

Statistique inférentielle

Antoine Godichon-Baggioni

Table des matières

1	Convergence de suites de variables aléatoires	9
1.1	Rappels de probabilités	9
1.2	Modes de convergence	10
1.2.1	Convergence en loi	10
1.2.2	Convergence en probabilité	12
1.2.3	Convergence presque sûre	14
1.2.4	Convergence en moyenne quadratique	15
1.2.5	Résumé	16
1.3	Quelques inégalités classiques	17
1.4	Théorèmes asymptotiques	20
1.5	Opérations sur les limites	23
2	Estimation	29
2.1	Définitions	29
2.2	Estimation de la moyenne et de la variance	31
2.2.1	Estimation de la moyenne	31
2.2.2	Estimation de la variance	32
2.3	Méthode des moments	35
2.4	Méthode du maximum de vraisemblance	40
2.4.1	Cas discret	40
2.4.2	Cas continu	41
2.5	Comparaison d'estimateurs	44
2.5.1	Comparaison des erreurs quadratiques moyennes	44
2.5.2	Biais d'un estimateur	46
2.5.3	L'approche asymptotique	47
3	Intervalles de confiance	49
3.1	Intervalle de confiance	49
3.1.1	Intervalle de confiance	49
3.1.2	Notion de quantile	53
3.2	Rappels sur la loi normale	54

3.3	Intervalles de confiance dans le cas gaussien	57
3.3.1	Estimation de la moyenne lorsque la variance est connue	57
3.3.2	Estimation de la moyenne lorsque la variance est inconnue	59
3.3.3	Estimation de la variance	60
3.4	Intervalles de confiance asymptotiques	61
4	Vecteurs Gaussiens et théorème de Cochran	65
4.1	Vecteurs aléatoires	66
4.2	Vecteurs aléatoires gaussiens	67
4.3	Théorème de Cochran et applications	71
5	Tests	75
5.1	Généralités	75
5.1.1	Généralités	75
5.1.2	Etapes d'un test statistique	76
5.2	Tests sur la moyenne et la variance	77
5.2.1	Test de conformité d'une moyenne	77
5.2.2	Test d'inégalité $\mu \leq \mu_0$	80
5.2.3	Test d'inégalité $\mu \geq \mu_0$	81
5.2.4	Application :	82
5.2.5	Test de conformité d'une variance	83
5.3	Tests de comparaison de deux moyennes	85
5.3.1	Introduction	85
5.3.2	Test d'égalité	86
5.3.3	Test d'inégalité $\mu_1 \leq \mu_2$	92
5.3.4	Test d'inégalité $\mu_1 \geq \mu_2$	93
5.4	Test de Fischer et test de Shapiro-Wilk	94
5.4.1	Test de Fischer	95
5.4.2	Test de Shapiro-Wilk	97
5.5	Test de Student dans le cas apparié	97
5.5.1	Introduction	97
5.5.2	Test d'égalité	98
5.5.3	Test d'inégalité $\mu_1 \leq \mu_2$	101
5.5.4	Test d'inégalité $\mu_1 \geq \mu_2$	102
5.6	Tests du Khi-deux	103
5.6.1	Test d'indépendance du Khi-deux	103
5.6.2	Test d'adéquation	105
5.7	Tests asymptotiques	108
5.7.1	Introduction	108
5.7.2	Test d'égalité	108

5.7.3	Test d'inégalité $\theta_0 \geq \theta$	110
5.7.4	Test d'inégalité $\theta \geq \theta_0$	110

Introduction

Le cadre

On s'intéresse à une population de N individus, et plus particulièrement à une caractéristique a priori inconnue de cette population. Par exemple, on peut s'intéresser à la taille moyenne d'une population, savoir si le sexe des individus a une influence sur cette taille... Pour cela, on peut soit :

- effectuer un recensement, i.e mesurer chaque individu sur l'ensemble de la population. Cependant, cette méthode peut s'avérer très exhaustive, notamment si on a affaire à un grand nombre d'individus.
- procéder à un échantillonnage, i.e on réalise l'étude sur une partie de la population seulement.

On s'intéresse ici aux méthodes basées sur l'échantillonnage. En effet, elle représente généralement un coût beaucoup plus réduit que le recensement, l'étude est plus rapide et les erreurs d'observations sont plus réduites. Cependant, on n'étudie qu'une partie de la population, et si l'échantillon étudié n'est pas représentatif du reste de la population, l'extrapolation des résultats à la population globale risque d'être erronée. De plus les résultats obtenus peuvent varier d'un échantillon à l'autre.

Statistique descriptive Vs statistique inférentielle : La statistique descriptive n'a pas vocation à s'intéresser à la population globale mais juste à la sous-population formée par l'échantillon. Elle a pour but de résumer, via des tableaux, des indicateurs de position (moyenne, médiane,...) et de dispersion (variance, écart-type, écart inter-quartile...) et des graphiques, par exemple, l'information contenue par cet échantillon. La différence fondamentale entre statistique descriptive et statistique inférentielle est que l'on va supposer que les données étudiées sont des réalisations de variables aléatoires (modélisation) et s'appuyer sur la théorie des probabilités pour pouvoir extrapoler sur la population globale.

Fluctuation d'échantillonnage : Si on étudie un même caractère sur plusieurs échantillons (éventuellement de même taille) d'une même population, on peut observer que les résultats ne sont pas identiques selon les échantillons : ce phénomène s'appelle la fluctuation d'échantillonnage. Donc, les conclusions que l'on pourra tirer à l'aide d'un échantillon seront plus ou moins justes et varieront d'un échantillon à l'autre.

Prenons l'exemple du lancer d'une pièce. On va effectuer 5 lancers, et compter le nombre de fois où la pièce retombe sur "Pile". Ceci représente notre premier échantillon. On réitère notre expérience

	Echantillon 1	Echantillon 2	Echantillon 3	Echantillon 4
Nombre de "Pile"	1	3	2	3

TABLE 1 –

3 fois (ce qui représente en tous 4 échantillons) et on résume les résultats dans le tableau 1.

Objectifs du cours

On reprend l'exemple du lancer de pièce. On souhaite vérifier que la fameuse pièce n'est pas truquée. Pour cela on va effectuer n lancers, et on obtient ainsi des réalisations x_1, \dots, x_n , où x_i vaut 1 si le i -ème lancer vaut "Pile" et 0 sinon. On peut donc voir x_1, \dots, x_n comme des réalisations des variables aléatoires indépendantes X_1, \dots, X_n suivant une loi de Bernoulli de paramètre θ , et on considère donc que la pièce est équilibrée (non truquée) si $\theta = 0.5$. L'objectif de ce cours est donc :

- de faire quelques rappels de probabilités pour étudier au mieux les variables aléatoires étudiées et introduire les notions de convergence de suites de variables aléatoires. Dans l'exemple du lancer de pièce, on a compté le nombre de piles, i.e on a regardé $\sum_{i=1}^5 x_i$ qui est donc la réalisation de la variable aléatoire $\sum_{i=1}^5 X_i$. Quelle est la loi de cette variable aléatoire? Que se passe-t-il si on effectue plus (une infinité) de lancers?
- d'introduire une ou plusieurs méthodes d'estimations des paramètres qui nous intéressent : comment estimer le paramètre θ de la loi de Bernoulli qui "régit" ma pièce?
- d'apprendre à comparer ces différents estimateurs pour savoir lequel choisir.
- de construire des intervalles de confiance pour les différents paramètres : avec un certain risque, dire que θ appartient à un certain intervalle.
- de pouvoir répondre, avec un certain risque, si la pièce est truquée ou non.

Chapitre 1

Convergence de suites de variables aléatoires

1.1 Rappels de probabilités

Dans ce qui suit, on considère un espace probabilisé $(\Omega, \mathcal{F}, \mathbb{P})$ et on considère une variable aléatoire $X : \Omega \rightarrow \mathbb{R}$. On rappelle que la fonction de répartition $F_X : \mathbb{R} \rightarrow [0, 1]$ de X est définie pour tout $x \in \mathbb{R}$ par

$$F_X(x) = \mathbb{P}[X \leq x].$$

Rappelons que cette fonction est croissante, continue à droite et

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

Exemples :

1. Si X suit une loi de Bernoulli de paramètre p ,

$$F_X(x) = \begin{cases} 0 & \text{si } x < 0 \\ 1 - p & \text{si } 0 \leq x < 1 \\ 1 & \text{si } x \geq 1 \end{cases}$$

2. Si X suit une loi uniforme sur $[0, 1]$,

$$F_X(x) = \begin{cases} 0 & \text{si } x < 0 \\ x & \text{si } 0 \leq x \leq 1 \\ 1 & \text{si } x \geq 1 \end{cases}$$

3. Si X suit une loi normale centrée réduite,

$$F_X(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt.$$

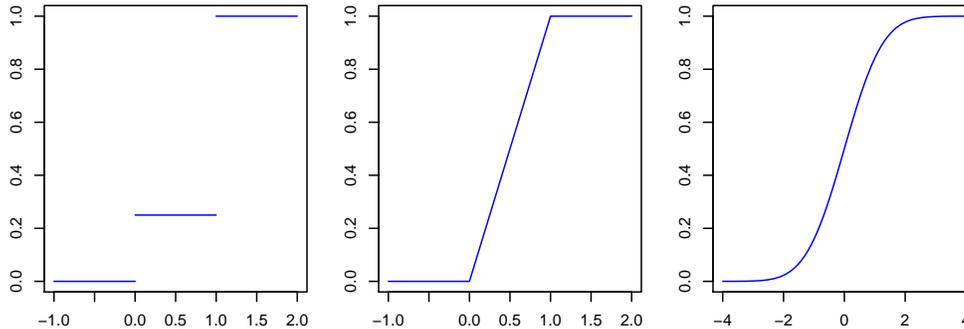


FIGURE 1.1 – Fonctions de répartition d’une loi de Bernoulli de paramètre $p = 0.75$ (à gauche), d’une loi uniforme (au centre) et d’une loi normale centrée réduite (à droite).

1.2 Modes de convergence

Dans ce qui suit, on s’intéresse aux différents modes de convergence d’une suite de variables aléatoires (X_n) . Plus précisément, on s’intéresse aux convergences en loi, en probabilité, presque sûre et en moyenne quadratique.

1.2.1 Convergence en loi

Definition 1.2.1 (Convergence en loi). On dit que la suite (X_n) converge en loi vers une variable aléatoire X et on note

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X$$

si pour toute fonction continue et bornée φ , on a

$$\mathbb{E}[\varphi(X_n)] \xrightarrow[n \rightarrow +\infty]{} \mathbb{E}[\varphi(X)].$$

De manière équivalente, (X_n) converge en loi vers X si pour toute fonction uniformément continue et bornée φ ,

$$\mathbb{E}[\varphi(X_n)] \xrightarrow[n \rightarrow +\infty]{} \mathbb{E}[\varphi(X)].$$

Exemple 1 : On considère une suite de variables aléatoires (X_n) , où pour tout $n \geq 1$, $X_n \sim \mathcal{B}\left(p + \frac{1-p}{n}\right)$, alors

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X \sim \mathcal{B}(p).$$

En effet, pour toute fonction continue et bornée φ , on a

$$\mathbb{E}[\varphi(X_n)] = \varphi(1) \left(p + \frac{1-p}{n} \right) + \varphi(0) \frac{n-1}{n} (1-p) \xrightarrow{n \rightarrow +\infty} p\varphi(1) + (1-p)\varphi(0) = \mathbb{E}[\varphi(X)].$$

Exemple 2 : On considère la suite (X_n) avec pour tout n , $X_n = -X$ où $X \sim \mathcal{N}(0, 1)$. Comme la loi normale est symétrique, on a $-X \sim \mathcal{N}(0, 1)$ et donc

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X.$$

Exemple 3 : Pour tout $n \geq 1$, on considère $X_n \sim \mathcal{N}(0, \sigma^2 + \frac{1}{n})$, on a pour toute fonction continue bornée φ , par convergence dominée

$$\mathbb{E}[\varphi(X_n)] = \int_{-\infty}^{\infty} \varphi(x) \frac{1}{\sqrt{2\pi(\sigma^2 + \frac{1}{n})}} \exp\left(-\frac{x^2}{2(\sigma^2 + \frac{1}{n})}\right) dx \xrightarrow{n \rightarrow +\infty} \int_{-\infty}^{\infty} \varphi(x) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} dx = \mathbb{E}[\varphi(X)],$$

où $X \sim \mathcal{N}(0, \sigma^2)$ et donc X_n converge en loi vers X .

Proposition 1.2.1. La suite de variables aléatoires (X_n) converge en loi vers X si et seulement si en tout point de continuité x de la fonction de répartition F_X de X ,

$$F_{X_n}(x) \xrightarrow[n \rightarrow +\infty]{} F_X(x).$$

Exemple 1 : On considère pour tout $n \geq 1$, $X_n \sim \mathcal{B}\left(p + \frac{1-p}{n}\right)$. On a alors pour tout $x \in \mathbb{R}$,

$$F_{X_n}(x) = \begin{cases} 0 & \text{si } x < 0 \\ \frac{n-1}{n}(1-p) & \text{si } 0 \leq x < 1 \\ 1 & \text{si } x \geq 1 \end{cases}$$

On a donc $F_{X_n}(x)$ qui converge vers $F_X(x)$ pour tout $x \notin \{0, 1\}$, où $X \sim \mathcal{B}(p)$.

Exemple 2 : Pour tout $n \geq 1$, on considère $X_n = -X$, où $X \sim \mathcal{N}(0, 1)$, alors pour tout $x \in \mathbb{R}$, par symétrie de la loi normale,

$$F_{X_n}(x) = \mathbb{P}[X_n \leq x] = \mathbb{P}[-X \leq x] = \mathbb{P}[X \geq -x] = \mathbb{P}[X \leq x] = F_X(x).$$

Exemple 3 : Pour tout n , $X_n \sim \mathcal{N}(0, \sigma^2 + \frac{1}{n})$. On a alors pour tout x ,

$$F_{X_n}(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi(\sigma^2 + \frac{1}{n})}} \exp\left(-\frac{t^2}{2(\sigma^2 + \frac{1}{n})}\right) dt.$$

Par convergence dominée,

$$F_{X_n}(x) \xrightarrow{n \rightarrow +\infty} \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{t^2}{2}} dt = F_X(x)$$

où $X \sim \mathcal{N}(0,1)$. On a donc la convergence en loi de X_n vers X .

Corollaire 1.2.1. *Supposons que les variables aléatoires X_n et X sont absolument continues de densités f_n et f . Alors, X_n converge en loi vers X si et seulement si pour tout $x \in \mathbb{R}$,*

$$f_n(x) \xrightarrow{n \rightarrow +\infty} f(x).$$

Si on reprend l'exemple 3, on obtient directement pour tout $x \in \mathbb{R}$,

$$f_n(x) \xrightarrow{n \rightarrow +\infty} f(x).$$

Enfin, rappelons qu'une variable aléatoire est caractérisée par sa fonction caractéristique, définie pour tout $t \in \mathbb{R}$ par $\Phi_X(t) = \mathbb{E}[e^{itX}]$. On peut alors réécrire la Proposition 1.2.1 comme :

Proposition 1.2.2. *La suite de variables aléatoires (X_n) converge en loi vers X si et seulement si pour tout $t \in \mathbb{R}$,*

$$\Phi_{X_n}(t) \xrightarrow{n \rightarrow +\infty} \Phi_X(t).$$

Passer par les fonctions caractéristiques peut représenter un grand intérêt lorsque l'on doit traiter des sommes de variables aléatoires indépendantes. En effet, rappelons que pour toutes variables aléatoires indépendantes X, Y ,

$$\Phi_{X+Y}(t) = \Phi_X(t)\Phi_Y(t).$$

Exemple : On considère l'exemple introductif du lancer d'une pièce de monnaie, et on considère les variables aléatoires X_i prenant comme valeur 1 si le i -ème lancer vaut "Pile" et 0 sinon. On a alors pour tout i , $X_i \sim \mathcal{B}(\theta)$ et un estimateur naturel de θ est $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$. On considère les lancers comme étant indépendants et on a

$$\Phi_{\hat{\theta}_n}(t) = \Phi_{\frac{X_1}{n}}(t) \dots \Phi_{\frac{X_n}{n}}(t) = \left(\Phi_{\frac{X_1}{n}}(t) \right)^n = \left(1 - \theta + \theta e^{i\frac{t}{n}} \right)^n \xrightarrow{n \rightarrow +\infty} e^{it\theta}$$

et on a donc la convergence en loi de $\hat{\theta}_n$ vers θ .

1.2.2 Convergence en probabilité

Definition 1.2.2 (Convergence en probabilité). *On dit que la suite (X_n) converge en probabilité vers une variable aléatoire X et on note*

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} X$$

si pour tout $\epsilon > 0$,

$$\mathbb{P}[|X_n - X| > \epsilon] \xrightarrow[n \rightarrow +\infty]{} 0$$

Exemple 1 : On considère pour tout $n \geq 1$, $X_n \sim \mathcal{B}\left(\frac{p}{n}\right)$ avec $p \in (0, 1)$. Alors X_n converge en probabilité vers 0. En effet pour tout $\epsilon \in (0, 1)$,

$$\mathbb{P}[|X_n - 0| > \epsilon] = \mathbb{P}[X_n > \epsilon] = \mathbb{P}[X_n = 1] = \frac{p}{n} \xrightarrow{n \rightarrow +\infty} 0.$$

Exemple 2 : On considère des variables aléatoires indépendantes et identiquement distribuées X_1, \dots, X_n avec $X_1 \sim \mathcal{U}([0, \theta])$, et on considère l'estimateur $X_{(n)} = \max_{i=1, \dots, n} X_i$. Alors $X_{(n)}$ converge presque sûrement vers θ . En effet, comme $X_{(n)} \leq \theta$, pour tout $\epsilon > 0$ on a

$$\mathbb{P}[|X_{(n)} - \theta| \geq \epsilon] = \mathbb{P}[\theta - X_{(n)} \geq \epsilon] = \mathbb{P}[X_{(n)} \leq \theta - \epsilon] = \mathbb{P}[\forall i, X_i \leq \theta - \epsilon]$$

A noter que cette probabilité est nulle si $\epsilon \geq \theta$, et on va donc prendre $\epsilon \in (0, \theta)$. De plus, comme les variables X_i sont indépendantes, on a

$$\mathbb{P}[|X_{(n)} - \theta| \geq \epsilon] = \prod_{i=1}^n \mathbb{P}[X_i \leq \theta - \epsilon] = (\mathbb{P}[X_1 \leq \theta - \epsilon])^n = \left(\frac{\theta - \epsilon}{\theta}\right)^n \xrightarrow{n \rightarrow +\infty} 0.$$

La proposition suivante donne le lien entre la convergence en loi et la convergence en probabilité.

Proposition 1.2.3. Soit (X_n) une suite de variables aléatoires. Si (X_n) converge en probabilité vers X , alors (X_n) converge en loi vers X .

