(x) = g_\theta^{(s)}(\widehat{S}(x)) \cdot h^{(s)}(x)$. Ainsi pour tout $x \in (\Omega')^n$ et $y \in (\Omega')^n$ tels que $\widehat{S}(x) = \widehat{S}(y)$, alors :

$$\frac{L_\theta(x)}{L_\theta(y)} = \frac{g_\theta^{(s)}(\widehat{S}(x)) \cdot h^{(s)}(x)}{g_\theta^{(s)}(\widehat{S}(y)) \cdot h^{(s)}(y)} = \frac{h^{(s)}(x)}{h^{(s)}(y)}, \text{ qui est indépendant de } \theta.$$

Mais d'après (2) ceci n'est possible que si $\widehat{T}(x) = \widehat{T}(y)$. Donc \widehat{T} est une fonction de \widehat{S} et la statistique \widehat{T} est donc minimale. ■

Qu'elle serait une sorte d'opposée de la notion de statistique exhaustive minimale ? Ce devrait être une statistique ne dépendant pas du paramètre, soit :

Définition. Une statistique \widehat{T} d'un modèle paramétrique est dite libre si sa loi ne dépend pas du paramètre.

Or, de façon assez surprenante il peut arriver qu'une statistique exhaustive minimale comprenne une statistique libre, qui intuitivement ne devrait pas être prise en compte pour donner toute l'information sur θ (soit par exemple la loi \mathbb{P}_θ discrète et équilibrée sur $\{\theta - 1, \theta, \theta + 1\}$; pour un échantillon de taille 2, la statistique $(X_{(2)} - X_{(1)}, X_1 + X_2)$ est exhaustive minimale, mais $X_{(2)} - X_{(1)}$ est libre). Aussi peut-on rajouter une autre caractérisation des statistiques exhaustives pour pouvoir atteindre une forme d'optimalité pour ces statistiques, qui serait qu'aucune fonctionnelle non constante de la statistique ne peut être libre. Cela peut également se traduire de la façon suivante :

Définition. Une statistique exhaustive \widehat{T} du modèle statistique paramétrique dominé avec \widehat{T} à valeur dans \mathbb{R}^d est dite complète si pour toute fonction borélienne $h : \mathbb{R}^d \rightarrow \mathbb{R}$ telle que $h(\widehat{T})$ soit intégrable, alors :

$$\forall \theta \in \Theta, \quad \mathbb{E}_\theta(h(\widehat{T})) = 0 \quad \implies \quad h(\widehat{T}) = 0.$$

Propriété. Soit un modèle statistique paramétrique dominé.

1. si \widehat{T} est une statistique exhaustive complète alors pour toute fonction borélienne h bijective $h(\widehat{T})$ est une statistique exhaustive complète.
2. si \widehat{T} est une statistique exhaustive complète alors \widehat{T} est une statistique exhaustive minimale.
3. (Théorème de Basu) si \widehat{T} est une statistique exhaustive complète alors \widehat{T} est indépendante de toute statistique libre sur le modèle.

Démonstration de la propriété : 3. Théorème de Basu. Soit \widehat{S} une statistique libre pour le modèle et soit f une fonction telle que $\mathbb{E}_\theta(f(\widehat{S}))$ existe. Comme \widehat{S} est libre, on peut noter $e(f) = \mathbb{E}_\theta(f(\widehat{S}))$ une application linéaire ne dépendant pas de θ . Par suite, la statistique $\mathbb{E}_\mu(f(\widehat{S}) | \widehat{T}) - e(f)$ est une fonction de \widehat{T} mesurable telle que $\mathbb{E}_\theta(\mathbb{E}_\mu(f(\widehat{S}) | \widehat{T}) - e(f)) = 0$ pour tout $\theta \in \Theta$. Comme on a supposé que \widehat{T} est exhaustive complète, alors $\mathbb{E}_\mu(f(\widehat{S}) | \widehat{T}) = e(f)$ presque-sûrement : les statistiques \widehat{S} et \widehat{T} sont indépendantes. ■

Définition. On suppose un modèle paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta \subset \mathbb{R}^p)$ dominé par une mesure μ . Si, pour tout $(x_1, \dots, x_n) \in (\Omega')^n$ et $\theta \in \Theta$, la vraisemblance de ce modèle par rapport à μ peut s'écrire sous la forme :

$$L_\theta(x_1, \dots, x_n) = \exp \left(\beta(\theta) + b(x_1, \dots, x_n) + \sum_{j=1}^p a_j(x_1, \dots, x_n) \cdot \alpha_j(\theta) \right), \quad (3)$$

avec les fonctions $a_j : (\Omega')^n \rightarrow \mathbb{R}$, $b : (\Omega')^n \rightarrow \mathbb{R}$, $\alpha_j : \Theta \subset \mathbb{R}^p \rightarrow \mathbb{R}$, et $\beta : \Theta \rightarrow \mathbb{R}$, alors on dit que le modèle est exponentiel (ou qu'il appartient à la famille exponentielle).

Exemple. Appartiennent à la famille exponentielle les lois :

- Loi discrètes : Lois de Bernoulli, binomiales, de Poisson,...
- Loi "continues" : Lois normales, exponentielles, gamma, du chi-deux,...

Remarque. Si (X_1, \dots, X_n) est un n -échantillon d'un modèle exponentiel (avec θ fixé) alors l'ensemble des valeurs prises par (X_1, \dots, X_n) ne dépend pas du paramètre θ .

Propriété. Soit un modèle exponentiel. Si pour tout $\theta \in \Theta$ on note $\alpha(\theta) = (\alpha_1(\theta), \dots, \alpha_p(\theta))$ et si l'ensemble $\alpha(\Theta)$ est d'intérieur non vide, alors $\widehat{T}(x_1, \dots, x_n) = (a_1(x_1, \dots, x_n), \dots, a_p(x_1, \dots, x_n))$ est une statistique exhaustive minimale et complète.

Démonstration de la propriété : Soit $g : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $\mathbb{E}_\theta(g(\widehat{T})) = 0$. Or, $\forall \theta \in \Theta$,

$$\mathbb{E}_\theta(g(\widehat{T})) = \int_{(\Omega')^n} g(\widehat{T}(x)) \cdot \exp(\beta(\theta) + b(x) + \langle \widehat{T}(x), \alpha(\theta) \rangle) d\mu(x),$$

où $\langle \cdot, \cdot \rangle$ désigne le produit scalaire. En considérant la mesure ν de densité $\exp(b(x))$ par rapport à μ , on obtient :

$$\begin{aligned} \mathbb{E}_\theta(g(\widehat{T})) = 0 &\implies \int_{(\Omega')^n} g(\widehat{T}(x)) \cdot \exp(\langle \widehat{T}(x), \alpha(\theta) \rangle) d\nu(x) = 0 \\ &\implies \int_{\widehat{T}((\Omega')^n)} g(y) \cdot \exp(\langle y, \alpha(\theta) \rangle) d\nu_{\widehat{T}}(y) = 0 \end{aligned}$$

pour tout $\theta \in \Theta$, en ayant noté $\nu_{\widehat{T}}$ la mesure image de ν par \widehat{T} et avec $\widehat{T}((\Omega')^n) \in \mathbb{R}^p$. Si on note g^+ et g^- les parties positives et négatives de g (donc $g = g^+ - g^-$), et π^+ et π^- les mesures de densités g^+ et g^- par rapport à $\nu_{\widehat{T}}$, alors, pour tout $\theta \in \Theta$:

$$\int_{\widehat{T}((\Omega')^n)} \exp(\langle y, \alpha(\theta) \rangle) d\pi^+(y) = \int_{\widehat{T}((\Omega')^n)} \exp(\langle y, \alpha(\theta) \rangle) d\pi^-(y).$$

En conséquence sur Θ , donc sur une partie d'intérieure non vide, les mesures π^+ et π^- ont des transformées de Laplace égales : ces deux mesures sont donc égales et donc $g^+ = g^-$ $\nu_{\widehat{T}}$ -presque partout (ce qui revient à $g = 0$). A partir des expressions des différentes mesures, on montre que $g = 0$, $\widehat{T}(\mathbb{P}_\theta)$ -presque partout. ■

3.3 Information de Fisher

Pour mesurer l'information fournit par un modèle paramétrique dominé (ou une statistique sur ce modèle) au sujet d'un paramètre, une idée naturelle serait de mesurer comment varie localement la mesure de probabilité, ou encore sa vraisemblance. Les fluctuations moyennes de cette vraisemblance serait donc un bon indicateur : pour ce faire on considèrera, lorsqu'il existe $\text{grad}_\theta(L_\theta(X_1, \dots, X_n))$, et on s'intéressera à la matrice de covariance de $\text{grad}_\theta(L_\theta(X_1, \dots, X_n))$, dont on peut montrer qu'elle ne dépend pas du choix de la mesure dominante choisie. Précisons d'abord la notion de modèle régulier qui nous permettra de définir cette quantité d'information.

Définition. Dans le cadre d'un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ , on dira que ce modèle est régulier lorsque :

1. Θ est un ouvert de \mathbb{R}^d ;
2. la vraisemblance $L_\theta(\cdot)$ vérifie $\forall (x_1, \dots, x_n) \in (\Omega')^n, \forall \theta \in \Theta, L_\theta(x_1, \dots, x_n) > 0$;
3. $\forall (x_1, \dots, x_n) \in (\Omega')^n$, la fonction $\theta \in \Theta \mapsto \log(L_\theta(\cdot))$ est différentiable sur Θ par rapport à θ , et son gradient appartient à $\mathbb{L}^2((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta) \forall \theta \in \Theta$;

4. $\forall \theta \in \Theta$, pour toute fonction $h : \mathbb{R}^n \rightarrow \mathbb{R}$ appartenant à $\mathbb{L}^1((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta)$, alors :

$$\frac{\partial}{\partial \theta} \int_{(\Omega')^n} h(x) \cdot L_\theta(x) d\mu(x) = \int_{(\Omega')^n} h(x) \cdot \frac{\partial}{\partial \theta} L_\theta(x) d\mu(x). \quad (4)$$

Conséquence. Pour un modèle régulier, $\mathbb{E}_\theta(\text{grad}_\theta(\log L_\theta(\cdot))) = 0$.

Démonstration : On a $\mathbb{E}_\mu(L_\theta(\cdot)) = 1$ donc $\mathbb{E}_\mu(\text{grad}_\theta L_\theta(\cdot)) = 0$. Par conséquent, $\mathbb{E}_\mu\left(\frac{\text{grad}_\theta(L_\theta(\cdot))}{L_\theta(\cdot)}\right) = 0$, soit $\mathbb{E}_\theta(\text{grad}_\theta(\log L_\theta(\cdot))) = 0$. ■

Définition. Pour un modèle statistique paramétrique dominé régulier, on appelle information de Fisher, la matrice :

$$I_n(\theta) = \left[\mathbb{E}_\theta \left(\frac{\partial(\log L_\theta(X_1, \dots, X_N))}{\partial \theta_i} \times \frac{\partial(\log L_\theta(X_1, \dots, X_N))}{\partial \theta_j} \right) \right]_{1 \leq i, j \leq p}.$$

Propriété. Pour un modèle statistique paramétrique dominé régulier, et si $\forall (x_1, \dots, x_n) \in (\Omega')^n$, la fonction $\theta \in \Theta \mapsto \log(L_\theta(\cdot))$ est $\mathcal{C}^2(\Theta)$, alors :

$$I_n(\theta) = - \left[\mathbb{E}_\theta \left(\frac{\partial^2(\log L_\theta(X_1, \dots, X_N))}{\partial \theta_i \cdot \partial \theta_j} \right) \right]_{1 \leq i, j \leq p}.$$

Définition. L'information de Fisher $I_n^{\widehat{T}}(\theta)$ associée à une statistique \widehat{T} , si elle existe, est la matrice de Fisher de la vraisemblance de \widehat{T} (déterminée à partir de la vraisemblance de \widehat{T}).

Propriété. Pour un modèle régulier, \widehat{T} est une statistique libre si et seulement si $I_n^{\widehat{T}}(\theta) = 0$.

Démonstration : \implies Si \widehat{T} est libre alors sa loi ne dépend pas de θ donc le gradient du logarithme de sa vraisemblance est nul; l'information de Fisher associée à \widehat{T} est nulle.

\impliedby Si $I_n^{\widehat{T}}(\theta) = 0$, donc la statistique $\text{grad}_\theta(\log L_\theta^{\widehat{T}}(\widehat{T}))$ est centrée et de matrice de covariance nulle. Ainsi, pour tout $\theta \in \Theta$, il existe un ensemble N_θ de mesure 1 pour la mesure de probabilité associée à \widehat{T} (donc, d'après la première hypothèse d'un modèle régulier, tel que $\mu(N_\theta) = 1$) et tel que pour tout $t \in N_\theta$, $\text{grad}_\theta(\log L_\theta^{\widehat{T}}(t)) = 0$. Pour montrer que $\text{grad}_\theta(\log L_\theta^{\widehat{T}}(t)) = 0$ est bien une variable aléatoire nulle μ -p.s., et donc que $\log L_\theta^{\widehat{T}}(\cdot)$ est une fonction constante en θ , il nous faut montrer que finalement les θ ne dépendent pas de θ . Soit $\Theta^{(d)} = \{\theta_i^{(d)}\}_{i \in \mathbb{N}}$ un sous-ensemble dénombrable de Θ , dense dans Θ . Comme $\Theta^{(d)}$ est dénombrable, il est clair que $N = \bigcup_{i \in \mathbb{N}} N_{\theta_i^{(d)}}$ est tel que $\mu(N) = 1$. De plus, pour tout $\theta \in \Theta$, il existe une sous-suite $(\theta_{\phi(n)}^{(d)})_n$ de $\Theta^{(d)}$ convergeant vers θ et telle que pour tout $t \in N$, pour tout $n \in \mathbb{N}$, $\text{grad}_{\theta_{\phi(n)}^{(d)}}(\log L_{\theta_{\phi(n)}^{(d)}}^{\widehat{T}}(t)) = 0$. Comme une telle fonction de $\theta_{\phi(n)}^{(d)}$ est continue, cette propriété passe à la limite, et donc pour tout $t \in N$, $\forall \theta \in \Theta$, $\text{grad}_\theta(\log L_\theta^{\widehat{T}}(t)) = 0$. Comme N ne dépend pas de θ , alors la fonction $\theta \rightarrow \log L_\theta^{\widehat{T}}(\cdot)$ est une constante ne dépendant pas de θ , μ -p.s. : la statistique \widehat{T} est bien libre. ■

Propriété. Pour un modèle régulier, si \widehat{T} est une statistique exhaustive : $I_n^{\widehat{T}}(\theta) = I_n(\theta)$ pour tout $\theta \in \Theta$.

Démonstration : Comme \widehat{T} est une statistique exhaustive, on peut écrire d'après la démonstration du Théorème de factorisation de Neyman que pour tout $(x_1, \dots, x_n) \in (\Omega')^n$ et tout $\theta \in \Theta$:

$$\frac{d\mathbb{P}_\theta}{d\mathbb{P}^*}(x_1, \dots, x_n) = g_\theta(\widehat{T}(x_1, \dots, x_n)).$$

On peut réécrire cela pour la densité de \widehat{T} sous la forme : $\frac{d\mathbb{P}_\theta^{\widehat{T}}}{d\mathbb{P}^{*\widehat{T}}}(t) = g_\theta(t)$, pour tout $t \in \widehat{T}((\Omega')^n)$ et tout $\theta \in \Theta$. En conséquence, pour tout $\theta \in \Theta$,

$$I(\theta) = \left[\mathbb{E}_\theta \left(\frac{\partial(\log L_\theta(X_1, \dots, X_N))}{\partial \theta_i} \times \frac{\partial(\log L_\theta(X_1, \dots, X_N))}{\partial \theta_j} \right) \right]_{1 \leq i, j \leq p}$$

$$\begin{aligned}
&= \left[\int_{(\Omega')^n} \left(\frac{\partial(\log L_\theta(x))}{\partial\theta_i} \times \frac{\partial(\log L_\theta(x))}{\partial\theta_j} \right) d\mathbb{P}_\theta(x) \right]_{1 \leq i, j \leq p} \\
&= \left[\int_{(\Omega')^n} \left(\frac{\partial(\log g_\theta(\widehat{T}(x)))}{\partial\theta_i} \times \frac{\partial(\log g_\theta(\widehat{T}(x)))}{\partial\theta_j} \right) g_\theta(\widehat{T}(x)) d\mathbb{P}^*(x) \right]_{1 \leq i, j \leq p} \quad \text{car } \log L_\theta(x) = \log g_\theta(\widehat{T}(x)) + \log h(x) \\
&= \left[\int_{\widehat{T}(\Omega')^n} \left(\frac{\partial(\log g_\theta(t))}{\partial\theta_i} \times \frac{\partial(\log g_\theta(t))}{\partial\theta_j} \right) g_\theta(t) d\mathbb{P}^{\widehat{T}}(x) \right]_{1 \leq i, j \leq p} \quad \text{d'après le théorème du transport} \\
&= \left[\int_{\widehat{T}(\Omega')^n} \left(\frac{\partial(\log g_\theta(t))}{\partial\theta_i} \times \frac{\partial(\log g_\theta(t))}{\partial\theta_j} \right) d\mathbb{P}_\theta(t) \right]_{1 \leq i, j \leq p} \\
&= I_n^{\widehat{T}}(\theta). \quad \blacksquare
\end{aligned}$$

Remarque. En rajoutant certaines hypothèses de continuité sur la vraisemblance de \widehat{T} , on peut montrer que la réciproque est également vraie, et donc que $I_n^{\widehat{T}}(\theta) = 0$ si et seulement si la statistique \widehat{T} est exhaustive.

Ainsi, on retrouve à l'aide de la notion d'information de Fisher les "intuitions" qui nous avaient guidées dans la section précédentes. Voyons maintenant les applications de la notion d'exhaustivité à l'estimation paramétrique.

3.4 Application à l'estimation paramétrique

On se place dans le cadre d'un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ . Par ailleurs, on suppose que Θ est un ouvert.

Définition. • Soit $g : \Theta \rightarrow \Theta'$, où $\Theta' \subset \mathbb{R}^{p'}$ avec $p' \in \mathbb{N}^*$, une fonction mesurable. On appelle estimateur de la fonction g du paramètre, donc de $g(\theta)$, une statistique \widehat{T} à valeurs dans $\mathbb{R}^{p'}$. En particulier, un estimateur du paramètre θ est une statistique à valeurs dans \mathbb{R}^p . Une estimation de $g(\theta)$ est une réalisation de \widehat{T} .

- On appelle biais d'un estimateur \widehat{T} de $g(\theta)$ le vecteur constant de $\mathbb{R}^{p'}$, $B(\theta) = \mathbb{E}_\theta(\widehat{T}) - g(\theta)$. On dira que l'estimateur est sans biais si $B(\theta) = 0$ pour tout $\theta \in \Theta$.
- On appelle risque quadratique de l'estimateur \widehat{T} de $g(\theta)$ le réel positif $R(\theta) = \mathbb{E}_\theta(\|\widehat{T} - g(\theta)\|^2)$, où $\|\cdot\|$ désigne usuellement la norme euclidienne (mais peut être une autre fonctionnelle positive et convexe). Si l'estimateur est sans biais alors, $R(\theta) = \text{Trace}(\text{cov}(\widehat{T}))$.