Démonstration. Soit φ une fonction uniformément continue, et soit $M > 0$ telle que pour tout x , $|\varphi(x)| \leq M$. Soit $\epsilon > 0$, par continuité uniforme de φ , il existe $\eta > 0$ tel que pour tout x, y avec $|x - y| \leq \eta$, on ait $|\varphi(x) - \varphi(y)| \leq \epsilon$. On a alors

$$\begin{aligned} |\mathbb{E}[\varphi(X_n)] - \mathbb{E}[\varphi(X)]| &\leq \mathbb{E}[|\varphi(X_n) - \varphi(X)|] \\ &= \mathbb{E}[|\varphi(X_n) - \varphi(X)| \mathbf{1}_{\{|X_n - X| < \eta\}}] + \mathbb{E}[|\varphi(X_n) - \varphi(X)| \mathbf{1}_{\{|X_n - X| \geq \eta\}}] \\ &\leq \epsilon + 2 \sup_x \varphi(x) \mathbb{E}[\mathbf{1}_{\{|X_n - X| \geq \eta\}}] \\ &\leq \epsilon + 2MP \mathbb{P}[|X_n - X| \geq \eta]. \end{aligned}$$

Comme (X_n) converge en probabilité vers X , il existe un rang n_ϵ tel que pour tout $n \geq n_\epsilon$,

$$|\mathbb{E}[\varphi(X_n)] - \mathbb{E}[\varphi(X)]| \leq 2\epsilon,$$

et on obtient donc la convergence en loi de (X_n) vers X . □

Remarque : Attention la réciproque de la proposition précédente est généralement fautive.

Contre-exemple : Reprenons l'exemple naïf précédent, où X suit une loi normale centrée réduite et pour tout n , $X_n = -X$. On a vu que (X_n) converge en loi vers X (par symétrie), mais il est évident que X_n ne converge pas en probabilité vers X . En effet, prenons $\epsilon = 3.92$, on a alors

$$\mathbb{P}[|X_n - X| > \epsilon] = \mathbb{P}[|2X| > \epsilon] = 0.05$$

et ne converge donc pas vers 0.

Proposition 1.2.4. Soit (X_n) une suite de variables aléatoires qui converge en loi vers une constante c . Alors (X_n) converge en probabilité vers c .

Démonstration. Pour tout $\epsilon > 0$, on a

$$\begin{aligned}\mathbb{P}[|X_n - c| > \epsilon] &= \mathbb{P}[X_n - c > \epsilon] + \mathbb{P}[c - X_n > \epsilon] \\ &= 1 - F_{X_n}[\epsilon + c] + F_{X_n}[X_n \leq c - \epsilon] - \mathbb{P}[X_n = c - \epsilon].\end{aligned}$$

Comme pour tout $x \in \mathbb{R}$, on a $F_c(x) = \mathbf{1}_{x \geq c}$, on a que $x = c$ est le seul point de discontinuité de F_c et comme $\epsilon > 0$,

$$F_{X_n}(\epsilon + c) \xrightarrow{n \rightarrow +\infty} F_c(\epsilon + c) = 1 \quad \text{et} \quad F_{X_n}(c - \epsilon) \xrightarrow{n \rightarrow +\infty} F_c(c - \epsilon) = 0.$$

Enfin, comme une fonction de répartition est continue à droite, on a

$$\lim_{n \rightarrow +\infty} \mathbb{P}[X_n = c - \epsilon] = \lim_{n \rightarrow +\infty} \lim_{h \rightarrow 0^+} F_{X_n}(c - \epsilon + h) = \lim_{h \rightarrow 0^+} \lim_{n \rightarrow +\infty} F_{X_n}(c - \epsilon + h) = \lim_{h \rightarrow 0^+} F_c(c - \epsilon + h) = 0.$$

□

1.2.3 Convergence presque sûre

Definition 1.2.3 (Convergence presque sûre). On dit que la suite (X_n) converge presque sûrement vers X et on note

$$X_n \xrightarrow[n \rightarrow +\infty]{p.s.} X$$

si

$$\mathbb{P}\left[\lim_{n \rightarrow +\infty} X_n = X\right] = \mathbb{P}\left[\left\{\omega \in \Omega; \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\right\}\right] = 1.$$

La proposition suivante donne le lien entre la convergence presque sûre et la convergence en probabilité.

Proposition 1.2.5. Soit (X_n) une suite de variables aléatoires. Si (X_n) converge presque sûrement vers X , alors (X_n) converge en probabilité vers X .

Attention, la réciproque est fautive.

Contre-exemple. On considère une variable X suivant une loi uniforme sur $[0, 1]$ et la suite d'événements (A_n) définie par $A_1 = \{X \in [0, 1]\}$, $A_2 = \{X \in [0, 1/2]\}$, $A_3 = \{X \in [1/2, 1]\}$, $A_4 = \{X \in [0, 1/4]\}$, ... et on considère la suite de variable aléatoire (Y_n) définie pour tout $n \geq 1$ par $Y_n = \mathbf{1}_{A_n}$. On peut montrer que

$$\mathbb{P}[A_n] = \left(\frac{1}{2}\right)^{\lfloor \log_2(n) \rfloor}$$

où $[\cdot]$ est la fonction partie entière, et on a donc Y_n qui converge en probabilité vers 0. Cependant, pour tout $\omega \in [0, 1]$, pour tout $n \geq 1$, il existe n_ω tel que $X(\omega) \in A_{n_\omega}$, i.e

$$Y_n(\omega) = \mathbf{1}_{X(\omega) \in A_n} \rightarrow 0$$

et on n'a donc pas convergence presque sûre.

Cependant, grâce au lemme de Borel-Cantelli, si pour tout $\epsilon > 0$,

$$\sum_{n=1}^{+\infty} \mathbb{P} [|X_n - X| \geq \epsilon] < +\infty,$$

alors (X_n) converge presque sûrement vers X .

Exemple : Soit (X_n) une suite de variables aléatoires telle que $X_n \sim \mathcal{B} \left(\frac{1}{n^2} \right)$. Alors

$$X_n \xrightarrow[n \rightarrow +\infty]{p.s.} 0.$$

En effet pour tout $\epsilon \in (0, 1)$, on a

$$\mathbb{P} [X_n \geq \epsilon] = \mathbb{P} [X_n = 1] = \frac{1}{n^2} \quad \text{et} \quad \sum_{n \geq 1} \mathbb{P} [X_n \geq \epsilon] = \sum_{n \geq 1} \frac{1}{n^2} < +\infty.$$

1.2.4 Convergence en moyenne quadratique

Definition 1.2.4 (Convergence en moyenne quadratique). Soit (X_n) une suite de variables aléatoires de carré intégrable. On dit que (X_n) converge en moyenne quadratique vers X et on note

$$X_n \xrightarrow[n \rightarrow +\infty]{L^2} X$$

si

$$\mathbb{E} [|X_n - X|^2] \xrightarrow[n \rightarrow +\infty]{} 0.$$

Exemple : Reprenons l'exemple du lancer de pièce, i.e où X_1, \dots, X_n sont i.i.d et pour tout i , $X_i \sim \mathcal{B}(\theta)$. On a $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ et par linéarité de l'espérance

$$\mathbb{E} [\hat{\theta}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X_i] = \theta.$$

On a alors

$$\mathbb{E} [|\hat{\theta}_n - \theta|^2] = \mathbb{E} [(\hat{\theta}_n - \mathbb{E} [\hat{\theta}_n])^2] = \mathbb{V} [\hat{\theta}_n].$$

Par indépendance, on a

$$\mathbb{V} [\hat{\theta}_n] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{V} [X_i] = \frac{\theta(1-\theta)}{n},$$

et on obtient donc la convergence en moyenne quadratique de $\hat{\theta}_n$ vers θ .

La proposition suivante donne le lien entre convergence en probabilité et convergence en moyenne quadratique.

Proposition 1.2.6. *Soit (X_n) une suite de variables aléatoires. Si (X_n) converge en moyenne quadratique vers X , alors (X_n) converge en probabilité vers X .*

Démonstration. Voir TD. □

Attention la réciproque est fausse.

Contre-exemple : On considère une suite de variables aléatoires discrètes (X_n) avec pour tout n ,

$$\mathbb{P}[X_n = 0] = \frac{n-1}{n}, \quad \mathbb{P}[X_n = n] = \frac{1}{n}.$$

Soit $\epsilon > 0$, pour tout $n \geq \epsilon^{-1}$ (et $n \geq \epsilon$), on a

$$\mathbb{P}[X_n \geq \epsilon] = \mathbb{P}[X_n = n] = \frac{1}{n},$$

et donc X_n converge en probabilité vers 0. Cependant,

$$\mathbb{E}[X_n^2] = n \xrightarrow{n \rightarrow +\infty} +\infty.$$

1.2.5 Résumé

Le tableau suivant résume les liens entre les différents types de convergence.

$$\begin{array}{l} X_n \xrightarrow[n \rightarrow +\infty]{p.s} X \implies X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} X \implies X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X \\ X_n \xrightarrow[n \rightarrow +\infty]{L^2} X \implies X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} X \implies X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X \end{array}$$

Remarque : Il n'y a pas d'implication générale entre la convergence en moyenne quadratique et la convergence presque sûre.

Contre-exemple 1 : Soit X une variable aléatoire suivant une loi uniforme sur $[0, 1]$ et la suite d'évènements (A_n) définie par $A_1 = \{X \in [0, 1]\}$, $A_2 = \{X \in [0, 1/2]\}$, $A_3 = \{X \in [1/2, 1]\}$, $A_4 = \{X \in [0, 1/4]\}$, ... et on considère la suite de variable aléatoire (Y_n) définie pour tout $n \geq 1$ par $Y_n = \mathbf{1}_{A_n}$. On rappelle que

$$\mathbb{P}[A_n] = \left(\frac{1}{2}\right)^{\lfloor \log_2(n) \rfloor}$$

et on a donc

$$\mathbb{E}[Y_n^2] = \mathbb{E}[\mathbf{1}_{A_n}] = \left(\frac{1}{2}\right)^{\lfloor \log_2(n) \rfloor} \xrightarrow{n \rightarrow +\infty} 0$$

et on a donc convergence en moyenne quadratique. Cependant, on a vu que l'on n'a pas convergence presque sûre.

Contre-exemple 2 : Soit (X_n) une suite de variables aléatoires telle que $\mathbb{P}[X_n = 0] = 1 - \frac{1}{n^2}$ et $\mathbb{P}[X_n = n] = \frac{1}{n^2}$. Alors (X_n) converge presque sûrement vers 0 mais pas en moyenne quadratique. En effet, pour tout $\epsilon \in (0, 1)$,

$$\sum_{n \geq 1} \mathbb{P}[X_n \geq \epsilon] = \sum_{n \geq 1} \frac{1}{n^2} < +\infty$$

d'où la convergence presque sûre. Cependant,

$$\mathbb{E}[X_n^2] = \frac{n^2}{n^2} = 1.$$

1.3 Quelques inégalités classiques

Lorsque l'on considère une suite de variables aléatoires (X_n) et que l'on ne connaît pas leurs lois, il n'est pas aisé de démontrer leur convergence en probabilité directement. Les inégalités suivantes peuvent alors servir.

Proposition 1.3.1 (Inégalité de Markov). Soit $c, p > 0$ et X une variable aléatoire admettant un moment d'ordre p , on a

$$\mathbb{P}[|X| \geq c] \leq \frac{\mathbb{E}[|X|^p]}{c^p}.$$

Démonstration. On a

$$|X|^p = |X|^p \mathbf{1}_{\{|X| < c\}} + |X|^p \mathbf{1}_{\{|X| \geq c\}} \geq c^p \mathbf{1}_{\{|X| \geq c\}}$$

En passant à l'espérance, on obtient

$$\mathbb{E}[|X|^p] \geq c^p \mathbb{E}[\mathbf{1}_{\{|X| \geq c\}}] = c^p \mathbb{P}[|X| \geq c].$$

□

L'inégalité de Markov est par exemple très utile pour montrer que la convergence en moyenne quadratique implique la convergence en probabilité (voir TD).

Corollaire 1.3.1 (Inégalité de Bienaymé-Tchebychev). Soit $c > 0$ et X une variable aléatoire admettant un moment d'ordre 2. Alors,

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq c] \leq \frac{\mathbb{V}[X]}{c^2}$$

Démonstration. C'est une application directe de l'inégalité de Markov, en considérant la variable aléatoire $Y = X - \mathbb{E}[X]$ et en prenant $p = 2$. En effet, on a alors

$$\mathbb{P}[|X - \mathbb{E}[X]| \geq c] = \mathbb{P}[|Y| \geq c] \leq \frac{\mathbb{E}[Y^2]}{c^2} = \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{c^2} = \frac{\mathbb{V}[X]}{c^2}.$$

□

Exemple : On considère l'exemple introductif du lancer d'une pièce. Rappelons que l'on considère les variables aléatoires X_i prenant comme valeur 1 si le i -ème lancé vaut "Pile" et 0 sinon. Rappelons qu'un estimateur naturel du paramètre θ est $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$, et comme $\mathbb{E} [\hat{\theta}_n] = \theta$ et $\mathbb{V} [\hat{\theta}_n] = \frac{\theta(1-\theta)}{n}$, on obtient pour tout $\epsilon > 0$, à l'aide de l'inégalité de Bienaymé-Tchebychev

$$\mathbb{P} [|\hat{\theta}_n - \theta| \geq \epsilon] \leq \frac{\mathbb{V} [\hat{\theta}_n]}{\epsilon^2} = \frac{\theta(1-\theta)}{\epsilon^2 n},$$

et donc la convergence en probabilité de $\hat{\theta}_n$ vers θ .

Théorème 1.3.1 (Inégalité de Cauchy-Schwarz). Soit X, Y deux variables aléatoires admettant un moment d'ordre 2. Alors

$$\mathbb{E} [|XY|] \leq \sqrt{\mathbb{E} [X^2] \mathbb{E} [Y^2]}.$$

Démonstration. On considère la fonction

$$\lambda \longmapsto \mathbb{E} [(\lambda|X| - |Y|)^2] = \lambda^2 \mathbb{E} [X^2] - 2\lambda \mathbb{E} [|XY|] + \mathbb{E} [Y^2].$$

On a un polynôme positif, et son discriminant $\Delta = 4(\mathbb{E} [|XY|])^2 - 4\mathbb{E} [X^2] \mathbb{E} [Y^2]$ est donc négatif ou nul, i.e

$$(\mathbb{E} [|XY|])^2 \leq \mathbb{E} [X^2] \mathbb{E} [Y^2].$$

□

L'inégalité de Cauchy Schwarz peut être particulièrement utile pour obtenir la convergence en moyenne quadratique de produits d'estimateurs.

Exemple : Soit $(X_n), (Y_n)$ deux suites d'estimateurs convergeant en moyenne quadratique vers 0. On suppose qu'ils convergent également à l'ordre 4, i.e ils admettent des moments d'ordre 4 et

$$\mathbb{E} [|X_n - 0|^4] \xrightarrow[n \rightarrow +\infty]{} 0 \quad \text{et} \quad \mathbb{E} [|Y_n - 0|^4] \xrightarrow[n \rightarrow +\infty]{} 0.$$

Alors $X_n Y_n$ converge en moyenne quadratique vers 0. En effet, on a

$$\mathbb{E} [(X_n Y_n)^2] \leq \sqrt{\mathbb{E} [X_n^4] \mathbb{E} [Y_n^4]} \xrightarrow[n \rightarrow +\infty]{} 0,$$

i.e on a donc bien convergence en moyenne quadratique. Cependant, l'inégalité de Cauchy-Schwarz est assez grossière il est parfois nécessaire d'avoir des inégalités un peu moins grossière, ce que nous donne l'inégalité de Hölder.

Théorème 1.3.2 (Inégalité de Hölder). Soit X, Y deux variables aléatoires, et $p, q > 1$ tel que $\frac{1}{p} + \frac{1}{q} = 1$. Si X admet un moment d'ordre p et si Y admet un moment d'ordre q , alors

$$\mathbb{E} [|XY|] \leq (\mathbb{E} [|X|^p])^{\frac{1}{p}} (\mathbb{E} [|Y|^q])^{\frac{1}{q}}.$$

Démonstration. Remarquons que pour tout $a, b \geq 0$, on a

$$ab \leq \frac{1}{p}a^p + \frac{1}{q}b^q.$$

En effet, si on considère la fonction $g_b : a \mapsto \frac{1}{p}a^p + \frac{1}{q}b^q - ab$, on a

$$g'_b(a) = a^{p-1} - b$$

et le minimum atteint en $a^* = b^{\frac{1}{p-1}}$. De plus, comme $q = \frac{p}{p-1}$,

$$g_b(a^*) = \frac{1}{p}b^{\frac{p}{p-1}} + \frac{1}{q}b^{\frac{p}{p-1}} - b^{\frac{p}{p-1}} = 0.$$

En prenant

$$a = \frac{|X|}{(\mathbb{E}[|X|^p])^{1/p}} \quad \text{et} \quad b = \frac{|Y|}{(\mathbb{E}[|Y|^q])^{1/q}},$$

on obtient

$$\frac{|XY|}{(\mathbb{E}[|X|^p])^{1/p} (\mathbb{E}[|Y|^q])^{1/q}} \leq \frac{|X|^p}{p\mathbb{E}[|X|^p]} + \frac{|Y|^q}{q\mathbb{E}[|Y|^q]}$$

et en passant à l'espérance

$$\frac{\mathbb{E}[|XY|]}{(\mathbb{E}[|X|^p])^{1/p} (\mathbb{E}[|Y|^q])^{1/q}} \leq \frac{1}{p} + \frac{1}{q} = 1.$$

□

Exemple : Si X admet un moment d'ordre r , alors pour tout $r' \in (0, r)$,

$$\mathbb{E}[|X|^{r'}] \leq (\mathbb{E}[|X|^r])^{\frac{r'}{r}}.$$

Exemple : Comme dit précédemment, l'inégalité de Hölder permet un certain "raffinement" dans les calculs. Si on reprend l'exemple précédent mais en supposant cette fois que X_n converge à l'ordre r avec $2 < r < 4$, les calculs précédents ne sont plus valables. Cependant, si on suppose que Y_n converge à l'ordre $\frac{2r}{r-2}$, on a à l'aide de l'inégalité de Hölder (en prenant $p = r/2$ et $q = \frac{r}{r-2}$)

$$\mathbb{E}[(X_n Y_n)^2] \leq (\mathbb{E}[|X_n|^r])^{\frac{2}{r}} \left(\mathbb{E}[|Y_n|^{\frac{2r}{r-2}}] \right)^{\frac{r-2}{r}}.$$

et ce terme converge donc vers 0, i.e on a bien convergence en moyenne quadratique.

1.4 Théorèmes asymptotiques

A travers l'exemple précédent, on a vu comment obtenir la convergence en probabilité d'une suite de variables aléatoires. On rappelle ici les théorèmes classiques qui permettent d'obtenir rapidement ce type de (ou de meilleurs) résultats. Dans ce qui suit, on considère une suite de variables aléatoires indépendantes et identiquement distribuées (X_i) et on note

$$S_n := X_1 + \dots + X_n.$$

Théorème 1.4.1 (Loi faible des Grands Nombres (LGN)). *Soient X_1, \dots, X_n des variables aléatoires indépendantes et identiquement distribuées admettant une moyenne $m = \mathbb{E}[X_1]$, alors*

$$\frac{S_n}{n} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} m.$$

A noter que pour la loi faible des grands nombres, il est souvent noté comme hypothèse $\mathbb{E}[|X_1|] < +\infty$. En réalité, les deux hypothèses sont équivalentes car on ne peut parler d'espérance de X_1 que si X_1 admet un moment d'ordre 1. De plus, il est fréquent de voir la loi faible des grands nombres avec des hypothèses sur l'existence de la variance de X_1 (permettant ainsi d'utiliser l'inégalité de Bienaymé-Tchebychev pour démontrer le résultat, à l'image de ce qui a été fait dans les exemples précédents). Cependant, cette hypothèse est inutile, notamment si on voit la loi faible des grands nombres comme un corollaire de la loi forte des grands nombres ci-dessous :

Théorème 1.4.2 (Loi Forte des Grands Nombres (LFGN)). *Soient X_1, \dots, X_n des variables aléatoires indépendantes et identiquement distribuées admettant une moyenne $m = \mathbb{E}[X_1]$, alors*

$$\frac{S_n}{n} \xrightarrow[n \rightarrow +\infty]{p.s.} m.$$

Exemple : Reprenons l'exemple du pile ou face. On a des variables aléatoires i.i.d et d'espérance $\mathbb{E}[X_1] = \theta$. En appliquant la loi forte des grands nombres, on obtient

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{p.s.} \theta,$$

et on a en particulier la convergence en probabilité.