Pour pouvoir parler du comportement asymptotique d'une statistique, on va devoir se placer dans un "gros" modèle, dans lequel un échantillon est une suite de v.a. En quelque sorte, ce gros modèle pourra s'écrire $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, \mathbb{P}_\theta^{\mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$ (la dimension du paramètre reste constante). Pour un n fixé, une statistique \widehat{T}_n sera d'abord une projection du "gros" modèle sur le modèle de taille n , puis une statistique "normale". On devra donc parler d'une suite d'estimateurs $(\widehat{T}_n)_n$.

Définition. Pour un modèle statistique paramétrique $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, \mathbb{P}_\theta^{\mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, et pour $(\widehat{T}_n)_n$ une suite d'estimateurs de $g(\theta)$:

- Si $\lim_{n \rightarrow \infty} B_n(\theta) = 0$, on dit que l'estimateur est asymptotiquement sans biais.
- On dit que $(\widehat{T}_n)_n$ est convergent lorsque $\widehat{T}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} g(\theta)$.
- S'il existe (a_n) une suite de réels positifs tels que $a_n(\widehat{T}_n - g(\theta)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Z_\theta$, où Z_θ est une loi centrée non nulle (ne dépend pas de n), on dit $(\widehat{T}_n)_n$ converge vers $g(\theta)$ à la vitesse a_n .

A priori, être sans biais n'est pas un bon critère pour garantir une certaine optimalité de la convergence d'un estimateur. On préférera plutôt discriminer entre de potentiels estimateurs à l'aide d'un critère portant sur le risque quadratique ou sur la matrice de variance-covariance. Cependant, il n'existe pas de résultats généraux pour trouver un "meilleur" estimateur en ce sens. Pour en obtenir, on devra se limiter à une certaine classe d'estimateurs, celle des estimateurs sans biais.

Définition. Soit un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, et soit \widehat{T} un estimateur sans biais de $g(\theta)$. On dit que \widehat{T} est de variance uniformément minimum parmi les estimateurs sans biais de $g(\theta)$ lorsque pour tout estimateur sans biais de $g(\theta)$, on a $\forall \theta \in \Theta, \text{cov}(\widehat{T}) \leq \text{cov}(\widehat{S})$ (au sens où $\text{cov}(\widehat{T}) - \text{cov}(\widehat{S})$ est une matrice positive).

Propriété. Si \widehat{T} est un estimateur de variance uniformément minimum parmi les estimateurs sans biais, alors il est unique \mathbb{P}_θ -p.s.

Démonstration : Soit \widehat{S} un autre estimateur que l'on suppose également de variance uniformément minimum parmi les estimateurs sans biais. Montrons d'abord que $E_\theta((\widehat{T} - \widehat{S}) \cdot {}^t\widehat{T}) = 0$. En effet, si $\alpha \in \mathbb{R}$, comme \widehat{T} est de variance minimum, en utilisant des inégalités sur les matrices symétriques :

$$\begin{aligned} \text{cov}(\widehat{T}) &\leq \text{cov}(\widehat{T} + \alpha(\widehat{T} - \widehat{S})) \\ &\leq \text{cov}(\widehat{T}) + \alpha^2 \text{cov}(\widehat{T} - \widehat{S}) + 2\alpha \cdot \mathbb{E}_\theta(\widehat{T} \cdot {}^t\widehat{S}) \\ \implies 0 &\leq \alpha \cdot (\alpha \cdot \text{cov}(\widehat{S}) + 2\mathbb{E}_\theta(\widehat{T} \cdot {}^t(\widehat{T} - \widehat{S}))) \quad \text{pour tout } \alpha \in \mathbb{R}. \end{aligned}$$

Comme $\text{cov}(\widehat{T} - \widehat{S})$ est une matrice positive, la seule possibilité pour avoir la dernière inégalité est que : $\mathbb{E}_\theta(\widehat{T} \cdot {}^t(\widehat{T} - \widehat{S})) = 0$. Par suite, comme $\text{cov}(\widehat{T} - \widehat{S}) = \mathbb{E}_\theta((\widehat{T} - \widehat{S}) \cdot {}^t(\widehat{T} - \widehat{S})) = \mathbb{E}_\theta(\widehat{T} \cdot {}^t(\widehat{T} - \widehat{S})) - \mathbb{E}_\theta(\widehat{S} \cdot {}^t(\widehat{T} - \widehat{S}))$, et que l'on a supposé \widehat{T} et \widehat{S} de variance minimum, $\text{cov}(\widehat{T} - \widehat{S}) = 0$. Donc $\widehat{T} = \widehat{S}$ sur un ensemble de \mathbb{P}_θ -mesure égale à 1. ■

Théorème (Rao-Blackwell). Si \widehat{T} est un estimateur sans biais de $g(\theta)$ et si \widehat{S} est une statistique exhaustive, alors $\widehat{R} = \mathbb{E}_\theta(\widehat{T} \mid \widehat{S})$, qui ne dépend pas de θ car \widehat{S} est exhaustive, est un estimateur sans biais de $g(\theta)$ de matrice de covariance inférieure ou égale à celle de \widehat{T} .

Démonstration : il est clair que $\mathbb{E}_\theta(\widehat{R}) = \mathbb{E}_\theta(\widehat{T}) = g(\theta)$. De plus, pour tout $u \in \mathbb{R}^{p'}$ (avec $g : \mathbb{R}^p \rightarrow \mathbb{R}^{p'}$),

$$\begin{aligned} \text{cov}({}^t u \cdot \widehat{T}) &= \mathbb{E}_\theta \left[({}^t u \cdot (\widehat{T} - g(\theta)))^2 \right] \\ &= \mathbb{E}_\theta \left(\mathbb{E}_\theta \left[({}^t u \cdot (\widehat{T} - g(\theta)))^2 \mid \widehat{S} \right] \right) \\ &\geq \mathbb{E}_\theta \left(\mathbb{E}_\theta \left[({}^t u \cdot (\widehat{T} - g(\theta)) \mid \widehat{S})^2 \right] \right) \quad \text{d'après l'inégalité de Jensen,} \\ &\geq \text{cov}({}^t u \cdot \widehat{R}). \end{aligned}$$

Cela revient bien à écrire que $\text{cov}(\widehat{T}) \geq \text{cov}(\widehat{R})$. ■

Théorème (Lehmann-Scheffé). Si \widehat{T} est un estimateur sans biais de $g(\theta)$ et si \widehat{S} est une statistique exhaustive et complète, alors l'unique estimateur de $g(\theta)$ sans biais uniformément de variance minimale est $\widehat{R} = \mathbb{E}_\theta(\widehat{T} \mid \widehat{S})$ (c'est-à-dire que \widehat{R} est une fonction de \widehat{S}).

Démonstration : Soit \widehat{T}' un autre estimateur sans biais de $g(\theta)$. Si $\widehat{R}' = \mathbb{E}_\theta(\widehat{T}' \mid \widehat{S})$, on sait que $\text{cov}(\widehat{T}') \geq \text{cov}(\widehat{R}')$ d'après le Théorème de Rao-Blackwell. Or $\mathbb{E}_\theta(\widehat{R} - \widehat{R}') = 0$ pour tout $\theta \in \Theta$ car les deux estimateurs sont sans biais. De plus comme \widehat{R} et \widehat{R}' sont des fonctions de \widehat{S} , $\widehat{R} - \widehat{R}'$ l'est aussi, et du fait que \widehat{S} est une statistique exhaustive et complète, alors pour tout $\theta \in \Theta$, $\widehat{R} = \widehat{R}'$, \mathbb{P}_θ -p.s. Par conséquent, pour tout $\theta \in \Theta$, $\text{cov}(\widehat{R}') = \text{cov}(\widehat{R})$ et donc $\text{cov}(\widehat{R}) \leq \text{cov}(\widehat{T}')$: \widehat{R} est bien l'estimateur sans biais de variance uniformément minimale. ■

Retenons donc de tout ceci que l'estimateur sans biais de $g(\theta)$ et de variance uniformément minimale est une unique fonction d'une statistique exhaustive et complète, lorsqu'une telle statistique existe. On aimerait maintenant connaître un peu mieux la covariance d'un tel estimateur.

Théorème (Inégalité de Cramer-Rao). Soit un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$ dominé et régulier, et soit \widehat{T} un estimateur sans biais de $g(\theta)$, tel que $\mathbb{E}_\theta \|\widehat{T}\|^2 < +\infty$. Si on suppose que l'information de Fisher est une matrice définie positive, alors, en notant $\frac{\partial g}{\partial \theta}(\theta)$ la matrice jacobienne de g , pour tout $\theta \in \Theta$:

$$\text{cov}(\widehat{T}) \geq \frac{\partial g}{\partial \theta}(\theta) \cdot (I_n(\theta))^{-1} \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) \quad (\text{au sens des matrices symétriques}).$$

En particulier, si \widehat{T} est un estimateur sans biais de θ , alors :

$$\text{cov}(\widehat{T}) \geq (I_n(\theta))^{-1} \quad (\text{au sens des matrices symétriques}).$$

Démonstration : Soit $Z_\theta(x) = \text{grad}(\log L_\theta(x))$ où $x \in (\Omega')^n$ suit \mathbb{P}_θ . On sait que comme le modèle est régulier, $\mathbb{E}_\theta(Z_\theta) = 0$ pour tout $\theta \in \Theta$ et donc :

$$\text{cov}(Z_\theta) = I(\theta) \quad \text{pour tout } \theta \in \Theta.$$

De plus, \widehat{T} est un estimateur sans biais de $g(\theta)$ donc pour tout $\theta \in \Theta$:

$$\begin{aligned} \mathbb{E}_\theta(\widehat{T}) = g(\theta) &\implies \int_{(\Omega')^n} \widehat{T}(x) \cdot \frac{\partial L_\theta}{\partial \theta}(x) d\mu(x) = \frac{\partial g}{\partial \theta}(\theta) \quad (\text{en dérivant}) \\ &\implies \int_{(\Omega')^n} \widehat{T}(x) \cdot \frac{\partial L_\theta}{\partial \theta}(x) \cdot (L_\theta(x))^{-1} d\mathbb{P}_\theta(x) = \frac{\partial g}{\partial \theta}(\theta) \\ &\implies \mathbb{E}_\theta(\widehat{T} \cdot {}^t Z_\theta) = \frac{\partial g}{\partial \theta}(\theta). \end{aligned}$$