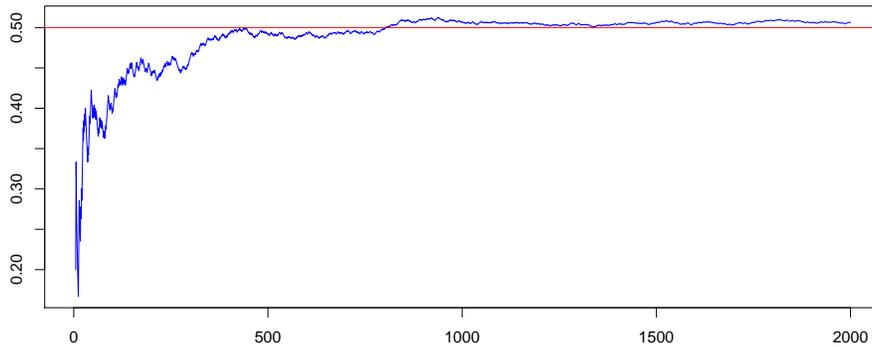


FIGURE 1.2 – Evolution en fonction de n d'une réalisation de $\hat{\theta}_n$ avec $\theta = 0.5$

La figure 1.2 laisse penser que l'estimateur $\hat{\theta}_n$ converge bien vers θ , et qu'il converge même rapidement. Les lois fortes des grands nombres, bien que très utiles pour obtenir "facilement" la convergence d'estimateurs, ne donnent malheureusement aucune indication sur la vitesse de convergence, contrairement au Théorème Central Limite :

Théorème 1.4.3 (Théorème Central Limite). Soit (X_n) une suite de variables aléatoires indépendantes et identiquement distribuées de moyenne $m = \mathbb{E}[X_1]$ et admettant une variance $\sigma^2 = \mathbb{V}[X_1]$, alors

$$\sqrt{n} \left(\frac{S_n}{n} - m \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

ce qui peut également s'écrire

$$\frac{\sqrt{n}}{\sigma} \left(\frac{S_n}{n} - m \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1),$$

Démonstration. A noter qu'en notant $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, on peut réécrire

$$Z_n := \frac{\sqrt{n}}{\sigma} \left(\frac{S_n}{n} - m \right) = \sqrt{n} \frac{\bar{X}_n - m}{\sigma} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \underbrace{\frac{X_i - m}{\sigma}}_{=: Y_i}.$$

A noter que les Y_i sont i.i.d et en regardant la fonction caractéristique, on a pour tout $t \in \mathbb{R}$,

$$\Phi_{Z_n}(t) = \mathbb{E} \left[e^{it \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i} \right] = \prod_{i=1}^n \mathbb{E} \left[e^{i \frac{t}{\sqrt{n}} Y_i} \right] = \left(\mathbb{E} \left[e^{i \frac{t}{\sqrt{n}} Y_1} \right] \right)^n.$$

On considère la fonction $\varphi : \mathbb{R} \rightarrow \mathbb{C}$ définie pour tout $u \in \mathbb{R}$ par $\varphi(u) = \mathbb{E} [e^{iuY_1}]$. Comme Y_1 admet un moment d'ordre 2, cette fonction est deux fois dérivable et

$$\varphi'(u) = i \mathbb{E} [Y_1 e^{iuY_1}] \quad \text{et} \quad \varphi''(u) = -\mathbb{E} [Y_1^2 e^{iuY_1}].$$

A l'aide d'un développement de Taylor, il existe une fonction g continue en 0 et telle que $g(0) = 0$ vérifiant

$$\varphi(u) = \varphi(0) + \varphi'(0)u + \frac{1}{2}u^2 (\varphi''(0) + g(u))$$

et on a donc

$$\varphi(u) = 1 + iu\mathbb{E}[Y_1] - \frac{1}{2}u^2 (\mathbb{E}[Y_1^2] + g(u)) = 1 - \frac{1}{2}u^2 (1 + g(u)).$$

Pour tout t , on a donc

$$\varphi\left(\frac{t}{\sqrt{n}}\right) = 1 - \frac{t^2}{2n} \left(1 + \epsilon\left(\frac{t}{\sqrt{n}}\right)\right) = 1 - \frac{t^2}{2n} + o\left(\frac{t}{2n}\right).$$

On a donc

$$\Phi_{Z_n}(t) = \varphi\left(\frac{t}{\sqrt{n}}\right)^n = \exp\left(n \log \varphi\left(\frac{t}{\sqrt{n}}\right)\right) \xrightarrow{n \rightarrow +\infty} \exp\left(-\frac{t^2}{2}\right)$$

et on retrouve ainsi la fonction caractéristique d'une loi normale centrée réduite, ce qui conclut la preuve.

□

Remarque : Le TCL est également appelé Théorème de la Limite Centrale (TLC).

Notons que dans l'exemple de la pièce équilibrée, on a $\mathbb{E}[X_1] = 1/2$, $\mathbb{V}[X_1] = 1/4$, et le TCL s'écrit alors

$$\sqrt{n} \left(\hat{\theta}_n - \frac{1}{2} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \frac{1}{4} \right),$$

ce que l'on peut également écrire

$$P_n := \sqrt{n} (2\hat{\theta}_n - 1) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

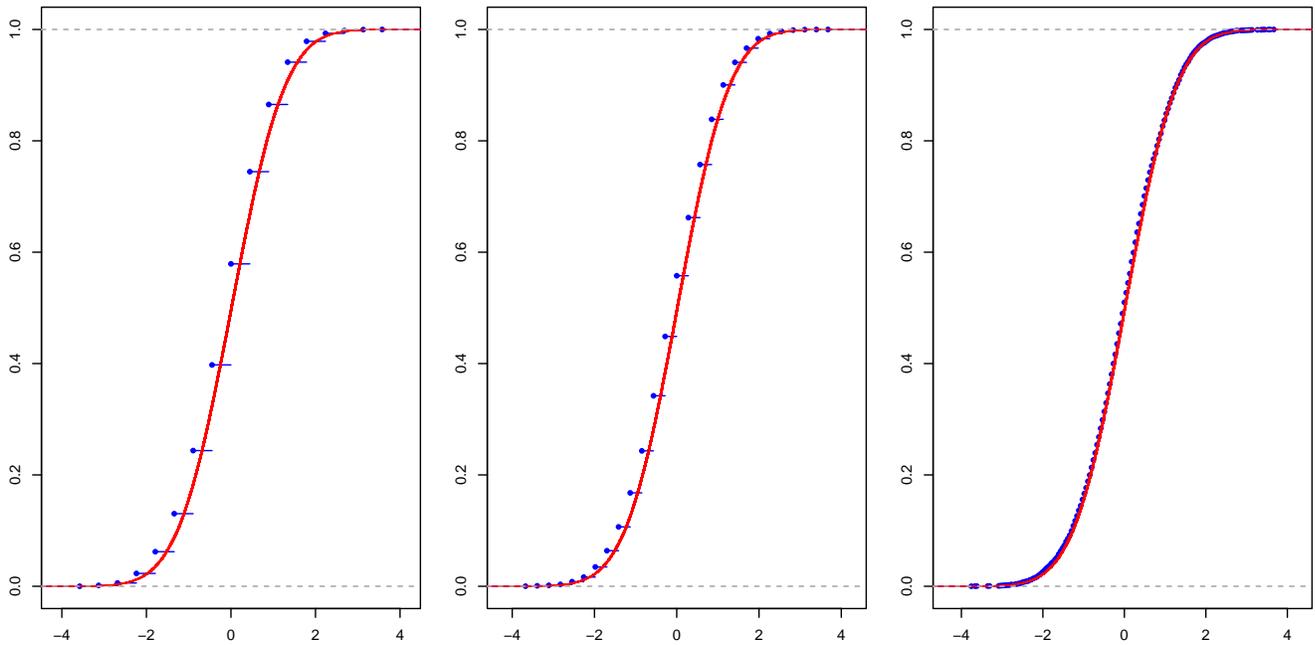


FIGURE 1.3 – Fonctions de répartition de P_n (en bleu) et d'une loi normale centrée réduite (en rouge) pour $n = 20$ (à gauche) $n = 50$ (au centre) et $n = 2000$ (à droite).

1.5 Opérations sur les limites

Si on reprend l'exemple du lancer de pièce, mais cette fois ci sans savoir si la pièce est équilibrée, i.e le paramètre θ est inconnu, le TLC nous donne alors

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

que l'on peut également écrire comme

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

Celui-ci est alors inexploitable en l'état (pour construire des intervalles de confiance par exemple) car le paramètre θ est inconnu, et on ne peut donc pas calculer $\sqrt{\theta(1 - \theta)}$. Cependant, un estimateur naturel serait $\sqrt{\hat{\theta}_n(1 - \hat{\theta}_n)}$, et la proposition suivante peut aider à montrer sa convergence.

Théorème 1.5.1 (Théorème de continuité). Soit (X_n) une suite de variables aléatoires et c une constante. Soit g une fonction une fonction continue en c , alors la suite $g(X_n)$ hérite du mode de convergence de la suite (X_n) , i.e

1. Si (X_n) converge presque sûrement vers c , alors $g(X_n)$ converge presque sûrement vers $g(c)$.
2. Si (X_n) converge en probabilité vers c , alors $g(X_n)$ converge en probabilité vers $g(c)$.
3. Si (X_n) converge en loi vers c , alors $g(X_n)$ converge en loi vers $g(c)$.

Démonstration. Preuve du point 3. Soit φ une fonction continue bornée. On a alors $\varphi \circ g$ qui est une fonction bornée et continue en c , et on a donc

$$\mathbb{E}[\varphi(g(X_n))] = \mathbb{E}[(\varphi \circ g)(X_n)] \xrightarrow{n \rightarrow +\infty} \mathbb{E}[(\varphi \circ g)c] = \varphi(g(c)).$$

Preuve du point 2 (version 1). On a vu que pour montrer qu'une suite de variables aléatoires converge en probabilité vers une constante, il suffit de montrer qu'elle converge en loi, ce que nous donne le point 1.

Preuve du point 2 (version 2). Comme g est continue en c , pour tout $\epsilon > 0$ il existe $\eta > 0$ tel que $|x - c| \leq \eta \implies |g(x) - g(c)| \leq \epsilon$. Ainsi,

$$\mathbb{P}[|g(X_n) - g(c)| > \epsilon] \leq \mathbb{P}[|X_n - c| > \eta] \xrightarrow{n \rightarrow +\infty} 0.$$

Preuve du point 1. On note

$$\Omega' = \left\{ \omega \in \Omega, X_n(\omega) \xrightarrow{n \rightarrow +\infty} c \right\}$$

et on rappelle que par définition de la convergence presque sûre, on a $\mathbb{P}[\Omega'] = 1$. De plus, pour tout $\omega \in \Omega'$, on a par continuité

$$g(X_n(\omega)) \xrightarrow{n \rightarrow +\infty} g(c).$$

□

Remarque : Attention, le théorème de continuité n'est pas vrai pour la convergence en moyenne quadratique. Le problème vient en partie du fait que $\mathbb{E}[g(X_n)]$ (où sa limite) n'est pas définie.

Contre-exemple : On considère une suite de variable aléatoire (X_n) telle que $\mathbb{P}[X_n = 0] = 1 - \frac{1}{n^3}$ et $\mathbb{P}[X_n = n] = \frac{1}{n^3}$. On a

$$\mathbb{E}[X_n^2] = \frac{1}{n} \xrightarrow{n \rightarrow +\infty} 0$$

et X_n converge donc en moyenne quadratique vers 0. On considère maintenant la fonction $g : x \mapsto x^2$, on a

$$\mathbb{E}[(g(X_n) - g(0))^2] = \mathbb{E}[X_n^4] = n \xrightarrow{n \rightarrow +\infty} +\infty.$$

Exemple : Reprenons l'exemple de la pièce avec $\theta \in (0, 1)$. On considère la fonction $g : x \mapsto \frac{1}{\sqrt{x(1-x)}}$ qui est continue en θ . On obtient donc, par le théorème de continuité,

$$\frac{1}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} \xrightarrow[n \rightarrow +\infty]{p.s.} \frac{1}{\sqrt{\theta(1-\theta)}}.$$

En particulier, comme la multiplication par une constante est également une fonction continue, on obtient

$$\frac{\sqrt{\theta(1-\theta)}}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} \xrightarrow[n \rightarrow +\infty]{p.s.} 1.$$

Attention! Le théorème de continuité ne nécessite pas que g soit continue voire même définie en X_n . Dans l'exemple précédent, si on note $g : x \longleftarrow \frac{1}{x(1-x)}$, elle n'est pas définie en 0 et 1 et donc pas nécessairement définie en $\hat{\theta}_n$. Cependant, on peut remarquer

$$\mathbb{P}[\hat{\theta}_n = 0] = \mathbb{P}[\forall i, X_i = 0] = (1 - \theta)^n \quad \text{et} \quad \mathbb{P}[\hat{\theta}_n = 1] = \mathbb{P}[\forall i, X_i = 1] = \theta^n.$$

Remarque 1.5.1. On peut donner une version moins restrictive du théorème de continuité, mais peut-être un peu moins digeste : soit (X_n) une suite de variables aléatoires et soit X une variable aléatoire. Soit g une fonction dont l'ensemble des points de discontinuité est noté D_g . Si $\mathbb{P}[X \in D_g] = 0$, alors la suite $g(X_n)$ hérite du mode de convergence de la suite (X_n) , i.e

1. Si (X_n) converge presque sûrement vers X , alors $g(X_n)$ converge presque sûrement vers $g(X)$.
2. Si (X_n) converge en probabilité vers X , alors $g(X_n)$ converge en probabilité vers $g(X)$.
3. Si (X_n) converge en loi vers X , alors $g(X_n)$ converge en loi vers $g(X)$.

Afin d'établir un TCL, il ne nous reste plus qu'à utiliser le théorème de Slutsky suivant :

Théorème 1.5.2 (Théorème de Slutsky). Soient $(X_n), (Y_n)$ des suites de variables aléatoires telles que (X_n) converge en loi vers X et (Y_n) converge en probabilité vers une constante c , alors

$$X_n + Y_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X + c \qquad X_n Y_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} cX.$$

Remarque : Le théorème précédent est généralement faux si Y_n converge en loi vers une variable aléatoire Y .

Contre-exemple : Soit $X \sim \mathcal{N}(0, 1)$ et $(X_n), (Y_n)$ deux suites de variables aléatoires telles que pour tout $n \geq 1$, $X_n = -X$ et $Y_n = X$. On a

$$X_n + Y_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} 0 \neq 2X$$

et

$$X_n Y_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} -X^2 \neq X^2.$$

Exemple : Reprenons l'exemple du lancer de pièce et rappelons que nous avons obtenu le TLC suivant :

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

Nous avons vu que $\frac{1}{\sqrt{\theta(1-\theta)}}$ est inconnu, ce qui pourrait (par la suite) nous empêcher de construire des intervalles de confiance. Naturellement, on va plutôt s'intéresser à la variable suivante

$$\sqrt{n} \frac{1}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} (\hat{\theta}_n - \theta),$$

que l'on peut également écrire comme

$$\frac{\sqrt{\theta(1-\theta)}}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} \sqrt{n} \frac{1}{\sqrt{\theta(1-\theta)}} (\hat{\theta}_n - \theta).$$

Comme on a

$$\frac{\sqrt{\theta(1-\theta)}}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} \xrightarrow[n \rightarrow +\infty]{p.s.} 1 \quad \text{et} \quad \sqrt{n} \frac{1}{\sqrt{\theta(1-\theta)}} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0,1),$$

en appliquant le Théorème de Slutsky (car la convergence p.s implique la convergence en probabilité), on obtient

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

Corollaire 1.5.1. Soit (X_n) une suite de variables aléatoires, a et σ^2 des constantes telles que

$$\sqrt{n} (X_n - a) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

alors (X_n) converge en probabilité vers a .

Démonstration. On va montrer, à l'aide du théorème de Slutsky, la convergence en loi de (X_n) vers a (qui implique alors la convergence en probabilité)

$$(X_n - a) = \frac{1}{\sqrt{n}} \sqrt{n} (X_n - a),$$

et comme $\frac{1}{\sqrt{n}}$ converge en probabilités vers 0, en appliquant Slutsky, on obtient le résultat. \square

Remarque 1.5.2. Notons que le corollaire précédent peut être généralisé, i. e on peut l'écrire comme suit : soit (v_n) une suite déterministe tendant vers $+\infty$, (X_n) une suite de variables aléatoires, X une variable aléatoire et a une constante telles que

$$v_n (X_n - a) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X$$

alors

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} a.$$

Théorème 1.5.3 (Delta méthode). Soit (X_n) une suite de variables aléatoires, a, σ^2 des constantes telles que

$$\sqrt{n} (X_n - a) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

alors

$$\sqrt{n} (g(X_n) - g(a)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, (g'(a))^2 \sigma^2\right).$$

Démonstration. D'après le Corollaire 1.5.1, on a

$$X_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} a.$$

De plus, comme g est dérivable en a , il existe (développement de Taylor) une fonction r telle que

$$g(x) = g(a) + (x - a)(g'(a) + r(x))$$

avec $r(a) = 0$ et r continue en a . Plus précisément, on considère la fonction définie pour tout $x \neq a$ par

$$r(x) = \frac{g(x) - g(a) - (x - a)g'(a)}{x - a}$$

et $r(a) = 0$. En particulier, comme X_n converge en probabilité vers a , en appliquant le théorème de continuité, on obtient que $r(X_n)$ converge en probabilité vers $r(a) = 0$. De plus,

$$g(X_n) = g(a) + (X_n - a)(g'(a) + r(X_n)),$$

ce que l'on peut écrire comme

$$\sqrt{n}(g(X_n) - g(a)) = (g'(a) + r(X_n))\sqrt{n}(X_n - a).$$

Comme

$$(g'(a) + r(X_n)) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} g'(a) \quad \text{et} \quad \sqrt{n}(X_n - a) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

on obtient le résultat en appliquant le Théorème de Slutsky. □

Exemple : Reprenons l'exemple du lancer de pièce. Supposons que l'on s'intéresse à l'estimation de $\frac{1}{\theta}$, avec $\theta > 0$. La fonction $g : x \mapsto \frac{1}{x}$ est continue en θ et par le théorème de continuité, on a donc

$$\frac{1}{\hat{\theta}_n} \xrightarrow[n \rightarrow +\infty]{p.s.} \frac{1}{\theta}.$$

De plus, comme

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

et comme g est dérivable en θ , on obtient

$$\sqrt{n}\left(\frac{1}{\hat{\theta}_n} - \frac{1}{\theta}\right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{1 - \theta}{\theta^3}\right).$$

Remarque 1.5.3. A noter que là aussi on a donné une version simplifiée de la Delta-méthode. En effet, on peut la réécrire comme : soit (X_n) une suite de variables aléatoires, (v_n) une suite de réels tendant vers $+\infty$,

X une variable aléatoire et a une constante telles que :

$$\sqrt{v_n} (X_n - a) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} X.$$

Soit g une fonction dérivable en a , alors

$$\sqrt{v_n} (g(X_n) - g(a)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} g'(a)X.$$

Chapitre 2

Estimation

Dans ce qui suit, on considère que l'on dispose de n données réelles x_1, \dots, x_n qui sont les mesures d'une variable quantitative. De plus, on supposera que ces données sont les réalisations de n variables aléatoires X_1, \dots, X_n que l'on suppose indépendantes et de même loi. On s'intéresse à une caractéristique de la loi des X_i , (espérance, variance,...) ou bien à un paramètre de cette loi, lorsque celle-ci est paramétrée (loi de Bernoulli, loi normale, loi exponentielle,...). On notera θ cette caractéristique ou ce paramètre supposé inconnu. En particulier, on s'intéressera aux deux exemples suivants :

Exemple 1 : le lancer de pièce. On considère des variables aléatoires X_i i.i.d suivant une loi de Bernoulli de paramètre θ . On s'intéressera à l'estimation du paramètre $\theta = \mathbb{E}[X_1]$.

Exemple 2 : la loi exponentielle. On considère des variables aléatoires X_i i.i.d suivant une loi exponentielle de paramètre θ . On s'intéressera à l'estimation du paramètre $\theta = (\mathbb{E}[X_1])^{-1}$.

2.1 Définitions

Dans ce qui suit, on considère un échantillon $\mathbf{X} = (X_1, \dots, X_n)$ dépendant d'un paramètre $\theta \in \Theta$, où Θ est un ouvert de \mathbb{R} .

Definition 2.1.1 (Statistique et estimateur). *Une statistique $T(\mathbf{X})$ est une fonction mesurable de l'échantillon \mathbf{X} ne dépendant pas de θ (mais dépendant éventuellement de paramètres connus). Un estimateur de θ est une statistique $\hat{\theta} = \theta(\mathbf{X})$ destinée à approcher θ .*

Exemple : Si on prend l'exemple du lancer de pièce, la variable

$$S_n = X_1 + \dots + X_n$$

est une statistique, mais pas un estimateur de θ contrairement à S_n/n .

Remarque : Attention! Ne pas confondre estimateur et estimation. Dans l'exemple précédent, on a S_n qui suit une loi binomiale de paramètre n, θ . On réalise 10 lancers de pièces et on obtient 3 "Pile" et 7 "Face". On obtient donc la réalisation $s_n = 3$ et l'estimation $\theta_n = 0.3$ (et pas $S_n = 3$ et $\hat{\theta}_n = 0.3$).

On verra par la suite qu'il peut exister plusieurs estimateurs d'un même paramètre. Une manière de comparer ces estimateurs est de comparer leur risque empirique ou erreur quadratique moyenne.

Definition 2.1.2 (Erreur quadratique moyenne). On suppose que θ est à valeurs dans $\Theta \subset \mathbb{R}$. L'erreur quadratique moyenne (ou risque quadratique) de l'estimateur $\hat{\theta}_n$ est définie pour tout $\theta \in \Theta$ par

$$EQM(\hat{\theta}_n, \theta) = \mathbb{E} [(\hat{\theta}_n - \theta)^2]$$

Notons que si on applique l'inégalité de Markov, on obtient pour toute constante $c > 0$,

$$\mathbb{P} [|\hat{\theta}_n - \theta| \geq c] \leq \frac{EQM(\hat{\theta}_n, \theta)}{c^2},$$

et donc, plus l'erreur quadratique moyenne est faible, plus la probabilité que $\hat{\theta}_n$ soit proche de θ est proche de 1. Pour comparer des estimateurs d'un même paramètre, on peut alors comparer leur erreur quadratique moyenne, et ce, pour toutes les valeurs possibles de θ .

Definition 2.1.3 (Biais d'un estimateur). On appelle biais d'un estimateur $\hat{\theta}_n$ de θ la quantité

$$B(\hat{\theta}_n, \theta) = \mathbb{E} [\hat{\theta}_n] - \theta.$$

1. S'il est nul, on dit que l'estimateur est sans biais ou non biaisé.
2. Si $\lim_{n \rightarrow \infty} B(\hat{\theta}_n, \theta) = 0$, on dit que l'estimateur est asymptotiquement sans biais.

Le biais permet de mesurer l'erreur moyenne de l'estimateur $\hat{\theta}_n$.

Exemple : Si on reprend l'exemple du lancer de pièce, on a vu que $\mathbb{E} [\hat{\theta}_n] = \theta$ et $\hat{\theta}_n$ est donc un estimateur sans biais de θ .

La décomposition suivante nous donne le lien entre le biais et l'erreur quadratique moyenne.

Proposition 2.1.1 (Décomposition Biais-Variance). On a

$$EQM(\hat{\theta}_n, \theta) = B(\hat{\theta}_n, \theta)^2 + \mathbb{V} [\hat{\theta}_n].$$

Démonstration. On a la décomposition suivante :

$$\begin{aligned} (\hat{\theta}_n - \theta)^2 &= ((\hat{\theta}_n - \mathbb{E} [\hat{\theta}_n]) - (\theta - \mathbb{E} [\hat{\theta}_n]))^2 \\ &= (\hat{\theta}_n - \mathbb{E} [\hat{\theta}_n])^2 - 2(\hat{\theta}_n - \mathbb{E} [\hat{\theta}_n])(\theta - \mathbb{E} [\hat{\theta}_n]) + (\theta - \mathbb{E} [\hat{\theta}_n])^2. \end{aligned}$$

En passant à l'espérance et par linéarité

$$\begin{aligned} \text{EQM}(\hat{\theta}_n, \theta) &= \mathbb{E} \left[(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2 - 2(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])(\theta - \mathbb{E}[\hat{\theta}_n]) + (\theta - \mathbb{E}[\hat{\theta}_n])^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2 \right] - 2\mathbb{E}[(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])(\theta - \mathbb{E}[\hat{\theta}_n])] + \mathbb{E}[(\theta - \mathbb{E}[\hat{\theta}_n])^2] \\ &= \mathbb{V}[\hat{\theta}_n] + B(\hat{\theta}_n, \theta)^2. \end{aligned}$$

□

Exemple : Prenons l'exemple du lancer de pièce, on a vu que $\hat{\theta}_n$ est un estimateur sans biais de θ . On obtient donc

$$\text{EQM}(\hat{\theta}_n, \theta) = \mathbb{V}[\hat{\theta}_n] = \frac{\theta(1-\theta)}{n}.$$

Cependant, il n'est pas toujours possible (ou facile) d'obtenir l'erreur quadratique moyenne d'un estimateur. Pour s'en convaincre, on peut considérer une variable aléatoire suivant une loi exponentielle de paramètre θ , proposer ensuite un estimateur de θ et essayer de calculer l'erreur quadratique moyenne. Il est alors important d'avoir d'autres critères pour comparer des estimateurs.

Definition 2.1.4 (Convergence, consistance et normalité asymptotique). On dit que l'estimateur $\hat{\theta}_n$ est

1. *convergent ou consistant* si

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \theta,$$

2. *fortement consistant* si

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{p.s.} \theta,$$

3. *asymptotiquement normal* si il existe $\sigma^2 > 0$ tel que

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2).$$

Remarque : A l'aide des nombreuses méthodes données précédemment (Slutsky, Delta méthode,...), on peut souvent montrer la normalité asymptotique des estimateurs. Ainsi, pour comparer deux estimateurs, on peut vérifier que ceux-ci convergent à la même vitesse avant de comparer les variances asymptotiques obtenues.

2.2 Estimation de la moyenne et de la variance

2.2.1 Estimation de la moyenne

On considère X_1, \dots, X_n des variables aléatoires indépendantes et de même loi, d'espérance $\mu = \mathbb{E}[X_1]$. Un estimateur naturel de la moyenne est donc

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Proposition 2.2.1.

1. \bar{X}_n est un estimateur sans biais et (fortement) consistant de μ .
2. Si $\sigma^2 = \mathbb{V}[X_1] < +\infty$, alors \bar{X}_n est asymptotiquement normal et

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

$$EQM(\bar{X}_n, \mu) = \mathbb{V}[\bar{X}_n] = \frac{\sigma^2}{n}$$

Dans la figure ci dessous, on voit que l'estimateur de la moyenne est bien convergent. De plus, on voit que les estimations sont centrées en la moyenne, ce qui illustre le fait que l'estimateur soit sans biais.

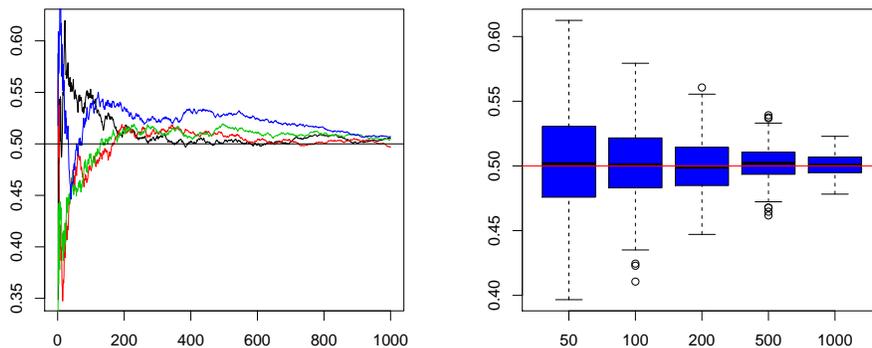


FIGURE 2.1 – Evolution de \bar{x}_n par rapport à n pour 4 échantillons (à gauche) et boxplots pour les \bar{x}_n obtenus pour $n = 50, 100, 200, 500, 1000$ à l'aide de 4000 échantillons (à droite).

2.2.2 Estimation de la variance

Lorsque μ est connu, un estimateur naturel de la variance est

$$\hat{V}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2,$$

et \hat{V}_n est alors un estimateur sans biais et fortement consistant de σ^2 car on peut alors le voir comme l'estimateur de la moyenne de la variable $Y = (X - \mu)^2$. Cependant, μ est généralement inconnu, et on le remplace alors par \bar{X}_n , obtenant ainsi la variance empirique

$$\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{2}{n} \sum_{i=1}^n X_i \bar{X}_n + \frac{1}{n} \sum_{i=1}^n \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2.$$

Proposition 2.2.2. $\hat{\sigma}_n^2$ est un estimateur biaisé de σ^2 mais asymptotiquement sans biais. Plus précisément,

on a

$$\mathbb{E} [\hat{\sigma}_n^2] = \frac{n-1}{n} \sigma^2.$$

Démonstration. On a

$$\mathbb{E} [\hat{\sigma}_n^2] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X_i^2] - \mathbb{E} [\bar{X}_n^2] = \mathbb{E} [X_1^2] - \mathbb{E} [\bar{X}_n^2]$$

Or,

$$\begin{aligned} \mathbb{E} [X_1^2] &= \mathbb{V} [X_1] + \mathbb{E} [X_1]^2 = \sigma^2 + \mu^2, \\ \mathbb{E} [\bar{X}_n^2] &= \mathbb{V} [\bar{X}_n] + \mathbb{E} [\bar{X}_n]^2 = \frac{1}{n} \sigma^2 + \mu^2, \end{aligned}$$

ce qui conclut la preuve. □

Definition 2.2.1. L'estimateur non biaisé de la variance S_n^2 est défini par

$$S_n^2 = \frac{n}{n-1} \hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n \bar{X}_n^2 \right).$$

Dans la figure ci-dessous, on voit que les estimations biaisées de la variance ne sont pas centrées en la variance, ce qui illustre le fait que l'estimateur soit biaisé, et ce décalage s'atténue (voire disparaît) lorsque la taille d'échantillon augmente, ce qui illustre le fait que l'estimateur soit asymptotiquement sans biais. Enfin, on voit que les estimations non biaisées de la variance sont bien centrées en la variance.

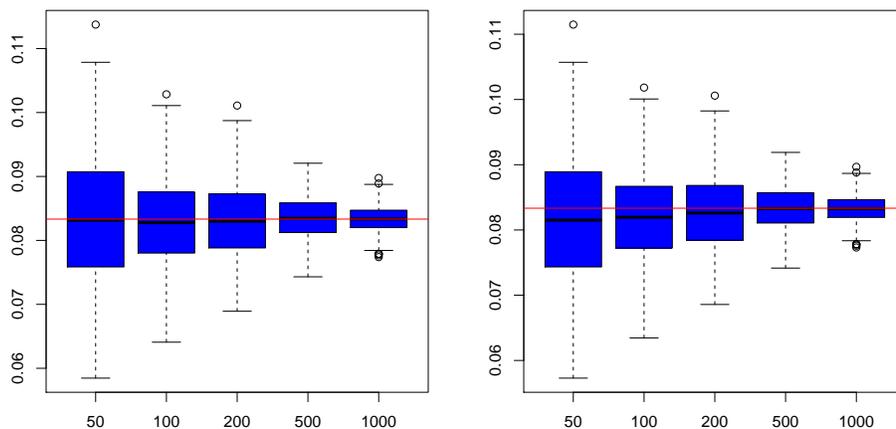


FIGURE 2.2 – Boxplots pour l'estimateur non biaisé (à gauche) et biaisé (à droite)

Proposition 2.2.3. Les estimateurs $\hat{\sigma}_n^2$ et S_n^2 de σ^2 sont consistants.

Démonstration. Rappelons que $\hat{\sigma}_n^2$ peut s'écrire

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2.$$

Par la loi des grands nombres, on a

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} E[X_1^2] \quad \text{et} \quad \bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mu,$$

et on obtient la consistance de $\hat{\sigma}_n^2$ à l'aide du théorème de Slutsky. De plus, on a

$$S_n^2 = \frac{n}{n-1} \hat{\sigma}_n^2$$

et on conclut là aussi grâce au théorème de Slutsky. \square

Proposition 2.2.4. *Si X_1 admet un moment d'ordre 4, alors $\hat{\sigma}_n^2$ et S_n^2 sont asymptotiquement normaux, et on a*

$$\sqrt{n} (\hat{\sigma}_n^2 - \sigma^2) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \tau^4 - \sigma^4) \quad \text{et} \quad \sqrt{n} (S_n^2 - \sigma^2) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \tau^4 - \sigma^4),$$

où

$$\tau^4 = \mathbb{E}[(X_1 - \mu)^4].$$

Démonstration. Posons $Y_i = X_i - \mu$. On a alors $\bar{Y}_n = \bar{X}_n - \mu$ et

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2.$$

On peut donc écrire

$$\sqrt{n} (\hat{\sigma}_n^2 - \sigma^2) = \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \sigma^2 \right) - \sqrt{n} \bar{Y}_n^2.$$

On va donc donner les convergences des deux termes à droite de l'égalité précédente. Notons d'abord que par la loi des grands nombres, \bar{Y}_n converge en probabilité vers 0 (car les Y_i sont d'espérance nulle). De plus, on peut écrire

$$\sqrt{n} \bar{Y}_n^2 = (\sqrt{n} \bar{Y}_n) \bar{Y}_n.$$

Le TCL appliqué au premier terme à droite de l'égalité précédente nous donne

$$\sqrt{n} \bar{Y}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

et comme \bar{Y}_n converge en probabilité vers 0, en appliquant le Théorème de Slutsky, on obtient

$$\sqrt{n}Y_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 0.$$

En appliquant le TCL à $\frac{1}{n} \sum_{i=1}^n Y_i^2$, on obtient

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \sigma^2 \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \tau^4 - \sigma^4 \right),$$

et donc, en appliquant Slutsky,

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \sigma^2 \right) - \sqrt{n}Y_n^2 \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \tau^4 - \sigma^4 \right).$$

Pour l'estimateur non biaisé, comme $S_n^2 = \frac{n}{n-1} \hat{\sigma}_n^2$, on a

$$\sqrt{n} (S_n^2 - \sigma^2) = \sqrt{n} \left(\frac{n}{n-1} \hat{\sigma}_n^2 - \sigma^2 \right) = \sqrt{n} (\hat{\sigma}_n^2 - \sigma^2) + \frac{1}{n} \sqrt{n} \hat{\sigma}_n^2$$

Pour avoir la normalité asymptotique, il faut donc montrer que le deuxième terme à droite de l'égalité précédente converge en probabilité vers 0. Notons que l'on peut écrire

$$\frac{1}{n} \sqrt{n} \hat{\sigma}_n^2 = \frac{1}{n} \sqrt{n} (\hat{\sigma}_n^2 - \sigma^2) + \frac{1}{\sqrt{n}} \sigma^2.$$

De plus, $\frac{1}{\sqrt{n}} \sigma^2$ converge vers 0, et par Slutsky,

$$\frac{1}{n} \sqrt{n} (\hat{\sigma}_n^2 - \sigma^2) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 0.$$

□

2.3 Méthode des moments

La méthode des moments repose sur la proposition suivante.

Proposition 2.3.1. *Soit Θ un intervalle ouvert de \mathbb{R} et $\theta \in \Theta$. Soit φ un C^1 -difféomorphisme de Θ dans $\varphi(\Theta)$, i.e une bijection de Θ dans $\varphi(\Theta)$, continûment dérivable et de réciproque continûment dérivable. Soit $\hat{\varphi}_n$ un estimateur consistant de $\varphi(\theta)$, alors $\hat{\theta}_n = \varphi^{-1}(\hat{\varphi}_n)$ est un estimateur consistant de θ , i.e*

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \theta.$$

De plus, si $\hat{\varphi}_n$ est un estimateur asymptotiquement normal de $\varphi(\theta)$, i.e si il existe $\sigma^2 > 0$ tel que

$$\sqrt{n} (\hat{\varphi}_n - \varphi(\theta)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2),$$

et si $\varphi'(\theta) \neq 0$, alors $\hat{\theta}_n$ est un estimateur asymptotiquement normal de θ , et plus précisément

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{\sigma^2}{(\varphi'(\theta))^2}\right).$$

Remarque : La propriété précédente reste vraie si on remplace \sqrt{n} par une suite v_n tendant vers $+\infty$.

Démonstration. Remarquons que comme $\hat{\varphi}_n$ converge en probabilité vers $\varphi(\theta)$ et comme Θ est un ouvert

$$\mathbb{P} [\hat{\varphi}_n \notin \varphi(\Theta)] \xrightarrow[n \rightarrow +\infty]{} 0,$$

et donc l'estimateur $\hat{\theta}_n = \varphi^{-1}(\hat{\varphi}_n)$ est bien défini avec une probabilité tendant vers 1. De plus comme φ^{-1} est continue en $\varphi(\theta)$, le théorème de continuité nous donne la consistance de $\hat{\theta}_n$.

Supposons maintenant que la normalité asymptotique est vérifiée, on va appliquer la delta méthode avec $g = \varphi^{-1}$. Rappelons que pour tout $x \in \varphi(\Theta)$,

$$\left(\varphi^{-1}\right)'(x) = \frac{1}{(\varphi' \circ \varphi^{-1})(x)},$$

et on obtient donc en particulier

$$\left(\varphi^{-1}\right)'(\varphi(\theta)) = \frac{1}{\varphi'(\varphi^{-1}(\varphi(\theta)))} = \frac{1}{\varphi'(\theta)}.$$

En appliquant la delta méthode, on obtient donc

$$\sqrt{n} \left(\varphi^{-1}(\hat{\varphi}_n) - \varphi^{-1}(\varphi(\theta)) \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{N}} \mathcal{N}\left(0, \left(\left(\varphi^{-1}\right)'(\varphi(\theta)) \right)^2 \sigma^2 \right),$$

et on obtient donc

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{1}{(\varphi'(\theta))^2} \sigma^2\right).$$

□

La méthode des moments consiste donc à trouver une fonction φ qui soit un difféomorphisme au voisinage de θ et un moment k tel que $\mathbb{E}[X_1^k] = \varphi(\theta)$. En effet, on a alors un estimateur "naturel" de $m_k = \mathbb{E}[X_1^k]$ donné par

$$\hat{m}_{n,k} = \frac{1}{n} \sum_{i=1}^n X_i^k$$

et on obtient alors l'estimateur $\hat{\theta}_n = \varphi^{-1}(\hat{m}_{n,k})$ de θ .