Ainsi, d'après ce qui précède,

$$\begin{aligned} \text{cov}_\theta(\widehat{T} - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta) &= \text{cov}_\theta(\widehat{T}) - 2 \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) + \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) \\ &= \text{cov}_\theta(\widehat{T}) - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta). \end{aligned}$$

En conséquence, comme $\text{cov}_\theta(\widehat{T} - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta)$ est une matrice positive, l'inégalité de Cramer-Rao est prouvée. \blacksquare

Corollaire. Deux cas particuliers méritent attention :

- Si le modèle est de la forme $((\Omega')^n, \mathcal{A}'_n, (f_\theta \cdot d\mu)^{\otimes n}, \theta \in \Theta)$, alors $I_n(\theta) = n \cdot I_1(\theta)$, où $I_1(\theta)$ est la matrice d'information de Fisher d'une seule variable aléatoire X distribuée suivant $f_\theta \cdot d\mu$ et l'Inégalité de Cramer-Rao devient donc :

$$\text{cov}(\widehat{T}) \geq \frac{1}{n} \cdot \left(\frac{\partial g}{\partial \theta}(\theta) \cdot (I_1(\theta))^{-1} \cdot {}^t \frac{\partial g}{\partial \theta}(\theta) \right) \quad (\text{au sens des matrices symétriques}).$$

On voit donc que pour un échantillon de variables indépendantes et identiquement distribuées, si la vraisemblance est régulière, alors la vitesse de convergence de tout estimateur sans biais est au mieux en \sqrt{n} .

- Si le modèle n'est pas régulier, mais que sous la probabilité \mathbb{P}_θ , la matrice d'information de Fisher existe et est inversible, et surtout si la propriété (4) est vérifiée, alors l'Inégalité de Cramer-Rao est vérifiée. **Cela exclut cependant les modèles dont le support de \mathbb{P}_θ dépend de θ , comme par exemple le simple modèle de v.a.i.i.d. de loi $\mathcal{U}(]0, \theta[)$, avec $\theta > 0$.**

Définition. Si un estimateur sans biais atteint (respectivement asymptotiquement) la borne de Cramer-Rao (qui ne dépend pas de l'estimateur), on dit qu'il est (resp. asymptotiquement) efficace.

Remarque. Un estimateur peut être sans biais, de variance minimale, mais ne pas atteindre la borne de Cramer-Rao, donc ne pas être efficace. De la même manière, il peut exister des estimateurs biaisés atteignant la borne de Cramer-Rao.

Nous allons voir que les modèles exponentiels jouent un rôle central pour l'estimation paramétrique puisque sous certaines conditions ils sont les seuls pour lesquels on aura une estimation sans biais efficace.

Théorème. Soit un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, avec $\Theta \subset \mathbb{R}^p$, dominé et régulier. Soit $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ de classe \mathcal{C}^1 sur Θ telle que la matrice carrée de taille p , $\frac{\partial g}{\partial \theta}(\theta)$ soit de rang p

pour tout $\theta \in \Theta$. Alors $\hat{T} = {}^t(\hat{T}_1, \dots, \hat{T}_d)$ est un estimateur sans biais de $g(\theta)$ atteignant la borne de Cramer-Rao si et seulement si le modèle est exponentiel et plus précisément s'il existe des fonctions $a : (\Omega')^n \rightarrow \mathbb{R}$, $\beta : \Theta \rightarrow \mathbb{R}$ et $\alpha_j : \Theta \rightarrow \mathbb{R}$ ($1 \leq j \leq p$), telles que pour tout $\theta \in \Theta$, $g(\theta) = - \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p}^{-1} \cdot \frac{\partial \beta}{\partial \theta}(\theta)$ et

$$L_\theta(x_1, \dots, x_n) = \exp \left(\beta(\theta) + b(x_1, \dots, x_n) + \sum_{j=1}^d T_j(x_1, \dots, x_n) \cdot \alpha_j(\theta) \right).$$

Démonstration : \Leftarrow On suppose donc le modèle exponentiel décrit dans le théorème. Si on dérive par rapport à θ un tel modèle, on obtient que pour μ -presque tout $x \in (\Omega')^n$:

$$\frac{\partial}{\partial \theta}(\log L_\theta(x)) = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \cdot \hat{T} + \frac{\partial \beta}{\partial \theta}(\theta), \quad \text{pour tout } \theta \in \Theta. \quad (5)$$

En conséquence, comme $I(\theta) = \mathbb{E}_\theta \left(\left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \cdot {}^t \left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \right)$, on en déduit que :

$$I(\theta) = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq d} \cdot \text{cov}_\theta(\hat{T}) \cdot {}^t \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \implies \text{cov}_\theta(\hat{T}) = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p}^{-1} \cdot I(\theta) \cdot {}^t \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p}^{-1}$$

Par ailleurs, comme \hat{T} est un estimateur sans biais de $g(\theta)$ d'après la preuve de l'Inégalité de Cramer-Rao,

$$\mathbb{E}_\theta \left(\hat{T}(\cdot) \cdot {}^t \left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \right) = \frac{\partial g}{\partial \theta}(\theta)$$

et en utilisant (5) que l'on multiplie par $\left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right)$, on obtient :

$$\mathbb{E}_\theta \left(\left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \cdot {}^t \left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \right) = \mathbb{E}_\theta \left(\left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \cdot \hat{T} \cdot {}^t \left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \right) + \mathbb{E}_\theta \left(\frac{\partial \beta}{\partial \theta}(\theta) \cdot {}^t \left(\frac{\partial}{\partial \theta}(\log L_\theta(\cdot)) \right) \right),$$

et donc $I(\theta) = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \cdot \frac{\partial g}{\partial \theta}(\theta)$. A l'aide de cette égalité, et en reprenant le calcul précédent, on en arrive à ce que :

$$\text{cov}_\theta(\hat{T}) = \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot {}^t \frac{\partial g}{\partial \theta}(\theta),$$

donc \hat{T} atteint bien la borne de Cramer-Rao. De plus, grâce à (5),

$$\mathbb{E}_\theta \left(\frac{\partial}{\partial \theta}(\log L_\theta(x)) \right) = \mathbb{E}_\theta \left(\left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \cdot \hat{T} + \frac{\partial \beta}{\partial \theta}(\theta) \right)$$

$$\text{soit} \quad 0 = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p} \cdot g(\theta) + \frac{\partial \beta}{\partial \theta}(\theta)$$

$$\text{et donc} \quad g(\theta) = - \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta) \right)_{1 \leq i, j \leq p}^{-1} \cdot \frac{\partial \beta}{\partial \theta}(\theta).$$

\implies D'après la preuve de l'Inégalité de Cramer-Rao, si \hat{T} est un estimateur sans biais de $g(\theta)$ atteignant la borne de Cramer-Rao, alors

$$\text{cov}_\theta(\hat{T} - \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta) = 0.$$

Ainsi, pour tout $\theta \in \Theta$, il existe un ensemble $N_\theta \subset (\Omega')^n$ tel que $\mathbb{P}_\theta(N_\theta) = 1$ et tel que pour tout $x \in N_\theta$, $\hat{T}(x) - g(\theta) = \frac{\partial g}{\partial \theta}(\theta) \cdot I^{-1}(\theta) \cdot Z_\theta(x)$. Par le même procédé que celui de la preuve de la nullité de l'information de Fisher pour une statistique libre, on peut déterminer un ensemble N ne dépendant pas de θ , tel que cette propriété soit également vraie, avec $\mu(N) = 1$, ce qui revient à écrire que $\forall x \in N$,

$$I(\theta) \cdot \left(\frac{\partial g}{\partial \theta}(\theta) \right)^{-1} \cdot (\hat{T}(x) - g(\theta)) = \frac{\partial}{\partial \theta}(\log L_\theta(x)), \quad \text{pour tout } \theta \in \Theta.$$

Alors en intégrant par rapport à θ , et en notant

$$\begin{cases} \alpha(\theta) \text{ le vecteur colonne "intégrant"} & I(\theta) \cdot \left(\frac{\partial g}{\partial \theta}(\theta)\right)^{-1} \\ \beta(\theta) \text{ la fonction "intégrant"} & -I(\theta) \cdot \left(\frac{\partial g}{\partial \theta}(\theta)\right)^{-1} \cdot g(\theta) \\ b(x) \text{ une fonction ne dépendant pas de } \theta & \end{cases}$$

on a $\log L_\theta(x) = \alpha(\theta) \cdot \widehat{T}(x) + \beta(\theta) + b(x)$, d'où l'écriture de la vraisemblance sous forme d'un modèle exponentiel, et on retrouve l'expression de $g(\theta)$ par le même raisonnement que plus haut. ■

Corollaire. *A l'inverse, si l'on dispose d'un modèle exponentiel régulier (3), alors il n'existe qu'une seule fonction (à une transformation affine près) du paramètre pouvant être estimée efficacement, il s'agit de*

$$g(\theta) = -\frac{1}{n} \cdot \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta)\right)^{-1}_{1 \leq i, j \leq p} \cdot \frac{\partial \beta}{\partial \theta}(\theta) \text{ (noter que cette fonction semble dépendre de } n; \text{ dans le cas de v.a.i.i.d.}$$

ce n'est pas le cas). L'estimateur est alors : $\widehat{T} = \frac{1}{n} \cdot (a_1(X_1, \dots, X_n), \dots, a_p(X_1, \dots, X_n))$ et sa matrice de covariance minimale est donnée par sa borne de Cramer-Rao, soit :

$$\text{cov}_\theta(\widehat{T}) = \frac{1}{n} \cdot \frac{\partial g}{\partial \theta}(\theta) \cdot \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta)\right)^{-1}_{1 \leq i, j \leq d}.$$

3.5 Estimateur du maximum de vraisemblance

Nous allons voir une méthode permettant d'obtenir aisément et dans la plupart des cas un estimateur possédant de très bonnes qualités... Par la suite on se place une nouvelle fois dans le cadre d'un modèle statistique paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, avec $\Theta \subset \mathbb{R}^p$, dominé.

Définition. *Pour $(x_1, \dots, x_n) \in (\Omega')^n$, soit $\theta \in \Theta \mapsto L_\theta(x_1, \dots, x_n)$ la vraisemblance du modèle. On appelle estimateur du maximum de vraisemblance une statistique $\widehat{\theta}_n$ telle que pour (X_1, \dots, X_n) un n -échantillon quelconque du modèle :*

$$L_{\widehat{\theta}_n}(X_1, \dots, X_n) = \sup_{\theta \in \Theta} L_\theta(X_1, \dots, X_n).$$

Remarque. *Il n'y a pas de garantie de l'unicité d'un tel estimateur. Une méthode pour l'obtenir (mais pas toujours) est de rechercher un extremum local de L_θ sur Θ , ce qui pourra être fait en annulant les dérivées partielles de L_θ par θ_i . De même, il est clair que l'estimateur du maximum de vraisemblance pourra être également obtenu en maximisant le logarithme de la vraisemblance, appelé encore la log-vraisemblance. Enfin, si l'on désire estimer $g(\theta)$ avec g une fonction bijective, alors $g(\widehat{\theta})$ sera l'estimateur du maximum de vraisemblance de $g(\theta)$.*

Propriété. *Si il existe une statistique exhaustive \widehat{T} pour le modèle, alors $\widehat{\theta}$ est une fonction mesurable de \widehat{T} pour tout $\theta \in \Theta$.*

Démonstration : Si \widehat{T} est exhaustive, d'après le théorème de factorisation, la vraisemblance du modèle par rapport à la mesure dominante P^* est $g_\theta(\widehat{T}(x_1, \dots, x_n))$ pour tout $\theta \in \Theta$ et \mathbb{P}_θ -presque tout $(x, \dots, x_n) \in (\Omega')^n$, ce qui revient à P^* -presque tout $(x, \dots, x_n) \in (\Omega')^n$ par la même démonstration que celle de la nullité de l'information de Fisher d'une statistique libre. Ainsi, prendre l'argument maximal de $\theta \rightarrow L_\theta$ revient à prendre l'argument maximal de $\theta \rightarrow g_\theta(\widehat{T}(x_1, \dots, x_n))$, et $\widehat{\theta}$ sera donc une fonction de \widehat{T} . ■

Propriété. *On suppose que le modèle est régulier. Si on suppose qu'il existe un estimateur sans biais efficace de θ alors c'est l'estimateur du maximum de vraisemblance de θ .*

Démonstration : D'après ce qui précède, si le modèle est régulier et que \widehat{T} est un estimateur sans biais efficace de θ , alors le modèle est exponentiel et l'égalité (5) a encore lieu, soit pour tout $\theta \in \Theta$,

$$\frac{\partial}{\partial \theta}(\log L_\theta(x)) = \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta)\right)_{1 \leq i, j \leq p} \cdot \widehat{T} + \frac{\partial \beta}{\partial \theta}(\theta) \implies \left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta)\right)_{1 \leq i, j \leq p} \cdot \mathbb{E}_\theta(\widehat{T}) + \frac{\partial \beta}{\partial \theta}(\theta) = 0.$$

Comme \widehat{T} est un estimateur sans biais de θ , on a donc $\left(\frac{\partial \alpha_j}{\partial \theta_i}(\theta)\right)_{1 \leq i, j \leq p} \cdot \theta + \frac{\partial \beta}{\partial \theta}(\theta) = 0$, pour tout $\theta \in \Theta$,

ce qui s'applique également à $\widehat{\theta}$ et donc :

$$\left(\frac{\partial \alpha_j}{\partial \theta_i}(\widehat{\theta})\right)_{1 \leq i, j \leq p} \cdot \widehat{\theta} + \frac{\partial \beta}{\partial \theta}(\widehat{\theta}) = 0.$$

Mais d'après sa définition, le modèle étant régulier $\hat{\theta}$ minimise la log-vraisemblance et annule donc sa dérivée, ce qui implique que :

$$\left(\frac{\partial \alpha_j}{\partial \theta_i}(\hat{\theta}) \right)_{1 \leq i, j \leq p} \cdot \hat{T} + \frac{\partial \beta}{\partial \theta}(\hat{\theta}) = 0.$$

En conséquence, obtient :

$$\left(\frac{\partial \alpha_j}{\partial \theta_i}(\hat{\theta}) \right)_{1 \leq i, j \leq p} \cdot (\hat{T} - \hat{\theta}) = 0 \implies \hat{T} = \hat{\theta},$$

car la matrice des dérivées des α_j est supposée de rang d . Enfin, l'unicité de $\hat{\theta}$ est liée à l'écriture du modèle exponentiel. \blacksquare

Nous allons nous intéresser maintenant au comportement asymptotique de l'estimateur du maximum de vraisemblance (lorsqu'il existe), donc quand la taille n de l'échantillon tend vers l'infini. Il est clair que pour chaque n l'expression de l'estimateur est différente et, surtout, le modèle statistique change. Pour palier à cela, on se placera dans un "gros" modèle, $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, \mathbb{P}'_{\theta}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$ (la dimension du paramètre reste constante) dans lequel un échantillon est une suite de v.a. Par ailleurs, on supposera désormais que **tout échantillon de ce modèle est constitué de v.a.i.i.d.**, et que $d\mathbb{P}'_{\theta} = (f_{\theta} \cdot d\mu)^{\otimes \mathbb{N}}$, le modèle étant dominé par la mesure μ , et f_{θ} étant la densité de chaque X_i par rapport à μ .

Théorème (Convergence de l'estimateur du maximum de vraisemblance). *On suppose le modèle paramétrique $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, (f_{\theta} \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^d$ dominé par une mesure μ et régulier. On suppose en plus que le modèle est identifiable (au sens où $f_{\theta_1} = f_{\theta_2}$, μ -presque partout, entraîne $\theta_1 = \theta_2$). Alors si la suite $(X_n)_{n \in \mathbb{N}}$ est issue du modèle avec pour paramètre $\theta_0 \in \Theta$,*

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{p.s.} \theta_0 \quad \text{pour la mesure } (f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}.$$

Démonstration : En premier lieu, pour n fixé, il est clair que pour tout $\theta \in \Theta$:

$$\log(L_{\theta}(x_1, \dots, x_n)) - \log(L_{\theta_0}(x_1, \dots, x_n)) = \sum_{i=1}^n \log \left(\frac{f_{\theta}(x_i)}{f_{\theta_0}(x_i)} \right).$$

Par ailleurs, pour tout $i \in \mathbb{N}$, les X_i ont tous la même loi et pour $\theta \in \Theta$,

$$\begin{aligned} \mathbb{E}_{\theta_0} \left[\log \left(\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right) \right] &\leq \log \left(\mathbb{E}_{\theta_0} \left[\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right] \right) \quad (\text{Inégalité de Jensen pour la fonction } -\log) \\ &\leq \log(\mathbb{E}_{\mu} [f_{\theta}(X_i)]) \\ &\leq 0. \end{aligned}$$

En fait, du fait que la fonction $-\log$ est strictement convexe, la borne 0 ne peut être atteinte que si $f_{\theta} = f_{\theta_0}$. Ainsi, avec la contrainte d'un modèle identifiable, dès que $\theta \neq \theta_0$, alors :

$$\mathbb{E}_{\theta_0} \left[\log \left(\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right) \right] < 0.$$

On peut appliquer la loi forte des grands nombres pour les variables aléatoires $\left(\log \left(\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right) \right)_{i \in \mathbb{N}}$ (qui sont bien i.i.d. et \mathbb{L}^1 car le modèle est régulier), et ainsi :

$$\begin{aligned} \frac{1}{n} (\log(L_{\theta}(X_1, \dots, X_n)) - \log(L_{\theta_0}(X_1, \dots, X_n))) &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right) \\ &\xrightarrow[n \rightarrow +\infty]{p.s.} \mathbb{E}_{\theta_0} \left[\log \left(\frac{f_{\theta}(X_i)}{f_{\theta_0}(X_i)} \right) \right] < 0, \end{aligned}$$

la convergence presque sûre ayant lieu pour la mesure $(f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}$. Considérons maintenant pour tout $\varepsilon > 0$ une famille dénombrable $(\theta_i^{(\varepsilon)})_{i \in I}$ dense sur la sphère de centre θ_0 et de rayon ε . Du fait du caractère dénombrable de cette famille, pour tout $\varepsilon > 0$, il existe n_{ε} tel que pour tout $n \geq n_{\varepsilon}$, pour tout $i \in I$:

$$\log(L_{\theta_i^{(\varepsilon)}}(X_1, \dots, X_n)) < \log(L_{\theta_0}(X_1, \dots, X_n)) \quad \text{p.s. pour la mesure } (f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}.$$

Comme le modèle est régulier, pour tout $n \in \mathbb{N}^*$, la log-vraisemblance de X_1, \dots, X_n est continue sur Θ . De plus pour tout n elle atteint son unique maximum en θ_0 . En conséquence, pour $n \geq n_\varepsilon$, $\hat{\theta}_n$ sera à l'intérieur de la boule de centre θ_0 et de rayon ε (toujours p.s. pour la mesure $(f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}$). Le raisonnement étant vrai pour tout $\varepsilon > 0$, le théorème s'en déduit. ■

Théorème (Normalité asymptotique de l'estimateur du maximum de vraisemblance). *On suppose le modèle paramétrique $((\Omega)^\mathbb{N}, \mathcal{A}'_\mathbb{N}, (f_\theta \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ et régulier. On suppose en plus que le modèle est identifiable et que la fonction $\theta \in \Theta \mapsto L_\theta$ est de classe $\mathcal{C}^2(\Theta)$. Alors si la suite $(X_n)_{n \in \mathbb{N}}$ est issue du modèle avec pour paramètre $\theta_0 \in \Theta$:*

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, I_1^{-1}(\theta_0)),$$

où $I_1(\theta)$ est la matrice de Fisher de taille p (supposée inversible) pour la variable X_1 .