Exemple : la loi uniforme. Soit X_1, \dots, X_n des variables aléatoires i.i.d suivant une loi uniforme sur $[0, \theta^2]$, avec $\theta > 0$, et donc de densité

$$f_\theta = \frac{1}{\theta^2} \mathbb{1}_{[0, \theta^2]}.$$

On a $\mathbb{E}[X_1] = \frac{\theta^2}{2}$. Par la loi des grands nombres, \bar{X}_n converge en probabilité vers $\theta^2/2$. On pose alors $\varphi : x \mapsto x^2/2$ qui est bien un C^1 -difféomorphisme de \mathbb{R}_+^* dans \mathbb{R}_+^* , et d'inverse $\varphi^{-1} : x \mapsto \sqrt{2x}$. On obtient alors la convergence en probabilité de $\hat{\theta}_n = \sqrt{2\bar{X}_n}$ vers θ . De plus, comme $\mathbb{V}[X_1] = \frac{\theta^4}{12}$ grâce au TLC, on a

$$\sqrt{n} \left(\bar{X}_n - \frac{\theta^2}{2} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \frac{\theta^4}{12} \right).$$

Comme $\varphi'(\theta) = \theta$, on obtient la normalité asymptotique suivante

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \frac{\theta^2}{12} \right)$$

Exemple : la loi exponentielle. On considère X_1, \dots, X_n des variables aléatoires i.i.d suivant une loi exponentielle de paramètre $\theta \in \mathbb{R}_+^*$, i.e ayant une densité définie pour tout $x \in \mathbb{R}$ par

$$f_\theta(x) = \theta \exp(-\theta x) \mathbb{1}_{\mathbb{R}_+}(x).$$

Rappelons que $\mathbb{E}[X_1] = \theta^{-1}$ et $\mathbb{V}[X_1] = \theta^{-2}$. La loi des grands nombres nous donne la convergence en probabilité de \bar{X}_n vers θ^{-1} . On pose alors $\varphi : x \mapsto x^{-1}$ qui est un C^1 difféomorphisme de \mathbb{R}_+^* dans \mathbb{R}_+^* . On a donc la convergence en probabilité de $\hat{\theta}_n = \bar{X}_n^{-1}$ vers θ . De plus, grâce au TCL, on a

$$\sqrt{n} \left(\bar{X}_n - \theta^{-1} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, \frac{1}{\theta^2} \right).$$

Comme $\varphi'(\theta) = \frac{-1}{\theta^2}$, on obtient

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} (0, \theta^2).$$

Remarque : Pour éviter les erreurs, il est plus judicieux (et cela revient à peu près au même) d'écrire θ comme une fonction de $\mathbb{E}[X^k]$ et d'appliquer ensuite le théorème de continuité pour avoir la consistance, et la delta méthode pour obtenir la normalité asymptotique...

Remarque : Attention! Il peut arriver qu'une variable aléatoire n'admette pas de moment d'ordre 1. Il faut alors essayer d'être malin!

Exemple : Soit $\theta > 0$, on considère une variable aléatoire X de densité définie pour tout $x \in \mathbb{R}$ par

$$f_\theta(x) = \frac{\theta}{x^2} \mathbb{1}_{x \geq \theta}.$$

On a ainsi

$$\mathbb{E}[X] = \int_{\theta}^{+\infty} \frac{\theta}{x} dx = +\infty.$$

Cependant, on peut s'intéresser à

$$\mathbb{E}\left[\frac{1}{X}\right] = \int_{\theta}^{+\infty} \frac{\theta}{x^3} dx = \frac{1}{2\theta}.$$

Ainsi, comme $\theta = \frac{1}{2\mathbb{E}\left[\frac{1}{X}\right]}$, on va considérer l'estimateur $\hat{\theta}_n = \frac{n}{2\sum_{i=1}^n \frac{1}{X_i}}$. Par la loi des grands nombres, on obtient

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{X_i} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \frac{1}{2\theta}$$

De plus, la fonction $\varphi : x \mapsto \frac{1}{2x}$ est continue en $\frac{1}{2\theta} \neq 0$, et on obtient donc la consistance via le théorème de continuité. De plus,

$$\mathbb{E}\left[\frac{1}{X^2}\right] = \int_{\theta}^{+\infty} \frac{\theta}{x^4} dx = \frac{1}{3\theta^2}$$

et

$$\mathbb{V}\left[\frac{1}{X}\right] = \frac{1}{12\theta^2}.$$

Le TLC nous donne donc

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{X_i} - \frac{1}{2\theta} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{1}{12\theta^2}\right).$$

De plus, $\varphi'(x) = -\frac{1}{2x^2}$, et φ est donc dérivable en $\frac{1}{2\theta}$ avec $\varphi'\left(\frac{1}{2\theta}\right) = 2\theta^2$. On obtient donc, par la delta méthode,

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{\theta^2}{3}\right).$$

Remarque : De plus, même si les moments d'ordre 1, 2, ... existent, il n'est pas toujours facile de les calculer. Il peut alors être intéressant de calculer la fonction génératrice des moments.

Définition 2.3.1 (Fonction génératrice des moments). Soit X une variable aléatoire, on appelle fonction génératrice des moments G_X de X la fonction définie par

$$G_X(t) = \mathbb{E}\left[e^{tX}\right].$$

Attention, la fonction génératrice des moments n'est pas nécessairement définie, et encore moins pour tout t , mais le théorème suivant nous indique quand elle peut l'être.

Théorème 2.3.1. On suppose que la variable aléatoire X admet des moments de tout ordre, i.e pour tout $k \in \mathbb{N}$ $\mathbb{E}[X^k] < +\infty$, et que la série $\sum_{k \geq 0} \frac{t^k \mathbb{E}[X^k]}{k!}$ admet un rayon de convergence R . Alors pour tout

$|t| < R$ on a

$$G_X(t) = \sum_{k \geq 0} \frac{t^k \mathbb{E}[X^k]}{k!}.$$

En particulier, pour tout $k \in \mathbb{N}$, on a

$$G_X^{(k)}(0) = \mathbb{E}[X^k].$$

Démonstration. On a, grâce à la décomposition en série entière de la fonction exponentielle,

$$e^{tX} = \sum_{k \geq 0} \frac{t^k X^k}{k!}$$

et on a donc, pour tout t tel que $|t| < R$, à l'aide du théorème de Fubini-Tonelli,

$$\mathbb{E}[e^{tX}] = \mathbb{E}\left[\sum_{k \geq 0} \frac{t^k X^k}{k!}\right] = \sum_{k \geq 0} \frac{t^k \mathbb{E}[X^k]}{k!}.$$

En particulier, pour tout $k \geq 0$, on a

$$G_X^{(k)}(t) = \mathbb{E}[X^k] + \sum_{k' \geq k+1} \frac{k'(k'-1)\dots(k'-k+1)}{k!} t^{k'-k} \mathbb{E}[X^{k'}]$$

et en particulier, on obtient $G_X^{(k)}(0) = \mathbb{E}[X^k]$. □

Exemple : la loi géométrique. Soit X une variable aléatoire suivant une loi géométrique de paramètre $p \in (0, 1)$. On a alors pour tout $t \geq 0$, en faisant le changement de variable $k' = k - 1$,

$$G_X(t) = \mathbb{E}[e^{tX}] = \sum_{k \geq 1} e^{tk} (1-p)^{k-1} p = e^t p \sum_{k'=0}^{+\infty} e^{tk'} (1-p)^{k'} = e^t p \sum_{k=0}^{+\infty} ((1-p)e^t)^k.$$

En prenant t tel que $(1-p)e^t < 1$, i.e $t < -\log(1-p)$ on a donc

$$G_X(t) = \frac{e^t p}{1 - (1-p)e^t} = \frac{p}{e^{-t} - (1-p)}.$$

Ainsi, on a pour tout $t < -\log(1-p)$

$$G_X'(t) = p \frac{e^{-t}}{(e^{-t} - (1-p))^2}$$

et on a ainsi $\mathbb{E}[X] = G_X'(0) = \frac{p}{p^2} = \frac{1}{p}$. De plus, on a

$$G_X''(t) = p \frac{-e^{-t} (e^{-t} - (1-p))^2 + 2e^{-2t} (e^{-t} - (1-p))}{(e^{-t} - (1-p))^4} = p \frac{-e^{-t} (e^{-t} - (1-p)) + 2e^{-2t}}{(e^{-t} - (1-p))^3}$$

et on a ainsi $\mathbb{E}[X^2] = G_X''(0) = p \frac{2-p}{p^3} = \frac{2-p}{p^2}$. Ainsi, on a

$$\mathbb{V}[X] = \frac{2-p}{p^2} - \frac{1}{p^2} = \frac{1-p}{p^2}.$$

2.4 Méthode du maximum de vraisemblance

Dans ce qui suit, on considère des variables aléatoires indépendantes et identiquement distribuées X_1, \dots, X_n .

2.4.1 Cas discret

On note f_θ la densité de X_1 , i.e la fonction définie pour tout $x \in \mathbb{R}$ par $f_\theta(x) = \mathbb{P}[X_1 = x]$.

Definition 2.4.1 ((log)-Vraisemblance). La vraisemblance de $\mathbf{X} = (X_1, \dots, X_n)$ est définie pour tout $\theta \in \Theta$ par

$$L_{\mathbf{X}}(\theta) = \prod_{i=1}^n f_\theta(X_i).$$

La log-vraisemblance de \mathbf{X} est définie pour tout $\theta \in \Theta$ par

$$l_{\mathbf{X}}(\theta) = \log(L_{\mathbf{X}}(\theta)) = \log\left(\prod_{i=1}^n f_\theta(X_i)\right).$$

Attention! La (log)-vraisemblance est une variable aléatoire.

Exemple 1 : loi de Bernoulli. On considère $\mathbf{X} = (X_1, \dots, X_n)$ où les X_i sont i.i.d et suivent une loi de Bernoulli de paramètre $\theta \in (0, 1)$, i.e la densité de X_1 est définie pour tout $x \in \{0, 1\}$ par

$$f_\theta(x) = \theta^x (1 - \theta)^{1-x}.$$

On obtient donc pour tout $\theta \in (0, 1)$,

$$L_{\mathbf{X}}(\theta) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1-X_i} = \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{\sum_{i=1}^n (1-X_i)} = \theta^{n\bar{X}_n} (1 - \theta)^{n - n\bar{X}_n}.$$

De plus, on a

$$l_{\mathbf{X}}(\theta) = n\bar{X}_n \ln(\theta) + (n - n\bar{X}_n) \log(1 - \theta).$$

Exemple 2 : loi de Poisson. On considère $\mathbf{X} = (X_1, \dots, X_n)$ où les X_i sont i.i.d et suivent une loi de Poisson de paramètre $\theta \in \mathbb{R}_+^*$, i.e la densité de X_1 est définie pour tout $x \in \mathbb{N}$ par

$$f_\theta(x) = \frac{\theta^x}{x!} e^{-\theta}.$$

On obtient donc pour tout $\theta \in \Theta$

$$L_n(\theta) = \prod_{i=1}^n \frac{\theta^{X_i}}{X_i!} e^{-\theta} = e^{-\theta n} \prod_{i=1}^n \frac{\theta^{X_i}}{X_i!}$$

$$l_n(\theta) = \log \left(e^{-\theta n} \right) \sum_{i=1}^n \log \left(\frac{\theta^{X_i}}{X_i!} \right) = -\theta n + \sum_{i=1}^n X_i \log(\theta) - \sum_{i=1}^n \log(X_i!)$$

Si on note x_1, \dots, x_n les réalisations de X_1, \dots, X_n , la réalisation de la vraisemblance est la probabilité d'obtenir une telle réalisation de l'échantillon si θ est le vrai paramètre qui régit X_1 . L'objectif est de trouver θ qui va maximiser cette probabilité, ce qui conduit à l'estimateur du maximum de vraisemblance.

Definition 2.4.2 (Estimateur du maximum de vraisemblance (EMV)). *Le maximum de vraisemblance, si il existe, est un élément $\hat{\theta}_n$ de Θ qui vérifie*

$$L_{\mathbf{X}}(\hat{\theta}_n) = \sup_{\theta \in \Theta} L_{\mathbf{X}}(\theta).$$

De manière équivalente, l'estimateur du maximum de vraisemblance, si il existe, vérifie

$$l_{\mathbf{X}}(\hat{\theta}_n) = \sup_{\theta \in \Theta} l_{\mathbf{X}}(\theta).$$

On voit donc que la réalisation θ_n de l'estimateur $\hat{\theta}_n$ maximise la probabilité d'obtenir un tel échantillon.

Attention! Ni l'existence ni l'unicité de l'estimateur du maximum de vraisemblance ne sont assurés.

Exemple : loi de Bernoulli. On va chercher à maximiser la log-vraisemblance. En notant $l'_{\mathbf{X}}(\theta) = \frac{\partial}{\partial \theta} l_{\mathbf{X}}(\theta)$, on a

$$l'_{\mathbf{X}}(\theta) = \frac{n\bar{X}_n}{\theta} - \frac{n - n\bar{X}_n}{1 - \theta} = \frac{n\theta - n\bar{X}_n}{\theta(1 - \theta)} = 0 \Leftrightarrow \theta = \bar{X}_n.$$

De plus, en dressant le tableau de variation, on voit que \bar{X}_n est l'unique maximiseur de la log-vraisemblance, et on a donc l'EMV $\hat{\theta}_n^{MV} = \bar{X}_n$.

2.4.2 Cas continu

On suppose maintenant que la variable aléatoire X_1 est continue et admet une densité f_{θ} (par rapport à la mesure de Lebesgue). La définition de la (log)-vraisemblance reste inchangée.

Definition 2.4.3 ((log)-Vraisemblance). *La vraisemblance de $\mathbf{X} = (X_1, \dots, X_n)$ est définie pour tout $\theta \in \Theta$ par*

$$L_{\mathbf{X}}(\theta) = \prod_{i=1}^n f_{\theta}(X_i).$$

La log-vraisemblance de \mathbf{X} est définie pour tout $\theta \in \Theta$ par

$$l_{\mathbf{X}}(\theta) = \log(L_{\mathbf{X}}(\theta)) = \log\left(\prod_{i=1}^n f_{\theta}(X_i)\right).$$

La définition de l'estimateur du maximum de vraisemblance reste donc inchangée dans le cas continue. Cependant, on peut noter que dans ce cas, la réalisation de la vraisemblance est la densité en la réalisation, et la réalisation de l'estimateur du maximum de vraisemblance est donc le paramètre θ qui maximise la densité en l'échantillon.

Exemple 1 : loi normale. On considère $\mathbf{X} = (X_1, \dots, X_n)$ où les X_i sont i.i.d et suivent une loi normale d'espérance $\theta \in \mathbb{R}$ et de variance 1, i.e la densité de X_1 est définie pour tout $x \in \mathbb{R}$ par

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

On obtient donc pour tout $\theta \in \Theta$

$$\begin{aligned} L'_{\mathbf{X}}(\theta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(X_i-\theta)^2}{2}\right) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (X_i-\theta)^2\right) \\ l'_{\mathbf{X}}(\theta) &= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n (X_i-\theta)^2 \end{aligned}$$

En dérivant par rapport à θ , on obtient

$$l'_{\mathbf{X}}(\theta) = \sum_{i=1}^n \frac{X_i - \theta}{\sigma^2}.$$

On cherche alors les zéros de la dérivée, et en multipliant par σ^2 , on obtient

$$\begin{aligned} \sum_{i=1}^n \frac{X_i - \theta}{\sigma^2} = 0 &\Leftrightarrow \sum_{i=1}^n (X_i - \theta) = 0 \\ &\Leftrightarrow n\bar{X}_n - n\bar{\theta} = 0 \\ &\Leftrightarrow \theta = \bar{X}_n. \end{aligned}$$

Reste à montrer que $\hat{\theta}_n = \bar{X}_n$ est bien l'unique maximum de la log vraisemblance. On regarde donc la dérivée seconde par rapport à θ , et on obtient

$$l''_{\mathbf{X}}(\theta) = \frac{-n}{\sigma^2} < 0$$

et la log-vraisemblance est donc strictement concave, et $\hat{\theta}_n$ est son unique maximum. On obtient donc l'EMV

$$\hat{\theta}_n = \bar{X}_n.$$

Exemple 2 : On considère que X_1 admet une densité f_θ définie pour tout $x \in \mathbb{R}$ par

$$f_\theta(x) = \frac{\theta}{x^2} \mathbf{1}_{x \geq \theta}.$$

avec $\theta > 0$. On a alors la vraisemblance qui est définie pour tout $\theta > 0$ par

$$L_n(\theta) = \prod_{i=1}^n \frac{\theta^n}{X_i^2} \mathbf{1}_{X_i \geq \theta} = \theta^n \prod_{i=1}^n \frac{1}{X_i^2} \prod_{i=1}^n \mathbf{1}_{X_i \geq \theta}.$$

De plus

$$\prod_{i=1}^n \mathbf{1}_{X_i \geq \theta} = 1 \Leftrightarrow \forall i, X_i \geq \theta \Leftrightarrow \theta \leq X_{(1)} := \min_i X_i.$$

On a alors

$$L_n(\theta) = \theta^n \mathbf{1}_{\theta \leq X_{(1)}} \prod_{i=1}^n \frac{1}{X_i^2}$$

qui est nulle sur $(X_{(1)}, +\infty)$, et positive et croissante sur $(0, X_{(1)})$, et le maximum est donc atteint en $X_{(1)}$ et est unique. On a donc l'estimateur du maximum de vraisemblance $\hat{\theta}_n = X_{(1)}$. Il reste alors à démontrer la convergence de l'estimateur. Comme $\hat{\theta}_n$ ne peut pas être vu comme une fonction de \bar{X}_n , on ne peut pas utiliser le cheminement classique, LGN, théorème de continuité, TLC, delta méthode... Cependant, pour tout $\epsilon > 0$, comme $X_{(1)} \geq \theta$,

$$\mathbb{P} \left[\left| X_{(1)} - \theta \right| \geq \epsilon \right] = \mathbb{P} \left[X_{(1)} - \theta \geq \epsilon \right] = \mathbb{P} \left[X_{(1)} \geq \theta + \epsilon \right] = \mathbb{P} \left[\forall i, X_i \geq \theta + \epsilon \right].$$

De plus, comme les X_i sont indépendants et identiquement distribués, on obtient

$$\mathbb{P} \left[\left| X_{(1)} - \theta \right| \geq \epsilon \right] = \prod_{i=1}^n \mathbb{P} \left[X_i \geq \theta + \epsilon \right] = (\mathbb{P} \left[X_1 \geq \theta + \epsilon \right])^n = (1 - F_{X_1}(\theta + \epsilon))^n,$$

où F_{X_1} est la fonction de répartition de X_1 qu'il faut donc calculer. On a pour tout $x \geq \theta$

$$F_{X_1}(x) = \int_{\theta}^x \frac{\theta}{t^2} dt = \theta \left[\frac{-1}{t} \right]_{\theta}^x = 1 - \frac{\theta}{x}.$$

On obtient donc

$$\mathbb{P} \left[\left| X_{(1)} - \theta \right| \geq \epsilon \right] = \left(\frac{\theta}{\theta + \epsilon} \right)^n \xrightarrow{n \rightarrow +\infty} 0.$$

Reste encore à trouver la vitesse de convergence. Pour "plagier" le TLC, on va considérer une suite positive $v_n \rightarrow +\infty$ telle que

$$V_n := v_n \left(X_{(1)} - \theta \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Z$$

où Z suit une loi connue (si possible). Pour obtenir la convergence en loi, on va donc regarder la fonction de répartition de V_n . Comme $X_{(1)} \geq \theta$, pour tout $t < 0$, on a $\mathbb{P} \left[V_n \leq t \right] = 0$. Pour tout

$t \geq 0$, on a (en utilisant les calculs précédents)

$$\mathbb{P}[V_n \leq t] = \mathbb{P}\left[X_{(1)} - \theta \leq \frac{t}{v_n}\right] = 1 - \mathbb{P}\left[X_{(1)} - \theta \geq \frac{t}{v_n}\right] = 1 - \left(\frac{\theta}{\theta + t/v_n}\right)^n.$$

On obtient donc

$$\mathbb{P}[V_n \leq t] = 1 - \left(1 - \frac{t/v_n}{\theta + t/v_n}\right)^n = 1 - \left(1 - \frac{t}{\theta v_n + t}\right)^n = 1 - \exp\left(n \ln\left(1 - \frac{t}{\theta v_n + t}\right)\right) \sim 1 - \exp\left(-\frac{tn}{\theta v_n + t}\right)$$

En choisissant $v_n = n$, on a

$$\mathbb{P}[V_n \leq t] \xrightarrow[n \rightarrow +\infty]{} 1 - \exp\left(-\frac{t}{\theta}\right),$$

i.e on a la convergence en loi

$$n(X_{(1)} - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{E}\left(\frac{1}{\theta}\right).$$

2.5 Comparaison d'estimateurs

On a vu, dans les sections précédentes, différentes façons de construire des estimateurs. On s'intéresse ici à la façon de les comparer et donc comment choisir le meilleur (si il existe).