Démonstration : Comme le modèle est régulier, on peut différencier la vraisemblance et pour tout $\theta \in \Theta$, noter :

$$M_\theta(X_1, \dots, X_n) = \frac{1}{n} \frac{\partial}{\partial \theta} \log L_\theta(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log (f_\theta(X_i)).$$

Un développement limité d'ordre 1 de M_θ autour de θ_0 est possible (toujours en raison du modèle régulier) et donc pour tout $\theta \in \Theta$:

$$M_\theta(X_1, \dots, X_n) = M_{\theta_0}(X_1, \dots, X_n) + (\theta - \theta_0) \cdot \frac{\partial}{\partial \theta} M_{\theta^*}(X_1, \dots, X_n),$$

avec θ^* dans le segment $[\theta, \theta_0]$ (remarquons que $\frac{\partial}{\partial \theta} M_{\theta^*}(X_1, \dots, X_n)$ est une matrice carrée de taille d). Ainsi en remplaçant θ par $\hat{\theta}_n$, on obtient pour chaque n l'existence de θ_n^* appartenant au segment $[\hat{\theta}_n, \theta_0]$ tel que :

$$M_{\hat{\theta}_n}(X_1, \dots, X_n) = M_{\theta_0}(X_1, \dots, X_n) + (\hat{\theta}_n - \theta_0) \cdot \frac{\partial}{\partial \theta} M_{\theta_n^*}(X_1, \dots, X_n). \quad (6)$$

Pour un modèle régulier, on a vu que $\mathbb{E}_{\theta_0} \left(\frac{\partial^2}{\partial \theta^2} \log f_{\theta_0}(X_i) \right) = -I_1(\theta_0)$, matrice de Fisher pour n'importe quelle variable X_i . Ainsi, $\frac{\partial}{\partial \theta} M_\theta(\cdot)$ étant une moyenne empirique, on a par la loi forte des grands nombres :

$$\frac{\partial}{\partial \theta} M_{\theta_0}(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f_{\theta_0}(X_i) \xrightarrow[n \rightarrow +\infty]{p.s.} -I_1(\theta_0) \text{ pour la mesure } (f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}.$$

Maintenant, en utilisant le fait que les densités f_θ sont de classe $\mathcal{C}^2(\Theta)$ et en utilisant la convergence presque sûre de $\hat{\theta}_n$ vers θ_0 démontrée au théorème précédent, on a :

$$\frac{\partial}{\partial \theta} M_{\theta_n^*}(X_1, \dots, X_n) \xrightarrow[n \rightarrow +\infty]{p.s.} -I_1(\theta_0) \text{ pour la mesure } (f_{\theta_0} \cdot d\mu)^{\otimes \mathbb{N}}.$$

Finalement, comme $\hat{\theta}_n$ est le maximum d'une fonction de classe \mathcal{C}^1 , cet estimateur annule $M_{\hat{\theta}_n}(X_1, \dots, X_n)$, et donc l'égalité (6) devient :

$$M_{\theta_0}(X_1, \dots, X_n) \cdot I_1^{-1}(\theta_0) = (\hat{\theta}_n - \theta_0).$$

Enfin, comme $M_{\theta_0}(X_1, \dots, X_n)$ est une moyenne empirique, ce vecteur aléatoire vérifie un théorème de la limite centrale :

$$\sqrt{n} \left(M_{\theta_0}(X_1, \dots, X_n) - \mathbb{E}_{\theta_0} \left(\frac{\partial}{\partial \theta} \log f_{\theta_0}(X_i) \right) \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, I_1(\theta_0)),$$

d'après la première définition de l'information de Fisher. Comme $\mathbb{E}_{\theta_0} \left(\frac{\partial}{\partial \theta} \log f_{\theta_0}(X_i) \right) = 0$ (voir les propriétés précédentes), on obtient la normalité asymptotique de $\hat{\theta}_n$. ■

Remarque. *Sous ces hypothèses, l'estimateur du maximum de vraisemblance est asymptotiquement sans biais et efficace. Cependant, à n fixé, il peut avoir un biais et ne pas être un estimateur efficace.*

3.6 Régions de confiance

En pratique, estimer un paramètre le plus souvent ne suffit pas. On aimerait connaître plus précisément quelle marge de sécurité on a sur la connaissance de ce paramètre.

Définition. On se place dans le cadre d'un modèle paramétrique $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$. Soit $\alpha \in]0, 1[$ un nombre fixé a priori. On appelle région de confiance du paramètre θ de niveau $1 - \alpha$ un sous-ensemble aléatoire $R_{1-\alpha}$ inclus dans \mathbb{R}^p et défini sur $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}})$, tel que pour tout $\theta \in \Theta$, $\{(x_1, \dots, x_n) \in (\Omega')^{\mathbb{N}}, \theta \in R_{1-\alpha}(x_1, \dots, x_n)\} \in \mathcal{A}'_n$ et :

$$\inf_{\theta \in \Theta} \{\mathbb{P}_\theta(\theta \in R_{1-\alpha})\} \geq 1 - \alpha. \quad (7)$$

Si un échantillon observé $(X_1(\omega), \dots, X_n(\omega))$ est connu, $R_{1-\alpha}(X_1(\omega), \dots, X_n(\omega))$ est appelé région de confiance observé. Dans le cas où le paramètre est un réel ($p = 1$), on pourra obtenir un intervalle de confiance.

Comment déterminer une région de confiance ? En premier lieu, il est clair que pour tout $\alpha \in]0, 1[$, $R_{1-\alpha} \subset \Theta$ (en général, on choisit α proche de 0, et en particulier $\alpha = 0.05$ est très souvent utilisé). Une démarche possible pour la construction de région de confiance est la suivante : naturellement, on désirerait utiliser un estimateur \hat{T} convergent de θ , mais sa loi dépend en général de θ ce qui rend difficile (à part quelques exceptions) son utilisation directe. On préférera donc utiliser ce que l'on appelle une fonction pivotale $\pi(\hat{T}, \theta)$, qui est une fonction mesurable d'un estimateur et de θ et qui est une statistique libre. On essaiera alors d'écrire la propriété (7) sous la forme

$$\inf_{\theta \in \Theta} \left\{ \mathbb{P}_\theta(\pi(\hat{T}, \theta) \in C_\alpha) \right\} \geq 1 - \alpha,$$

où C_α est une région déterministe. Aussi pourra-t-on ensuite construire la région de confiance en fonction des quantiles (souvent à $\alpha/2$ et $1 - \alpha/2$) de la loi de la fonction pivotale.

Exemple. Si le modèle est régulier, sous les conditions du théorème de normalité asymptotique du maximum de vraisemblance, on peut également montrer (théorème de Slutski) que

$$\pi(\hat{\theta}_n, \theta_0) = \sqrt{n} \cdot (I_1(\hat{\theta}_n))^{1/2} \cdot (\hat{\theta}_n - \theta_0) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, I_p),$$

où I_d est la matrice identité de taille p et $(I_1(\theta))^{1/2} \cdot (I_1(\theta))^{1/2} = I_1(\theta)$ pour tout $\theta \in \Theta$. Ainsi, si n est grand, on pourra assimiler la loi de $\pi(\hat{\theta}_n, \theta_0)$ avec la loi normale centrée réduite multidimensionnelle. Or si $Z \sim \mathcal{N}_p(0, I_p)$, avec $q_{1-\alpha/2} > 0$ le quantile d'une loi normale centrée réduite réelle de niveau $1 - \alpha/2$, tel que $P(Z \in [-q_{1-\alpha/2}, q_{1-\alpha/2}]^d) \geq 1 - \alpha$. Aussi le polyèdre $n^{-1/2} \cdot (I_1(\hat{\theta}_n))^{-1/2} \cdot [-q_{1-\alpha/2}, q_{1-\alpha/2}]^d$ recentré autour de $\hat{\theta}_n$ formera la région de confiance cherchée.

4 Tests paramétriques

4.1 Principes d'un test

Un test permet, à partir d'une réalisation d'un échantillon, de décider entre deux hypothèses, en mettant en avant une hypothèse privilégiée, appelée hypothèse H_0 , et une hypothèse alternative, appelée H_1 . On associe à un test un niveau α (avec souvent $\alpha \simeq 0.05$) et une puissance $1 - \beta$. La plupart du temps, α est fixé a priori et β s'en déduit. Plus précisément,

Définition. On se place dans le cadre d'un modèle paramétrique dominé $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$ et soit θ la "vraie" valeur du paramètre. Un problème de test est un choix entre deux hypothèses :

$$\begin{cases} H_0 : \theta \in \Theta_0 & : \text{hypothèse dite nulle} \\ H_1 : \theta \in \Theta_1 & : \text{hypothèse dite alternative,} \end{cases} \quad (8)$$

où $\Theta_0 \subset \mathbb{R}^p$, $\Theta_1 \subset \mathbb{R}^d$ et $\Theta_0 \cap \Theta_1 = \emptyset$.

Ceci posé, on peut préciser deux types de problèmes de tests suivant les constitutions de Θ_0 et Θ_1 :

Définition. Une hypothèse (H_0 ou H_1) est dite simple si elle est associée à un singleton (Θ_0 ou Θ_1). Sinon, elle sera dite composite. Dans le cas réel ($\Theta \subset \mathbb{R}$), si H_0 est simple de la forme $\theta = \theta_0$, et si H_1 est composite de la forme $\theta > \theta_0$ ou $\theta < \theta_0$, on parlera de test unilatéral; si H_1 est composite de la forme $\theta \neq \theta_0$, on parlera de test bilatéral.