2.5.1 Comparaison des erreurs quadratiques moyennes

Une façon de quantifier la qualité d'un estimateur $\hat{\theta}_n$ de θ , est de considérer son risque quadratique

$$\text{EQM}(\hat{\theta}_n, \theta) = \mathbb{E}\left[(\hat{\theta}_n - \theta)^2\right].$$

On considère que $\hat{\theta}_n$ est un meilleur estimateur de θ que l'estimateur $\tilde{\theta}_n$ si

$$\forall \theta \in \Theta, \quad \text{EQM}(\hat{\theta}_n, \theta) \leq \text{EQM}(\tilde{\theta}_n, \theta).$$

Exemple : loi uniforme. Soit $\theta > 0$, on considère une variable aléatoire X suivant une loi uniforme sur $[0, \theta]$ et on se donne X_1, \dots, X_n de même loi que X . On a $\mathbb{E}[X] = \frac{\theta}{2}$ et donc, par la méthode des moments, on obtient l'estimateur

$$\hat{\theta}_n = 2\bar{X}_n.$$

qui est un estimateur sans biais de θ . De plus, comme $\mathbb{V}[X_1] = \frac{\theta^2}{12}$, on a, grâce à la décomposition biais variance,

$$\text{EQM}(\hat{\theta}_n, \theta) = \mathbb{V}[\hat{\theta}_n] = 4 \frac{\mathbb{V}[X_1]}{n} = \frac{\theta^2}{3n}.$$

On s'intéresse maintenant à l'estimateur du maximum de vraisemblance. On a la vraisemblance

qui est définie pour tout $\theta > 0$ par

$$L_{\mathbf{X}}(\theta) = \prod_{i=1}^n \frac{1}{\theta} \mathbf{1}_{[0,\theta]}(X_i) = \frac{1}{\theta^n} \prod_{i=1}^n \mathbf{1}_{\theta \geq X_i}.$$

Or, $\prod_{i=1}^n \mathbf{1}_{\theta \geq X_i} = 1 \Leftrightarrow \forall i, X_i \leq \theta \Leftrightarrow X_{(n)} \leq \theta$. On a donc

$$L_{\mathbf{X}}(\theta) = \frac{1}{\theta^n} \mathbf{1}_{\theta \geq X_{(n)}}$$

qui est nulle sur $(0, X_{(n)})$ et positive et décroissante sur $[X_{(n)}, +\infty[$, et son maximum est donc atteint en $X_{(n)}$ qui est donc l'estimateur du maximum de vraisemblance. De plus, pour tout $t \in [0, \theta]$,

$$\begin{aligned} F_{X_{(n)}}(t) &= \mathbb{P} [X_{(n)} \leq t] \\ &= \mathbb{P} [\{\forall i, X_i \leq t\}] \\ &= \prod_{i=1}^n \mathbb{P} [X_i \leq t] \\ &= (\mathbb{P} [X_1 \leq t])^n \\ &= \left(\frac{t}{\theta}\right)^n \end{aligned}$$

On obtient donc

$$f_{X_{(n)}}(t) = \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1} \mathbf{1}_{[0,\theta]}(t).$$

Ainsi,

$$\mathbb{E} [\hat{\theta}_n] = \int_0^\theta t \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1} dt = \frac{n}{n+1} \theta \int_0^\theta \frac{n+1}{\theta} \left(\frac{t}{\theta}\right)^n dt = \frac{n}{n+1} \theta$$

et l'estimateur est donc biaisé, avec en particulier

$$B(\hat{\theta}_n, \theta) = -\frac{1}{n+1} \theta$$

De plus, on a

$$\mathbb{E} [\hat{\theta}_n^2] = \int_0^\theta t^2 \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1} dt = \frac{n}{n+2} \theta^2 \int_0^\theta \frac{n+2}{\theta} \left(\frac{t}{\theta}\right)^{n+1} dt = \frac{n}{n+2} \theta^2$$

On obtient donc

$$\begin{aligned}\mathbb{V} [\hat{\theta}_n] &= \mathbb{E} [\hat{\theta}_n^2] - (\mathbb{E} [\hat{\theta}_n])^2 \\ &= \frac{n}{n+2}\theta^2 - \frac{n^2}{(n+1)^2}\theta^2 \\ &= \frac{n(n+1)^2 - n^2(n+2)}{(n+2)(n+1)^2}\theta^2 \\ &= \frac{n}{(n+2)(n+1)^2}\theta^2.\end{aligned}$$

On a alors,

$$\text{EQM} (\hat{\theta}_n, \theta) = \mathbb{V} [\hat{\theta}_n] + B (\hat{\theta}_n, \theta)^2 = \frac{n}{(n+2)(n+1)^2}\theta^2 + \frac{n^2}{(n+1)^2}\theta^2 = \frac{2\theta^2}{(n+1)(n+2)}$$

Et c'est donc un meilleur estimateur que l'estimateur des moments dès que $\frac{2}{(n+1)(n+2)} \leq \frac{1}{3n}$, i.e dès que $n \geq 2$.

2.5.2 Biais d'un estimateur

On donne souvent trop d'importance à l'absence de biais d'un estimateur, voire même on peut penser qu'un estimateur sera meilleur si il est non biaisé, ce qui est faux! Débiaiser un estimateur ne donne pas forcément de meilleurs résultats!

Exemple : loi uniforme. Soit $\theta > 0$, on considère une variable aléatoire X suivant une loi uniforme sur $[0, \theta]$. On rappelle que l'estimateur du maximum de vraisemblance est $X_{(n)}$ et qu'il est biaisé avec

$$B (\hat{\theta}_n, \theta) = -\frac{1}{n+1}\theta.$$

On considère maintenant l'estimateur non biaisé $\tilde{\theta}_n$ défini par

$$\tilde{\theta}_n = \frac{n+1}{n}X_{(n)}.$$

On a

$$\mathbb{V} [\tilde{\theta}_n] = \frac{(n+1)^2}{n^2}\mathbb{V} [X_{(n)}] = \frac{\theta^2}{n(n+2)} = \text{EQM} (\tilde{\theta}_n, \theta),$$

et on obtient donc un meilleur estimateur que $X_{(n)}$. Cependant, si on considère un estimateur de la forme $\alpha X_{(n)}$ avec α un réel positif, on obtient

$$\mathbb{E} [\alpha X_{(n)}] = \alpha \frac{n}{n+1}\theta \quad \text{et} \quad B (\alpha X_{(n)}, \theta) = \frac{\alpha n - n - 1}{n+1}\theta$$

et

$$\begin{aligned} \text{EQM}(\alpha X_n, \theta) &= \alpha^2 \mathbb{E} \left[X_{(n)}^2 \right] - 2\alpha\theta \mathbb{E} \left[X_{(n)} \right] + \theta^2 \\ &= \theta^2 \left(\frac{n}{n+2} \alpha^2 - \frac{2n}{n+1} \alpha + 1 \right) \end{aligned}$$

On prend alors α qui minimise la quantité précédente, i.e $\alpha = \frac{n+2}{n+1}$ et on obtient ainsi un meilleur estimateur.

2.5.3 L'approche asymptotique

On a vu qu'il n'était pas toujours évident (voire possible) de calculer les erreurs quadratiques moyennes. De plus, il est relativement atypique d'obtenir des vitesses de convergence en moyenne quadratique de l'ordre de $1/n^2$. Très souvent, on dispose d'estimateurs $\hat{\theta}_n, \tilde{\theta}_n$ asymptotiquement normaux, i.e il existe σ_1^2, σ_2^2 telles que

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma_1^2) \quad \text{et} \quad \sqrt{n} (\tilde{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma_2^2)$$

Si on a $\sigma_1^2 \leq \sigma_2^2$, on choisira alors l'estimateur $\hat{\theta}_n$. En effet, on verra (chapitre suivant) que l'on a alors, par exemple,

$$\mathbb{P} \left[\theta \in \left[\hat{\theta} \pm 1.96 \frac{\sigma_1}{\sqrt{n}} \right] \right] \xrightarrow[n \rightarrow +\infty]{} 0.95 \quad \text{et} \quad \mathbb{P} \left[\theta \in \left[\tilde{\theta} \pm 1.96 \frac{\sigma_2}{\sqrt{n}} \right] \right] \xrightarrow[n \rightarrow +\infty]{} 0.95$$

et le premier intervalle est alors plus précis. En effet, les deux intervalles sont alors de longueurs

$$2 \times 1.96 \frac{\sigma_1}{\sqrt{n}} \quad \text{et} \quad 2 \times 1.96 \frac{\sigma_2}{\sqrt{n}}.$$

Chapitre 3

Intervalle de confiance

Dans ce chapitre, on considère des variables aléatoires i.i.d X_1, \dots, X_n suivant une loi F . On s'intéresse à l'estimation d'une caractéristique ou d'un paramètre θ de la loi F . On a vu dans les chapitres précédents comment obtenir et comparer des estimateurs $\hat{\theta}_n$ de θ . On s'intéresse ici à la quantification de la confiance que l'on peut accorder à l'estimation, i.e à l'obtention d'intervalles de confiance.

3.1 Intervalle de confiance

3.1.1 Intervalle de confiance

Soit $\alpha \in (0, 1)$, un intervalle de confiance pour le paramètre θ au niveau de confiance $1 - \alpha$ est un intervalle de la forme

$$IC_{1-\alpha}(\theta) = [a(X_1, \dots, X_n), b(X_1, \dots, X_n)]$$

avec

$$\mathbb{P}[\theta \in [a(X_1, \dots, X_n), b(X_1, \dots, X_n)]] = 1 - \alpha.$$

On parle alors d'intervalles de confiance bilatères.

Attention! Cela ne signifie pas que θ appartienne à $IC_{1-\alpha}(\theta)$. On peut seulement affirmer que $\theta \in IC_{1-\alpha}(\theta)$ en prenant un risque α .

Attention! On ne peut en aucun cas dire que la probabilité que θ appartienne à la réalisation de $IC_{1-\alpha}(\theta)$ est égale à $1 - \alpha$. En effet, la réalisation de cet intervalle n'est en aucun cas aléatoire, et donc la probabilité que θ appartienne à cet intervalle est soit 0 soit 1.

Exemple 1 : On considère des variables aléatoires X_1, \dots, X_n i.i.d de densité f_θ définie pour tout $x \in \mathbb{R}$ par

$$f_\theta(x) = \frac{\theta}{x^2} \mathbf{1}_{x \geq \theta},$$

avec $\theta \in \mathbb{R}$. On a vu que l'estimateur du maximum de vraisemblance est défini par $X_{(1)}$. Comme

$\theta \leq X_{(1)}$, on va chercher un intervalle de la forme

$$IC_{1-\alpha}(\theta) = [X_{(1)} - c_\alpha, X_{(1)}].$$

Rappelons que l'on a vu que pour tout $t \geq 0$,

$$\mathbb{P} [X_{(1)} - \theta \geq t] = \left(\frac{\theta}{\theta + t} \right)^n.$$

De plus, on cherche c_α tel que

$$\begin{aligned} 1 - \alpha &= \mathbb{P} [X_{(1)} - c_\alpha \leq \theta \leq X_{(1)}] = \mathbb{P} [X_{(1)} - c_\alpha \leq \theta] \Leftrightarrow \alpha = \mathbb{P} [X_{(1)} - \theta \geq c_\alpha] \\ &\Leftrightarrow \alpha = \left(\frac{\theta}{\theta + c_\alpha} \right)^n \\ &\Leftrightarrow \alpha^{1/n} = \frac{\theta}{\theta + c_\alpha} \\ &\Leftrightarrow c_\alpha = \frac{\theta}{\alpha^{1/n}} - \theta = \theta (\alpha^{-1/n} - 1). \end{aligned}$$

Si on se contente de cela, on a un intervalle de la forme

$$IC_{1-\alpha}(\theta) = [X_{(1)} - \theta (\alpha^{-1/n} - 1); X_{(1)}],$$

donc dépendant de θ , et on ne peut donc pas le calculer. Il faut donc continuer de travailler. On a

$$\alpha = \mathbb{P} [X_{(1)} - \theta \geq \theta (\alpha^{-1/n} - 1)] = \mathbb{P} [X_{(1)} \geq \theta \alpha^{-1/n}] = \mathbb{P} [X_{(1)} \alpha^{1/n} \geq \theta]$$

et on obtient donc l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = [X_{(1)} \alpha^{1/n}; X_{(1)}].$$

Pour éviter les "subtilités" de la fin, on aurait pu se douter, vu la fonction de répartition de $X_{(1)}$ que l'on allait rencontrer des problèmes et donc directement regarder des intervalles de la forme

$$IC_{1-\alpha}(\theta) = [X_{(1)} c_\alpha; X_{(1)}].$$

On résout alors

$$\begin{aligned} 1 - \alpha &= \mathbb{P} [X_{(1)} c_\alpha \leq \theta] = \mathbb{P} \left[X_{(1)} \leq \frac{\theta}{c_\alpha} \right] \Leftrightarrow \alpha = \mathbb{P} \left[X_{(1)} \geq \frac{\theta}{c_\alpha} \right] = \mathbb{P} \left[X_{(1)} - \theta \geq \frac{\theta}{c_\alpha} - \theta \right] \\ &\Leftrightarrow \alpha = \left(\frac{\theta}{\theta/c_\alpha} \right)^n = c_\alpha^n \\ &\Leftrightarrow c_\alpha = \alpha^{1/n}, \end{aligned}$$

et on retrouve ainsi l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[X_{(1)}\alpha^{1/n}; X_{(1)} \right].$$

Exemple 2 : la loi uniforme. On considère des variables aléatoires i.i.d avec $X_i \sim \mathcal{U}([0, \theta])$. On sait que l'estimateur du maximum de vraisemblance est $X_{(n)}$. Comme $X_{(n)} \leq \theta$, on va chercher un intervalle de confiance de la forme $[X_{(n)}, X_{(n)}c_\alpha]$. On a pour tout $t \in [0, \theta]$,

$$\mathbb{P} \left[X_{(n)} \leq t \right] = \mathbb{P} [X_1 \leq t]^n = \left(\frac{t}{\theta} \right)^n$$

On cherche donc c_α tel que

$$\begin{aligned} 1 - \alpha = \mathbb{P} \left[X_{(n)} \leq \theta \leq X_{(n)}c_\alpha \right] &= \mathbb{P} \left[\theta \leq X_{(n)}c_\alpha \right] \Leftrightarrow 1 - \alpha = 1 - \left(\frac{1}{c_\alpha} \right)^n \\ &\Leftrightarrow \alpha^{-1/n} = c_\alpha \end{aligned}$$

On obtient donc l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[X_{(n)}; X_{(n)}\alpha^{-1/n} \right].$$

Remarque : On cherche souvent des intervalles de confiance tels que

$$\alpha/2 = \mathbb{P} [\theta \leq a(X_1, \dots, X_n)] = \mathbb{P} [\theta \geq b(X_1, \dots, X_n)].$$

Dans le cas de l'exemple 1, on a vu que pour tout $t \in (0, 1)$, on a $\mathbb{P} [\theta \leq X_{(1)}t^{1/n}] = t$. Ainsi on obtient

$$\mathbb{P} \left[\theta \leq X_{(1)} \left(\frac{\alpha}{2} \right)^{1/n} \right] = \alpha/2 \quad \text{et} \quad \mathbb{P} \left[\theta \geq X_{(1)} \left(1 - \frac{\alpha}{2} \right)^{1/n} \right] = \alpha/2.$$

On obtient donc, pour $\alpha \in (0, 1)$, l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[X_{(1)} \left(\frac{\alpha}{2} \right)^{1/n}; X_{(1)} \left(1 - \frac{\alpha}{2} \right)^{1/n} \right].$$

Dans le cas de l'exemple 2, on a vu que pour tout $t \in (0, 1)$, on a $\mathbb{P} [\theta \geq X_{(n)}t^{-1/n}] = t$. On a donc en particulier

$$\mathbb{P} \left[\theta \geq X_{(n)} \left(\frac{\alpha}{2} \right)^{-1/n} \right] = \alpha/2 \quad \text{et} \quad \mathbb{P} \left[\theta \leq X_{(n)} \left(1 - \frac{\alpha}{2} \right)^{-1/n} \right] = \alpha/2.$$

On obtient donc, pour $\alpha \in (0, 1/2)$, l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[X_{(n)} \left(1 - \frac{\alpha}{2}\right)^{-1/n}; X_{(n)} \left(\frac{\alpha}{2}\right)^{-1/n} \right].$$

Il n'est malheureusement pas toujours facile de calculer la fonction de répartition de l'estimateur et ainsi construire les intervalles de confiance. On peut par contre utiliser des inégalités usuelles (Markov, Bienaymé-Tchebycheff, Chernoff...) pour obtenir des intervalles de confiance de niveau au moins $1 - \alpha$, i.e des intervalles de la forme

$$IC_{1-\alpha}(\theta) = [a(X_1, \dots, X_n), b(X_1, \dots, X_n)]$$

avec

$$\mathbb{P}[\theta \in [a(X_1, \dots, X_n), b(X_1, \dots, X_n)]] \geq 1 - \alpha.$$

Exemple 1 : loi de Bernoulli. On considère des variables aléatoires i.i.d X_i suivant une loi de Bernoulli de paramètre θ . On rappelle que l'on a l'estimateur \bar{X}_n , et on cherche c_α tel que

$$\mathbb{P}[\bar{X}_n - c_\alpha \leq \theta \leq \bar{X}_n + c_\alpha] \geq 1 - \alpha.$$

En utilisant l'inégalité de Bienaymé-Tchebycheff, on a

$$\mathbb{P}[\bar{X}_n - c_\alpha \leq \theta \leq \bar{X}_n + c_\alpha] = \mathbb{P}[|\bar{X}_n - \theta| \leq c_\alpha] \geq 1 - \frac{\theta(1-\theta)}{nc_\alpha^2}$$

On peut alors chercher à résoudre

$$\alpha = \frac{\theta(1-\theta)}{nc_\alpha^2} \Leftrightarrow c_\alpha = \sqrt{\frac{\theta(1-\theta)}{n\alpha}}$$

et on obtiendrait l'intervalle suivant

$$IC_{1-\alpha}(\theta) = \left[\bar{X}_n \pm \sqrt{\frac{\theta(1-\theta)}{n\alpha}} \right].$$

Cependant, les bornes de l'intervalle dépendent de θ et ne sont donc pas calculables. Cependant, comme $\theta(1-\theta) \leq 1/4$, on obtient

$$\mathbb{P}[|\bar{X}_n - \theta| \geq c_\alpha] \leq \frac{1}{4nc_\alpha^2}$$

et donc l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[\bar{X}_n \pm \frac{1}{2\sqrt{n\alpha}} \right].$$

Exemple 2 : loi uniforme. On considère des variables aléatoires indépendantes et identiquement distribuées X_1, \dots, X_n avec $X_1 \sim \mathcal{U}([0, \theta])$, mais on considère cette fois-ci l'estimateur obtenu via la méthode des moments, i.e $\hat{\theta}_n = 2\bar{X}_n$. On cherche c_α tel que

$$\mathbb{P} [\theta \in [\hat{\theta}_n \pm c_\alpha]]$$

ce qui revient à chercher c_α tel que

$$\mathbb{P} [|\hat{\theta}_n - \theta| \geq c_\alpha] \leq \alpha.$$

En utilisant l'inégalité de Bienaymé-Tchebycheff, on obtient

$$\mathbb{P} [|\hat{\theta}_n - \theta| \geq c_\alpha] \leq \frac{\mathbb{V}[\hat{\theta}_n]}{c_\alpha^2} = \frac{\theta^2}{3nc_\alpha^2} = \alpha \Leftrightarrow c_\alpha = \frac{\theta}{\sqrt{3n\alpha}}.$$

On obtient donc

$$IC_{1-\alpha}(\theta) = \left[\hat{\theta}_n \pm \frac{\theta}{\sqrt{3n\alpha}} \right]$$

mais cet intervalle est inexploitable car dépendant de θ .

Remarque : On peut également chercher des intervalles de confiance de la forme

$$IC_{1-\alpha}(\theta) = [a(X_1, \dots, X_n), +\infty[\quad \text{ou} \quad IC_{1-\alpha}(\theta) =]-\infty, b(X_1, \dots, X_n)].$$

On parle alors d'intervalles de confiance unilatères.

3.1.2 Notion de quantile

On considère une variable aléatoire X et on note F sa fonction de répartition.

Definition 3.1.1. Pour tout $\alpha \in (0, 1)$, on appelle quantile d'ordre α le réel q_α tel que

$$q_\alpha = \inf \{x \in \mathbb{R}, F(x) \geq \alpha\}.$$

Lorsque la fonction de répartition de X est strictement croissante, alors cette dernière est bijective et le quantile q_α est alors la solution de

$$F(q_\alpha) = \alpha \Leftrightarrow q_\alpha = F^{-1}(\alpha).$$

Exemple 1 : loi uniforme. Soit $X \sim \mathcal{U}([a, b])$ et $\alpha \in (0, 1)$, alors le quantile q_α d'ordre α de X est défini par

$$q_\alpha = a + (b - a)\alpha$$

En effet, on a pour tout $t \in [a, b]$, $F_X(t) = \frac{t-a}{b-a}$. On résout donc

$$\alpha = \frac{t-a}{b-a} \Leftrightarrow t = a + (b-a)\alpha.$$

Exemple 2 : loi exponentielle. Soit $X \sim \mathcal{E}(1)$ et $\alpha \in (0, 1)$, alors le quantile q_α d'ordre α de X est défini par

$$q_\alpha = -\ln(1 - \alpha)$$

En effet, on a pour tout $t \geq 0$, $F_X(t) = 1 - e^{-t}$ et on résout donc

$$\alpha = 1 - e^{-t} \Leftrightarrow e^{-t} = 1 - \alpha \Leftrightarrow t = -\ln(1 - \alpha).$$

Exemple 3 : loi de Bernoulli. Soit $X \sim \mathcal{B}(\theta)$, alors le quantile q_α d'ordre α de X est défini par

$$q_\alpha = \begin{cases} 0 & \text{si } \alpha \in (0, 1 - \theta] \\ 1 & \text{sinon} \end{cases} \quad (3.1)$$

En effet, si $\alpha \leq 1 - \theta$, la plus petite valeur x telle que $F(x) \geq \alpha$ est 0 (et on a alors $F(0) = 1 - \theta$). De la même façon, si $\alpha \in (1 - \theta, 1]$, la plus petite valeur x telle que $F(x) \geq \alpha$ est 1 (et on a alors $F(1) = 1$).

Lorsque la fonction inverse est explicite (loi exponentielle, loi de Weibull, loi uniforme,...) on peut donc très facilement calculer les quantiles. Plus généralement, l'inverse de F n'est pas explicitement calculable et on utilise alors des tables statistiques pour retrouver le quantile cherché (loi normale, loi de Student, loi du Chi deux, loi de Fisher,...).

3.2 Rappels sur la loi normale

Dans ce qui suit, on considère des variables aléatoires X_1, \dots, X_n suivant des lois normales d'espérances μ_i et de variance σ_i^2 . Rappelons que la fonction caractéristique d'une variable aléatoire X suivant une loi normale d'espérance μ et de variance σ^2 est définie pour tout t par

$$\Phi_X(t) = \exp\left(\mu it - \frac{t^2 \sigma^2}{2}\right).$$

Proposition 3.2.1. *Si les X_i sont indépendants, alors toute combinaison linéaire des X_i suit une loi normale. Plus précisément, soit $\lambda_1, \dots, \lambda_n \in \mathbb{R}$, alors*

$$\sum_{i=1}^n \lambda_i X_i \sim \mathcal{N}(\mu, \sigma^2),$$

avec

$$\mu = \sum_{i=1}^n \lambda_i \mu_i \quad \text{et} \quad \sigma^2 = \sum_{i=1}^n \lambda_i^2 \sigma_i^2.$$

Démonstration. On note $X = \sum_{i=1}^n \lambda_i X_i$. Sa fonction caractéristique est définie pour tout t par (par indépendance des X_i)

$$\Phi_X(t) = \Phi_{\sum_{i=1}^n \lambda_i X_i}(t) = \prod_{i=1}^n \Phi_{\lambda_i X_i}(t) = \prod_{i=1}^n \Phi_{X_i}(\lambda_i t)$$

On obtient donc

$$\Phi_X(t) = \prod_{i=1}^n \left(\exp \left(i \mu_i \lambda_i t - \frac{\sigma_i^2 \lambda_i^2 t^2}{2} \right) \right) = \exp \left(i t \sum_{i=1}^n \lambda_i \mu_i - \frac{t^2 \sum_{i=1}^n \lambda_i^2 \sigma_i^2}{2} \right),$$

qui est la fonction caractéristique d'une loi normale de paramètres $\mu = \sum_{i=1}^n \lambda_i \mu_i$ et $\sigma^2 = \sum_{i=1}^n \lambda_i^2 \sigma_i^2$. \square

Definition 3.2.1 (Loi du Chi-deux). Soient X_1, \dots, X_n des variables aléatoires indépendantes suivant une loi normale centrée réduite, alors la variable aléatoire

$$Z_n = \sum_{i=1}^n X_i^2$$

suit une loi du Chi deux à n degrés de liberté (χ_n^2).

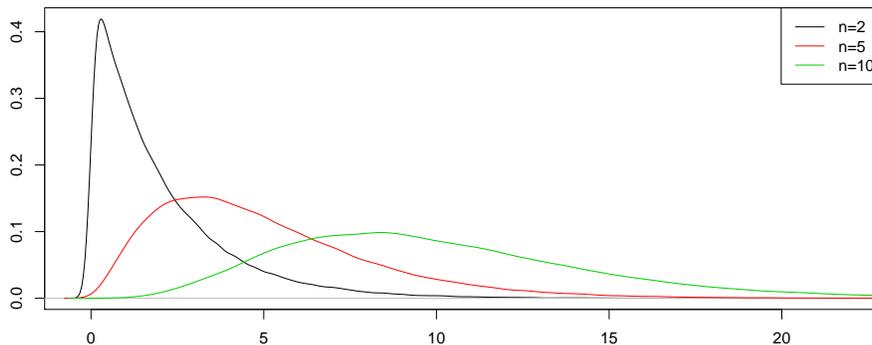


FIGURE 3.1 – Densité d'une chi deux à $n = 2, 5, 10$ degrés de liberté

Afin d'obtenir des intervalles de confiance pour l'estimateur de la moyenne, on définit maintenant la loi de Student.

Definition 3.2.2 (Loi de Student). Soit Z, U deux variables aléatoires indépendantes telles que $Z \sim \mathcal{N}(0, 1)$ et $U \sim \chi_n^2$, alors

$$\frac{Z}{\sqrt{U/n}} \sim T_n,$$

où T_n est une loi de Student à n degrés de liberté.

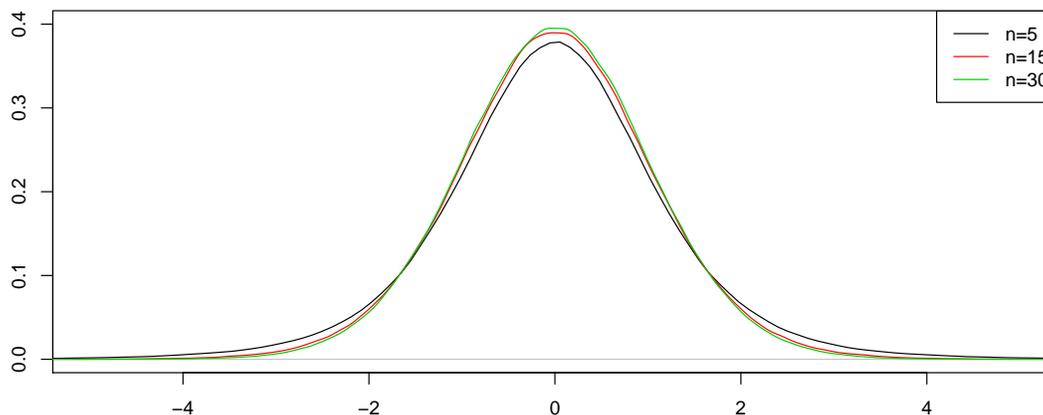


FIGURE 3.2 – Densité d'une loi de Student à $n = 5, 15, 30$ degrés de liberté.

Proposition 3.2.2. Soit T_n une loi de Student à n degrés de liberté. Si $n \geq 2$, alors T_n admet un moment d'ordre 1 et $\mathbb{E}[T_n] = 0$. De plus, la loi de Student est symétrique en 0. Enfin,

$$T_n \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

Démonstration. A noter que si $n = 1$, on a alors, par indépendance,

$$\mathbb{E}[|T_1|] = \mathbb{E}[|Z|] \mathbb{E}\left[\frac{1}{|X_1|}\right] = \mathbb{E}[|Z|] \int_{\mathbb{R}} \frac{1}{|x|} e^{-\frac{x^2}{2}} dx = +\infty.$$

On admettra que pour tout $n \geq 2$, on a

$$\mathbb{E}\left[\frac{1}{\sqrt{\sum_{i=1}^n X_i^2}}\right] < +\infty.$$

Le moment d'ordre 1 de T_n existe alors (cf calculs précédents) et

$$\mathbb{E}[T_n] = \mathbb{E}[Z] \mathbb{E}\left[\frac{\sqrt{n}}{\sqrt{\sum_{i=1}^n X_i^2}}\right] = 0.$$

De plus, pour tout $t \in \mathbb{R}$, on a par symétrie de la loi normale

$$\mathbb{P}[T_n \leq t] = \mathbb{P}\left[\frac{Z}{\sqrt{U/n}} \leq t\right] = \mathbb{P}\left[Z \leq t\sqrt{U/n}\right] = \mathbb{P}\left[Z \geq -t\sqrt{U/n}\right] = \mathbb{P}\left[\frac{Z}{\sqrt{U/n}} \geq -t\right].$$

Enfin, par la loi des grands nombres, on a

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mathbb{E}[X_1^2] = 1$$

et par le théorème de continuité, on a donc

$$\frac{1}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 1.$$

Le théorème de Slutsky nous donne donc

$$\frac{Z}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

□

3.3 Intervalles de confiance dans le cas gaussien

3.3.1 Estimation de la moyenne lorsque la variance est connue

On considère dans ce qui suit des variables aléatoires X_1, \dots, X_n i.i.d suivant une loi normale d'espérance μ inconnue et de variance σ^2 connue.

Proposition 3.3.1. *On a*

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

Démonstration. Voir TD.

□

Corollaire 3.3.1. *Pour tout $\alpha \in (0, 1)$,*

$$\mathbb{P}\left[\bar{X}_n - q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right] = 1 - \alpha,$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite, i.e si $Z \sim \mathcal{N}(0, 1)$,

$$\mathbb{P}[Z \leq q_{1-\alpha/2}] = 1 - \alpha/2.$$

On obtient donc l'intervalle de confiance

$$IC_{1-\alpha}(\mu) = \left[\bar{X}_n - q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}; \bar{X}_n + q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right].$$

Démonstration. Comme $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$, en centrant et en réduisant, on obtient

$$\frac{\sqrt{n}}{\sigma} (\bar{X}_n - \mu) \sim \mathcal{N}(0, 1).$$

De plus, soit $Z \sim \mathcal{N}(0, 1)$, par définition de $q_{1-\alpha/2}$, et par symétrie, on a

$$\mathbb{P}[Z \leq -q_{1-\alpha/2}] = \mathbb{P}[Z \geq q_{1-\alpha/2}] = 1 - \mathbb{P}[Z \leq q_{1-\alpha/2}] = \frac{\alpha}{2},$$

i.e $q_{\alpha/2} = -q_{1-\alpha/2}$. On a donc

$$\begin{aligned} \mathbb{P}[|Z| \leq q_{1-\alpha/2}] &= \mathbb{P}[Z \leq q_{1-\alpha/2}] - \mathbb{P}[Z \leq -q_{1-\alpha/2}] \\ &= \mathbb{P}[Z \leq q_{1-\alpha/2}] - \mathbb{P}[Z \leq q_{\alpha/2}] \\ &= 1 - \alpha. \end{aligned}$$

On obtient donc

$$\begin{aligned} 1 - \alpha &= \mathbb{P}\left[\frac{\sqrt{n}}{\sigma} |\bar{X}_n - \mu| \leq q_{1-\alpha/2}\right] \\ &= \mathbb{P}\left[-q_{1-\alpha/2} \leq \frac{\sqrt{n}}{\sigma} (\mu - \bar{X}_n) \leq q_{1-\alpha/2}\right] \\ &= \mathbb{P}\left[-q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu - \bar{X}_n \leq q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right] \\ &= \mathbb{P}\left[\bar{X}_n - q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right]. \end{aligned}$$

□

Remarque : On peut également trouver les intervalles de confiances unilatères suivant :

$$IC_{1-\alpha}(\mu) = \left[\bar{X}_n - q_{1-\alpha} \frac{\sigma}{\sqrt{n}}, +\infty\right] \quad \text{et} \quad IC_{1-\alpha}(\mu) = \left[-\infty, \bar{X}_n + q_{1-\alpha} \frac{\sigma}{\sqrt{n}}\right]$$

En effet, par définition de $q_{1-\alpha}$, on a

$$1 - \alpha = \mathbb{P}\left[\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \leq q_{1-\alpha}\right] = \mathbb{P}\left[-\mu \leq -\bar{X}_n + \frac{\sigma}{\sqrt{n}} q_{1-\alpha}\right] = \mathbb{P}\left[\mu \geq \bar{X}_n - q_{1-\alpha} \frac{\sigma}{\sqrt{n}}\right].$$

Pour trouver le dernier intervalle, il suffit de remarquer que $\sqrt{n} \frac{\mu - \bar{X}_n}{\sigma} \sim \mathcal{N}(0, 1)$, et donc

$$1 - \alpha = \mathbb{P}\left[\sqrt{n} \frac{\mu - \bar{X}_n}{\sigma} \leq q_{1-\alpha}\right].$$

3.3.2 Estimation de la moyenne lorsque la variance est inconnue

On considère maintenant que la variance est inconnue. La proposition suivante ainsi que son corollaire sont cruciaux pour obtenir les intervalles de confiance pour la moyenne.

Proposition 3.3.2 (Admis pour le moment). Soient $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ et $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. Alors

1. $\frac{n-1}{\sigma^2} S_n^2 \sim \chi_{n-1}^2$.
2. S_n et \bar{X}_n sont indépendants.

Le corollaire suivant est crucial pour obtenir des intervalles de confiance pour la moyenne dans le cas gaussien, ainsi que pour construire des tests.

Corollaire 3.3.2. On a

$$\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim T_{n-1},$$

où T_{n-1} est une loi de Student à $n - 1$ degrés de liberté.

Démonstration. On a la décomposition

$$\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim T_{n-1} = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \frac{\sigma}{S_n} = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \frac{1}{\sqrt{\frac{(n-1)S_n^2}{(n-1)\sigma^2}}},$$

et on conclut grâce à la proposition précédente. □

Le corollaire suivant donne les intervalles de confiances pour la moyenne.

Corollaire 3.3.3 (Intervalles de confiance). Soit $\alpha \in (0, 1)$, alors

$$\mathbb{P} \left[\bar{X}_n - t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \leq \mu \leq \bar{X}_n + t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \right] = 1 - \alpha,$$

où $t_{n-1, 1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi de Student à $n - 1$ degrés de liberté, i.e si $T \sim T_{n-1}$,

$$\mathbb{P} [T \leq t_{n-1, 1-\alpha/2}] = 1 - \alpha/2.$$

On obtient donc l'intervalle de confiance

$$IC_{1-\alpha}(\mu) = \left[\bar{X}_n - t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}}; \bar{X}_n + t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \right].$$

Démonstration. Comme la loi de Student est symétrique, on a, si $T \sim T_{n-1}$,

$$\mathbb{P} [|T| \leq t_{n-1, 1-\alpha/2}] = 1 - \alpha.$$

Ainsi,

$$\begin{aligned}
 1 - \alpha &= \mathbb{P} \left[\left| \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \right| \leq t_{n-1, 1-\alpha/2} \right] \\
 &= \mathbb{P} \left[-t_{n-1, 1-\alpha/2} \leq \sqrt{n} \frac{\mu - \bar{X}_n}{S_n} \leq t_{n-1, 1-\alpha/2} \right] \\
 &= \mathbb{P} \left[-t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \leq \mu - \bar{X}_n \leq t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \right] \\
 &= \mathbb{P} \left[\bar{X}_n - t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \leq \mu \leq \bar{X}_n + t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \right].
 \end{aligned}$$

□

Remarque : Anoter que l'on peut également trouver les intervalles de confiance unilatères suivant :

$$IC_{1-\alpha} = \left[\bar{X}_n - t_{n-1, 1-\alpha} \frac{S_n}{\sqrt{n}}; +\infty \right] \quad \text{et} \quad \left[-\infty, \bar{X}_n + t_{n-1, 1-\alpha} \frac{S_n}{\sqrt{n}} \right].$$

3.3.3 Estimation de la variance

La proposition suivante nous donne les intervalles de confiance pour la variance.