Comment faire pour choisir entre les deux hypothèses H_1 et H_2 ? Il faudra partir de ce que l'on peut connaître du modèle, c'est-à-dire généralement un échantillon observé (X_1, \dots, X_n) . Pour cela, on définit une statistique qui sera la clé de voûte du test :

Définition. Dans le cadre du problème de test (8, soit \widehat{T} une statistique (donc une fonction mesurable d'un échantillon (X_1, \dots, X_n) issu du modèle) à valeurs dans \mathbb{R}^d , qui sera appelée statistique du test. Le test sera défini par la fonction $\widehat{\phi} = \mathbb{I}_{\widehat{T} \in W}$, où W est une partie de \mathbb{R}^p appelée région critique du test (et sa partie complémentaire dans \mathbb{R}^p est appelée région d'acceptation du test). Si $\widehat{\phi} = 1$, on choisira H_1 , sinon on décidera plutôt H_0 .

Donc, à chaque hypothèse H_0 et H_1 , on associe une partie de \mathbb{R}^p pour la statistique de test \widehat{T} . En général, ces parties ne sont pas Θ_0 et Θ_1 . Pour pouvoir précisément déterminer la région W , dans un cadre théorique (qui n'est pas le même que le cadre pratique, voir plus bas), on peut commencer par associer une fonction puissance à la statistique de test, puis définir les erreurs de premier espèce α et de deuxième espèce β :

Définition. Pour la statistique de test \widehat{T} , on associe :

- une fonction puissance, qui est la probabilité de choisir H_1 : $\theta \in \Theta_1 \mapsto \mathbb{P}_\theta(\widehat{T} \notin W)$.
- une erreur de première espèce : $P_{H_0}(\text{Choisir } H_1) = \alpha = \sup_{\theta \in \Theta_0} \mathbb{P}_\theta(\widehat{T} \in W)$;
- une erreur de seconde espèce : $P_{H_1}(\text{Choisir } H_0) = \beta = \sup_{\theta \in \Theta_1} \mathbb{P}_\theta(\widehat{T} \notin W)$.

La puissance du test est $1 - \beta$.

Cependant, ce qui vient d'être écrit reste théorique. En pratique, on utilisera plutôt la démarche suivante :

Construction concrète d'un test : On suppose le problème de test (8). On pose également a priori α qui dépend du problème posé (mais en général $\alpha = 0.05$), et $1 - \alpha$ est appelé le niveau du test. Par la suite, on réalise :

1. L'expression quantitative des hypothèses H_0 et H_1 .
2. Le choix de la statistique \widehat{T} du test.
3. La construction d'une région critique W à l'hypothèse H_1 par rapport à \widehat{T} .
4. La détermination explicite de W en fonction de α .
5. Le calcul (si possible) de la puissance du test $1 - \beta$.
6. Pour la réalisation de l'échantillon, rejet ou acceptation de H_0 .

Remarque : Cependant, en pratique on ne procède pas ainsi. On a donc deux types d'erreur. Le choix de l'hypothèse privilégiée est donc fondamental car le résultat d'un test n'est pas symétrique. Par exemple, supposons que l'on ait pour modèle $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{N}(\theta, 1)^{\otimes n}, \theta \in \mathbb{R})$ et que l'on veuille tester $H_0 : \theta = 0$ contre $H_1 : \theta = 1$ à partir d'un échantillon (X_1, \dots, X_n) du modèle. Nous verrons pourquoi un peu plus loin, \overline{X}_n est une statistique de test pertinente. Par exemple, si $n = 1$, et $X_1(\omega) = \overline{X}_1(\omega) = 0.8$, que va-t-on choisir entre H_0 et H_1 ? Naturellement, une région critique sera de la forme $[s, +\infty[$, où $s \in \mathbb{R}$, car \overline{X}_n est un estimateur de θ . On détermine s à l'aide de α , puisque $P_{H_0}(\text{Choisir } H_1) = \alpha = P_0(\overline{X}_1 \geq s)$, donc par exemple, si $\alpha = 0.05$, $s \simeq 1.65$. Par suite, si $\overline{X}_1(\omega) = 0.8$, on accepte H_0 et l'erreur de seconde espèce est $P_1(\overline{X}_1 < s) \simeq 0.74$, donc très élevée : le test n'est pas très discriminant. Maintenant, si on inverse H_0 et H_1 , soit $H_0 : \theta = 1$ contre $H_1 : \theta = 0$, le même résultat $X_1(\omega) = 0.8$, conduit à accepter H_0 , avec une erreur de second espèce encore $\simeq 0.74$. On obtient donc deux résultats opposés pour la même expérience aléatoire. Les hypothèses H_0 et H_1 ne sont clairement pas interchangeable.

La question qui se pose maintenant est de savoir comment trouver une statistique de test. Une idée naturelle dans ce cadre paramétrique serait d'utiliser un estimateur du maximum de vraisemblance.

4.2 Test de Wald

Un estimateur du maximum de vraisemblance permet d'associer à chaque hypothèse du test un ensemble de même "forme" que Θ_0 et Θ_1 . Cependant, la difficulté est trouver la loi de l'estimateur du maximum de vraisemblance $\hat{\theta}$ à n fixé. Si cela est possible, on utilisera directement $\hat{\theta}$ comme statistique de test.

Sinon, de manière plus générale, on connaît la loi asymptotique de $\hat{\theta}_n$ quand le modèle est régulier. Donc quand n est grand, on pourrait utiliser une loi normale comme approximation de la loi de $\hat{\theta}_n$. Mais, un nouvel obstacle apparaît : la matrice de covariance asymptotique, qui est la matrice d'information de Fisher inverse, dépend du paramètre θ . Aussi va-t-on préférer utiliser la statistique de test \hat{T} suivante :

Définition. Pour un modèle paramétrique dominé régulier $((\Omega')^n, \mathcal{A}'_n, \mathbb{P}_\theta, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$. La statistique de Wald \hat{T} pour le test $H_0 : \theta = \theta_0$ contre $H_1 : \theta \in \Theta_1$ est : $\hat{T}_n = n \cdot {}^t(\hat{\theta}_n - \theta) \cdot I(\theta) \cdot (\hat{\theta}_n - \theta)$.

Pour montrer "théoriquement" la pertinence de ce test, on va donc considérer la suite de tests (\hat{T}_n) en se plaçant dans le "grand" modèle asymptotique :

Théorème. Dans le cadre d'un modèle paramétrique $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, (f_\theta \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ et régulier, pour le problème de test $H_0 : \theta = \theta_0$ contre $H_1 : \theta \neq \theta_0$, alors, en notant \hat{T}_n la statistique de test de Wald pour le modèle projeté de taille n sous l'hypothèse H_0 ,

$$\hat{T}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi^2(p).$$

La région de rejet asymptotique du test sera donc de la forme $\hat{T}_n > s_\alpha$, où s_α est le quantile d'ordre $1 - \alpha$ de la loi du $\chi^2(p)$. La suite de test $(\hat{T}_n)_n$ a donc une puissance qui tend vers 1 lorsque α est fixé.

Démonstration : La loi asymptotique de $\hat{\theta}_n$ induit la loi asymptotique de \hat{T}_n , car $\sqrt{n} \cdot I(\theta)^{1/2} \cdot (\hat{\theta}_n - \theta)$ suit asymptotiquement une loi $\mathcal{N}(0, I_d)$ sous l'hypothèse H_0 et $\hat{T}_n = \|\sqrt{n} \cdot I(\theta)^{1/2} \cdot (\hat{\theta}_n - \theta)\|^2$. ■

Voici donc un premier type de test, qui sous certaines conditions de régularités du modèle et pour certaines hypothèses de tests est intéressant. Mais pourrait-on faire mieux ? Et en quel sens ? Désormais, il nous faut donc définir un moyen de comparaison entre deux tests.

4.3 Test du rapport de vraisemblance

Définition. Sous les hypothèses et notations précédentes, on dira qu'un test ϕ est uniformément le plus puissant (U.P.P.) au seuil α si le niveau de $\hat{\phi}$ associé à la statistique \hat{T} est inférieur ou égal à α et si pour tout autre test $\hat{\phi}'$ associé à la statistique \hat{T}' de niveau inférieur ou égal à α , $\forall \theta \in \Theta_1$,

$$\mathbb{E}_\theta(\hat{\phi}) = 1 - \mathbb{P}_\theta(\hat{T} \notin W) \leq 1 - \mathbb{P}_\theta(\hat{T}' \notin W') = \mathbb{E}_\theta(\hat{\phi}').$$

Définition. Sous les hypothèses précédentes, si $L_\theta(\cdot)$ est la vraisemblance, on appellera test du rapport de vraisemblance (test de Neyman-Person dans le cas d'hypothèses simples) un test de statistique \hat{T} telle que :

$$\hat{T} = \frac{\sup_{\theta \in \Theta_0} L_\theta(X_1, \dots, X_n)}{\sup_{\theta \in \Theta_1} L_\theta(X_1, \dots, X_n)}.$$

La région critique W associée à un tel test est de la forme $W =] + \infty, K[$ (donc si $\hat{T} < K$, on rejette H_0).

Une des vertus du test du rapport de vraisemblance par rapport au test de Wald est qu'il peut être utilisé dans un modèle non régulier (mais la question de sa loi, ou de la loi d'une fonctionnelle de ce test, demeure). De plus, la propriété suivante confirme l'intérêt de cette statistique de test :

Propriété (Principe de Lehmann). Dans le cas du test de deux hypothèses simples, ou d'un test unilatéral ($\Theta \subset \mathbb{R}$), ce test est U.P.P. Dans le cas d'un test bilatéral, il n'existe pas forcément de test U.P.P.

Démonstration : ■

Enfin, un tel test pour un modèle régulier, va pouvoir être traité de manière générale grâce à la normalité asymptotique de l'estimateur du maximum de vraisemblance :

Théorème. Dans le cadre d'un modèle paramétrique $((\Omega')^{\mathbb{N}}, \mathcal{A}'_{\mathbb{N}}, (f_{\theta} \cdot d\mu)^{\otimes \mathbb{N}}, \theta \in \Theta)$, où $\Theta \subset \mathbb{R}^p$, dominé par une mesure μ et régulier, pour le problème de test $H_0 : \theta = \theta_0$ contre $H_1 : \theta \neq \theta_0$, alors, en notant \widehat{T}_n la statistique du rapport de vraisemblance pour le modèle projeté de taille n ,

$$-2 \log(\widehat{T}_n) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi^2(p).$$

La région de rejet asymptotique du test sera donc de la forme $-2 \log(\widehat{T}_n) > s_{\alpha}$, où s_{α} est le quantile d'ordre $1 - \alpha$ de la loi du $\chi^2(p)$. La suite de test $(\widehat{T}_n)_n$ a donc une puissance qui tend vers 1 lorsque α est fixé.

Démonstration : la démonstration reprend un peu celle de la normalité asymptotique du maximum de vraisemblance. ■

5 Introduction à la statistique non-paramétrique

On se place donc un modèle semi-paramétrique, soit $((\Omega')^n, \mathcal{A}'_n, P_{(\theta, f)}, \theta \in \Theta, f \in \mathcal{F})$, où $\Theta \subset \mathbb{R}^p$ et \mathcal{F} n'est pas de dimension finie, ou dans un modèle non-paramétrique, soit $((\Omega')^n, \mathcal{A}'_n, P_f, f \in \mathcal{F})$, où \mathcal{F} n'est pas de dimension finie. **Dans un tel cadre, les estimations et tests construits à partir de la vraisemblance ne sont plus directement utilisables.**

5.1 Estimation semi-paramétrique et non-paramétrique

• Quelques estimations semi-paramétriques

Soit (X_1, \dots, X_n) un n -échantillon issu d'une suite de v.a.i.i.d. réelles.

1. Estimation de la moyenne : si on suppose que $\mathbb{E}X_i = m$ existe, un estimateur convergent (grâce à la L.F.G.N.) de m est : $\bar{X}_n = \frac{1}{n} \sum_{j=1}^n X_j$. De plus, si $\text{var}X_i < +\infty$, le T.L.C. de Slutsky montre qu'avec

$$\bar{\sigma}_n^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2,$$

$$\sqrt{n} \frac{(\bar{X}_n - m)}{\bar{\sigma}_n} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

Ceci permet d'obtenir des intervalles de confiance sur m .

2. Estimation de la variance : si on suppose que $\text{var}X_i = \sigma^2$ existe, un estimateur convergent (grâce à la L.F.G.N.) de σ^2 est : $\bar{\sigma}_n^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2$. De plus, si $\mathbb{E}X_i^4 < +\infty$, le T.L.C. de Slutski montre

$$\text{qu'avec } \bar{\mu}_{4,n} = \frac{1}{n} \sum_{j=1}^n ((X_j - \bar{X}_n)^2 - \bar{\sigma}_n^2)^2,$$

$$\sqrt{n} \frac{(\bar{\sigma}_n^2 - \sigma^2)}{\sqrt{\bar{\mu}_{4,n}}} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

Ceci permet d'obtenir des intervalles de confiance sur σ^2 .

• Quelques estimations non-paramétriques

Soit (X_1, \dots, X_n) un n -échantillon issu d'une suite de v.a.i.i.d. réelles.

1. Estimation de la fonction de répartition :

$$F_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}_{X_j \leq x},$$

pour $x \in \mathbb{R}$ est un estimateur convergent de $F(x)$.

2. Estimation de la densité par histogramme :

- si X est une v.a. discrète, on utilise les valeurs $(x_i)_{i \in I}$ possibles des X_j , et

$$f_n(x) = \sum_{i \in I} \frac{1}{x_{i+1} - x_i} \mathbb{I}_{x \in [x_i, x_{i+1}[} \left(\frac{1}{n} \sum_{j=1}^n \mathbb{I}_{x_i \leq X_j < x_{i+1}} \right), \quad \text{pour } x \in \mathbb{R}.$$

- si X est une v.a. continue, on découpe \mathbb{R} en m -classe **a priori**, $[x_i, x_{i+1}]$. On estime f par

$$f_n(x) = \sum_{i=1}^m \frac{1}{x_{i+1} - x_i} \mathbb{I}_{x \in [x_i, x_{i+1}[} \left(\frac{1}{n} \sum_{j=1}^n \mathbb{I}_{x_i \leq X_j \leq x_{i+1}} \right), \quad \text{pour } x \in \mathbb{R}.$$

3. Estimation de la densité par noyau pour une variable **continu**. On utilise un noyau $K(x)$ qui vérifie :

- $K(x) \geq 0$ pour $x \in \mathbb{R}$ et $K(0) > 0$.
- $\int_{\mathbb{R}} K(u) du = 1$ et $\int_{\mathbb{R}} K^2(u) du < +\infty$.

Par exemple, $K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ (dit "noyau gaussien") ou $K(x) = (1 - |x|) \cdot \mathbb{I}_{|x| \leq 1}$ (fonction "triangle"). On définit alors un estimateur de f , densité de X_j par rapport à la mesure de Lebesgue sur \mathbb{R} , par

$$f_{n,h}(x) = \frac{1}{n} \sum_{j=1}^n \frac{1}{h} K\left(\frac{X_j - x}{h}\right),$$

pour $x \in \mathbb{R}$ et $h > 0$. Si $h = h_n$ avec $h_n \rightarrow 0$ et $n \cdot h_n \rightarrow +\infty$ quand $n \rightarrow \infty$, alors $f_n(x)$ est un estimateur convergent de $f(x)$.

5.2 Tests semi et non-paramétriques

Quelques tests semi-paramétriques

- Test sur la moyenne (ou test de Student) : on suppose un échantillon (X_1, \dots, X_n) issu d'une suite (X_i) de v.a.i.i.d. d'espérance m et de variance finie. On veut tester $H_0 : m = m_0$ contre $H_1 : m \in \Theta_1$. Alors la statistique

$$T_n = \sqrt{n} \frac{\bar{X}_n - m}{\hat{\sigma}_n}$$

est la statistique du test de Student, et :

1. pour tout $n \geq 2$, on a $T_n \sim t(n-1)$ si les X_i sont gaussiennes (auquel cas le modèle est paramétrique);
 2. $T_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1)$ (dans le cadre semi-paramétrique).
- Test de comparaison de moyenne (ou test de Student) : on suppose deux échantillons (X_1, \dots, X_n) et (Y_1, \dots, Y_n) issus de suites (X_i) et (Y_i) de v.a.i.i.d. d'espérance m_X et m_Y et de même variance σ finie. On veut tester $H_0 : m_X = m_Y$ contre $H_1 : m_X - m_Y \in \Theta_1$. Alors la statistique :

$$T_n = \sqrt{n-1} \frac{(\bar{X}_n - \bar{Y}_n) - (m_X - m_Y)}{\sqrt{\hat{\sigma}_X^2 + \hat{\sigma}_Y^2}}$$

est la statistique du test de comparaison de Student, et :

1. pour tout $n \geq 2$, on a $T_n \sim t(2n-2)$ si les X_i sont gaussiennes (auquel cas le modèle est paramétrique);;
2. $T_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1)$ (dans le cadre semi-paramétrique).

- Test sur la variance : on suppose un échantillon (X_1, \dots, X_n) issu d'une suite (X_i) de v.a.i.i.d. de variance σ^2 et de moment d'ordre 4 fini. On veut tester $H_0 : \sigma^2 = \sigma_0$ contre $H_1 : \sigma^2 \in \Theta_1$. Alors la statistique de test,

$$T_n = \sqrt{n} \frac{(\hat{\sigma}_n^2 - \sigma^2)}{\sqrt{\bar{\mu}_{4,n}}},$$

où $\bar{\mu}_{4,n} = \frac{1}{n} \sum_{j=1}^n ((X_j - \bar{X}_n)^2 - \bar{\sigma}_n^2)^2$ est telle que $T_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1)$, d'où l'utilisation du test.

Quelques tests non-paramétriques

- Test de signe : on suppose n individus dont on peut connaître les valeurs de deux variables indépendantes X et Y , donc deux échantillons (X_1, \dots, X_n) et (Y_1, \dots, Y_n) de v.a.i.i.d. On veut tester $H_0 : P(X_i > Y_i) = P(X_i < Y_i)$ contre $H_1 : P(X_i < Y_i) > 1/2$. Pour cela, on considère la statistique :

$$T_n = \sum_{i=1}^n \mathbb{I}_{X_i < Y_i}.$$

Sous l'hypothèse H_0 , $T_n \sim \mathcal{B}(n, 1/2)$, d'où le test.

- Test de signe et de rang : on suppose n individus dont on peut connaître les valeurs de deux variables indépendantes X et Y , donc deux échantillons (X_1, \dots, X_n) et (Y_1, \dots, Y_n) de v.a.i.i.d. On veut tester $H_0 : P(X_i > Y_i) = P(X_i < Y_i)$ contre $H_1 : P(X_i < Y_i) > 1/2$. Pour cela, on considère la statistique :

$$T_n = \sum_{i=1}^n \text{rang}(|X_i - Y_i|) \cdot \mathbb{I}_{X_i < Y_i},$$

où $\text{rang}(|X_i - Y_i|)$ est le classement dans l'ordre croissant (du plus petit au plus grand) des différentes valeurs de $|X_i - Y_i|$. Sous l'hypothèse H_0 , on peut montrer que $T_n \sim \mathcal{L}_n$, où \mathcal{L}_n est une loi dont on connaît les quantiles. Ceci permet l'utilisation du test.

- Test d'ajustement : Tests du χ^2 et tests de Kolmogorov-Smirnov.