Proposition 3.3.3. Soit $\alpha \in (0, 1)$. Alors

$$\mathbb{P} \left[\frac{(n-1)S_n^2}{k_{1-\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)S_n^2}{k_{\alpha/2}} \right] = 1 - \alpha,$$

où $k_{\alpha/2}$ et $k_{1-\alpha/2}$ sont les quantiles d'ordres $\alpha/2$ et $1 - \alpha/2$ d'une loi du chi deux à $n - 1$ degrés de libertés, i.e si $Z \sim \chi_{n-1}^2$,

$$\mathbb{P} [Z \leq k_{\alpha/2}] = \alpha/2$$

$$\mathbb{P} [Z \leq k_{1-\alpha/2}] = 1 - \alpha/2.$$

On obtient donc l'intervalle de confiance

$$IC_{1-\alpha} (\sigma^2) = \left[\frac{(n-1)S_n^2}{k_{1-\alpha/2}}, \frac{(n-1)S_n^2}{k_{\alpha/2}} \right].$$

Démonstration. On a vu que

$$\frac{n-1}{\sigma^2} S_n^2 \sim \chi_{n-1}^2.$$

On a donc

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left[k_{\alpha/2} \leq \frac{n-1}{\sigma^2} S_n^2 \leq k_{1-\alpha/2} \right] \\ &= \mathbb{P} \left[\frac{k_{\alpha/2}}{(n-1)S_n^2} \leq \frac{1}{\sigma^2} \leq \frac{k_{1-\alpha/2}}{(n-1)S_n^2} \right] \\ &= \mathbb{P} \left[\frac{(n-1)S_n^2}{k_{\alpha/2}} \geq \sigma^2 \geq \frac{(n-1)S_n^2}{k_{1-\alpha/2}} \right] \end{aligned}$$

□

3.4 Intervalles de confiance asymptotiques

Dans cette section, on s'intéresse à l'estimation du paramètre θ d'une variable aléatoire X . On considère un estimateur $\hat{\theta}_n$ asymptotiquement normal, i.e il existe $\sigma^2 > 0$ tel que

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2).$$

On supposera également que l'on a un estimateur $\hat{\sigma}_n^2$ consistant de σ^2 .

Proposition 3.4.1. Soit $\alpha \in (0, 1)$,

$$\mathbb{P} \left[\hat{\theta}_n - q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}} \right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha,$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite.

On obtient donc l'intervalle de confiance

$$IC_{1-\alpha}(\theta) = \left[\hat{\theta}_n - q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}}; \hat{\theta}_n + q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}} \right]$$

Démonstration. On a

$$Z_n := \sqrt{n} \frac{\hat{\theta}_n - \theta}{\hat{\sigma}_n} = \frac{\sigma}{\hat{\sigma}_n} \sqrt{n} \frac{\hat{\theta}_n - \theta}{\sigma},$$

et par le théorème de continuité, $\frac{\sigma}{\hat{\sigma}_n}$ converge en probabilité vers 1. On obtient, à l'aide du Théorème de Slutsky,

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

On note F_{Z_n} la fonction de répartition de Z_n et F_Z celle de Z où $Z \sim \mathcal{N}(0, 1)$. Donc, pour tout x , $F_{Z_n}(x)$ converge vers $F(x)$ et en particulier

$$\mathbb{P} \left[\frac{\sqrt{n}}{\hat{\sigma}_n} |\hat{\theta}_n - \theta| \leq q_{1-\alpha/2} \right] \xrightarrow[n \rightarrow +\infty]{} \mathbb{P} [|Z| \leq q_{1-\alpha/2}] = 1 - \alpha.$$

□

Remarque : Attention, ceci ne reste qu'un résultat asymptotique, et la question est de savoir à partir de quelle taille d'échantillon l'approximation par une loi normale est de bonne qualité. Si on a la chance de pouvoir avoir des intervalles non asymptotiques, il faut toujours privilégier cette option !

Exemple 1 : le lancer de pièce. On considère X_1, \dots, X_n des variables aléatoires i.i.d suivant une loi de Bernoulli de paramètre $\theta \in (0, 1)$. On a vu

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

On a donc

$$\mathbb{P} \left[-q_{1-\alpha/2} \leq \sqrt{n} \frac{\bar{X}_n - \theta}{\sqrt{\bar{X}_n (1 - \bar{X}_n)}} \leq q_{1-\alpha/2} \right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite. On obtient alors

$$\begin{aligned} & \mathbb{P} \left[-q_{1-\alpha/2} \leq \sqrt{n} \frac{\bar{X}_n - \theta}{\sqrt{\bar{X}_n (1 - \bar{X}_n)}} \leq q_{1-\alpha/2} \right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha \\ & \mathbb{P} \left[-q_{1-\alpha/2} \frac{\sqrt{\bar{X}_n (1 - \bar{X}_n)}}{\sqrt{n}} \leq \bar{X}_n - \theta \leq q_{1-\alpha/2} \frac{\sqrt{\bar{X}_n (1 - \bar{X}_n)}}{\sqrt{n}} \right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha \\ & \mathbb{P} \left[\bar{X}_n - q_{1-\alpha/2} \frac{\sqrt{\bar{X}_n (1 - \bar{X}_n)}}{\sqrt{n}} \leq \theta \leq \bar{X}_n + q_{1-\alpha/2} \frac{\sqrt{\bar{X}_n (1 - \bar{X}_n)}}{\sqrt{n}} \right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha \end{aligned}$$

Exemple 2 : la loi exponentielle. On considère une variable aléatoire X suivant une loi exponentielle de paramètre θ . Rappelons que l'on a alors $\mathbb{E}[X] = \theta^{-1}$ (et $\mathbb{V}[X] = \theta^{-2}$) et qu'un estimateur naturel de θ est alors $\hat{\theta}_n = \frac{n}{\sum_{i=1}^n X_i} = \frac{1}{\bar{X}_n}$. On a alors, par la delta-méthode (voir TD)

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \theta^2),$$

ce que l'on peut réécrire comme

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

A l'aide du théorème de Slutsky, on obtient

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

et on a donc, pour tout $\alpha \in (0, 1)$

$$\mathbb{P} \left[-q_{1-\alpha/2} \leq \sqrt{n} \frac{\hat{\theta}_n - \theta}{\hat{\theta}_n} \leq q_{1-\alpha/2} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha,$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite. On obtient alors

$$\begin{aligned} & \mathbb{P} \left[-q_{1-\alpha/2} \leq \sqrt{n} \frac{\hat{\theta}_n - \theta}{\hat{\theta}_n} \leq q_{1-\alpha/2} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha \\ \Leftrightarrow & \mathbb{P} \left[-q_{1-\alpha/2} \frac{\hat{\theta}_n}{\sqrt{n}} \leq \hat{\theta}_n - \theta \leq q_{1-\alpha/2} \frac{\hat{\theta}_n}{\sqrt{n}} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha, \\ \Leftrightarrow & \mathbb{P} \left[\hat{\theta}_n - q_{1-\alpha/2} \frac{\hat{\theta}_n}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + q_{1-\alpha/2} \frac{\hat{\theta}_n}{\sqrt{n}} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha. \end{aligned}$$

Exemple 2bis : la loi exponentielle. En réalité, on aurait pu être plus malin pour l'intervalle de confiance. En effet, le TLC peut se réécrire comme

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

En notant $q_{1-\alpha/2}$ le quantile de la loi normale centrée réduite, on obtient donc

$$\begin{aligned} & \mathbb{P} \left[-q_{1-\alpha/2} \leq \sqrt{n} \left(\frac{\hat{\theta}_n}{\theta} - 1 \right) \leq q_{1-\alpha/2} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha \\ & \mathbb{P} \left[1 - \frac{q_{1-\alpha/2}}{\sqrt{n}} \leq \frac{\hat{\theta}_n}{\theta} \leq 1 + \frac{q_{1-\alpha/2}}{\sqrt{n}} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha \\ & \mathbb{P} \left[\frac{\hat{\theta}_n}{1 + \frac{q_{1-\alpha/2}}{\sqrt{n}}} \geq \theta \geq \frac{\hat{\theta}_n}{1 - \frac{q_{1-\alpha/2}}{\sqrt{n}}} \right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha. \end{aligned}$$

On obtient donc l'intervalle de confiance asymptotique

$$IC_{1-\alpha}(\theta) = \left[\frac{\hat{\theta}_n}{1 + \frac{q_{1-\alpha/2}}{\sqrt{n}}}; \frac{\hat{\theta}_n}{1 - \frac{q_{1-\alpha/2}}{\sqrt{n}}} \right].$$

Chapitre 4

Vecteurs Gaussiens et théorème de Cochran

Une généralisation du Théorème Central Limite lorsque l'on doit traiter des données multivariées, i.e à valeurs dans \mathbb{R}^d est la suivante :

Théorème 4.0.1 (Théorème Central Limite multivarié). Soit X_1, \dots, X_n des vecteurs aléatoires à valeurs dans \mathbb{R}^d , indépendants et identiquement distribués. On suppose $\mathbb{E} [\|X_1\|^2] < +\infty$ et on note $m = \mathbb{E} [X_1]$ son vecteur moyen (que l'on définira par la suite) et $\Gamma = \text{Var} [X_1]$ sa matrice de variance-covariance (que l'on définira par la suite). On note $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ et on a alors la convergence en loi

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

Ainsi, connaître et maîtriser les propriétés des vecteurs gaussiens peut être crucial pour l'étude d'estimateurs lorsque le paramètre à estimer appartient à \mathbb{R}^d , mais pas seulement ! Sans le savoir, on a déjà vu une application du théorème de Cochran (issu de l'étude des vecteurs gaussiens que l'on verra par la suite) pour obtenir la loi de l'estimateur de la variance dans le cadre gaussien. Plus précisément, dans ce cours (et en TD), on utilisera les vecteurs gaussiens pour les problèmes suivants :

- Obtention d'intervalles de confiance non asymptotique de la moyenne et de la variance d'une variable aléatoire suivant une loi normale (cf chapitre précédent).
- Estimation des paramètres dans le cadre du modèle linéaire gaussien, i.e on considère :
 - des variables aléatoires Y_1, \dots, Y_n que lon cherche à expliquer
 - pour tout i , des variables explicatives $x_{i,1}, \dots, x_{i,p}$
 - des variables aléatoires indépendantes $\epsilon_i \sim \mathcal{N} (0, \sigma^2)$
 - le modèle de régression

$$Y_i = a_1 x_{i,1} + \dots + a_p x_{i,p} + \epsilon_i$$

et l'objectif est donc d'estimer a_1, \dots, a_p , construire des intervalles de confiance... (voir TD).

- Construction de tests dans le cadre gaussien. On verra par la suite comment tester la va-

leur d'une moyenne, d'une variance, comparer les moyennes de différentes variables aléatoires....

4.1 Vecteurs aléatoires

Soit $X = \begin{pmatrix} X_1 \\ \dots \\ X_d \end{pmatrix} \in \mathbb{R}^d$ un vecteur aléatoire tel que $\mathbb{E} [X_j^2] < +\infty$. Pour tout i, j , la covariance entre X_i et X_j est définie par

$$\text{Cov} (X_i, X_j) = \mathbb{E} [(X_i - \mathbb{E} [X_i]) (X_j - \mathbb{E} [X_j])] = \mathbb{E} [X_i X_j] - \mathbb{E} [X_i] \mathbb{E} [X_j].$$

Remarque : Si X_i et X_j sont indépendantes, $\text{Cov} (X_i, X_j) = 0$ mais la réciproque est généralement fautive (sauf pour les vecteurs gaussiens, ce que l'on verra par la suite).

Definition 4.1.1. Soit X un vecteur aléatoire de \mathbb{R}^d tel que pour tout $i = 1, \dots, d$, $\mathbb{E} [X_i^2] < +\infty$. Le vecteur moyen de X est défini par

$$\mathbb{E} [X] = \begin{pmatrix} \mathbb{E} [X_1] \\ \vdots \\ \mathbb{E} [X_d] \end{pmatrix}$$

et sa matrice de covariance est définie par

$$\begin{aligned} \text{Var} [X] &= \mathbb{E} [(X - \mathbb{E} [X]) (X - \mathbb{E} [X])^T] \\ &= \begin{pmatrix} \text{Var} [X_1] & \text{Cov} (X_1, X_2) & \dots & \text{Cov} (X_1, X_d) \\ \text{Cov} (X_2, X_1) & \text{Var} [X_2] & \dots & \text{Cov} (X_2, X_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov} (X_d, X_1) & \text{Cov} (X_d, X_2) & \dots & \text{Var} [X_d] \end{pmatrix} \end{aligned}$$

A noter que pour que le vecteur moyen soit défini, il suffit de supposer que pour tout i , la variable aléatoire X_i est intégrable.

Théorème 4.1.1. La matrice de covariance est une matrice symétrique semi-définie positive.

La symétrie vient de la symétrie de la covariance, i.e pour tout i, j ,

$$\text{Cov} (X_i, X_j) = \text{Cov} (X_j, X_i).$$

De plus, on a pour tout $v = (v_1, \dots, v_d)^T$

$$v^T \text{Var} (X) = \left(\sum_{i=1}^d v_i \text{Cov} (X_i, X_1), \sum_{i=1}^d v_i \text{Cov} (X_i, X_2), \dots, \sum_{i=1}^d v_i \text{Cov} (X_i, X_d) \right)$$

et on obtient donc

$$\begin{aligned}
 v^T \text{Var}(X)v &= \sum_{j=1}^d \sum_{i=1}^d v_i \text{Cov}(X_i, X_j) v_j \\
 &= \sum_{j=1}^d \sum_{i=1}^d \text{Cov}(v_i X_i, v_j X_j) \\
 &= \sum_{i=1}^d \mathbb{V}[v_i X_i] + \sum_{i \neq j} \text{Cov}(v_i X_i, v_j X_j) \\
 &= \text{Var}(v_1 X_1 + \dots + v_d X_d) \geq 0.
 \end{aligned}$$

Théorème 4.1.2. Soit X un vecteur aléatoire de \mathbb{R}^d de moyenne m et de matrice de covariance C . Alors, si A est une matrice réelle $p \times d$, le vecteur aléatoire $AX \in \mathbb{R}^p$ a pour vecteur moyen Am et pour matrice de covariance ACA^T .

Démonstration. Par linéarité de l'espérance, on a

$$\mathbb{E}[AX] = A\mathbb{E}[X] = Am$$

Si vraiment on veut s'en convaincre, on a

$$AX = \left(\sum_{i=1}^d a_{1,i} X_i, \dots, \sum_{i=1}^d a_{p,i} X_i \right) \quad \text{et} \quad Am = \left(\sum_{i=1}^d a_{1,i} \mathbb{E}[X_i], \dots, \sum_{i=1}^d a_{p,i} \mathbb{E}[X_i] \right)$$

et par définition de l'espérance multidimensionnelle et par linéarité, on obtient le résultat. \square

Definition 4.1.2. Soit X un vecteur aléatoire de \mathbb{R}^d . Sa fonction caractéristique Φ_X est définie pour tout $u \in \mathbb{R}^d$ par

$$\Phi_X(u) = \mathbb{E} \left[e^{iu^T X} \right] = \mathbb{E} \left[e^{i(u_1 X_1 + \dots + u_d X_d)} \right].$$

4.2 Vecteurs aléatoires gaussiens

Definition 4.2.1. Soit X un vecteur aléatoire de \mathbb{R}^d . X est un vecteur gaussien de \mathbb{R}^d si et seulement si toute combinaison linéaire de ses composantes est une variable aléatoire réelle gaussienne, i.e

$$\forall u \in \mathbb{R}^d, \quad u^T X \sim \mathcal{N}(m_u, \sigma_u^2).$$

Exemple : Un vecteur X dont les composantes X_i sont indépendantes et suivent une loi normale est un vecteur gaussien.

Théorème 4.2.1. Soit X un vecteur aléatoire de \mathbb{R}^d . Les assertions suivantes sont équivalentes :

1. Le vecteur X est un vecteur gaussien de moyenne m et de matrice de covariance Γ .
2. La fonction caractéristique de X est donnée pour tout $u \in \mathbb{R}^d$ par

$$\Phi_X(u) = \mathbb{E} \left[e^{iu^T X} \right] = \exp \left(iu^T m - \frac{1}{2} u^T \Gamma u \right).$$

Démonstration. Soit $X \sim \mathcal{N}(m, \Gamma)$. A noter que pour tout $u \in \mathbb{R}^d$, on a par définition $u^T X$ qui suit une loi normale et plus précisément

$$u^T X \sim \mathcal{N} \left(u^T m, u^T \Gamma u \right).$$

On a donc

$$\Phi_{u^T X}(t) = \mathbb{E} \left[e^{itu^T X} \right] = \exp \left(itu^T m - \frac{1}{2} u^T \Gamma u t^2 \right).$$

On obtient donc pour tout u ,

$$\Phi_X(u) = \mathbb{E} \left[e^{iu^T X} \right] = \Phi_{u^T X}(1) = \exp \left(iu^T m - \frac{1}{2} u^T \Gamma u \right).$$

□

Corollaire 4.2.1. Soit $X \sim \mathcal{N}(m, \Gamma)$ à valeurs dans \mathbb{R}^d . Soit A une matrice $p \times d$ et $b \in \mathbb{R}^p$, alors

$$AX + b \sim \mathcal{N} \left(Am + b, A\Gamma A^T \right).$$

Démonstration. On a pour tout $u \in \mathbb{R}^p$,

$$\begin{aligned} \Phi_{AX+b}(u) &= \mathbb{E} \left[\exp \left(iu^T AX + u^T b \right) \right] \\ &= \exp \left(iu^T b \right) \mathbb{E} \left[\exp \left(iu^T AX \right) \right] \\ &= \exp \left(iu^T b \right) \Phi_X \left(A^T u \right) \\ &= \exp \left(iu^T b \right) \exp \left(iu^T Am - \frac{1}{2} u^T A\Gamma A^T u \right) \\ &= \exp \left(iu^T (Am + b) - \frac{1}{2} u^T (A\Gamma A^T) u \right). \end{aligned}$$

□

Corollaire 4.2.2. Soit $X \sim \mathcal{N}(m, \Gamma)$, si Γ est inversible, alors

$$\Gamma^{-1/2} (X - m) \sim \mathcal{N} (0, I_d).$$

Démonstration. Soit $Y = \Gamma^{-1/2} (X - m)$, d'après le corollaire précédent, Y est un vecteur gaussien

et par linéarité

$$\mathbb{E}[Y] = \Gamma^{-1/2} \mathbb{E}[X - m] = 0.$$

De plus, comme $Y = \Gamma^{-1/2}X - \Gamma^{-1/2}m$,

$$\text{Var}[Y] = \Gamma^{-1/2} \Gamma \left(\Gamma^{-1/2} \right)^T = I_d.$$

□

Théorème 4.2.2. Soit $X \in \mathbb{R}^d$ avec $X \sim \mathcal{N}(m, \Gamma)$. X admet une densité par rapport à la mesure de Lebesgue sur \mathbb{R}^d si et seulement si Γ est inversible. Dans ce cas, la densité de f est donnée pour tout $x \in \mathbb{R}^d$ par

$$f(x) = \frac{1}{\sqrt{(2\pi)^d |\det(\Gamma)|}} \exp\left(-\frac{1}{2}(x - m)^T \Gamma^{-1}(x - m)\right)$$

Théorème 4.2.3. Pour tout vecteur gaussien X de \mathbb{R}^d , les propriétés suivantes sont équivalentes :

- Les composantes X_1, \dots, X_d de X sont indépendantes.
- La matrice de covariance Γ de X est diagonale.

Démonstration. L'implication (1) \Rightarrow (2) est évidente. Dans l'autre sens, on rappelle que des variables aléatoires Y_1, \dots, Y_d sont indépendantes si et seulement si pour tout $u \in \mathbb{R}^d$,

$$\Phi_{(Y_1, \dots, Y_d)}(u) = \prod_{i=1}^d \Phi_{Y_i}(u_i)$$

On note

$$\Gamma = \text{Diag}(\sigma_1^2, \dots, \sigma_d^2) = \text{Diag}(\text{Var}(X_1), \dots, \text{Var}(X_d))$$

la matrice de covariance de X . On a alors pour tout $u \in \mathbb{R}^d$,

$$\begin{aligned} \Phi_X(u) &= \exp\left(iu^T m - \frac{1}{2}u^T \Gamma u\right) \\ &= \exp\left(i \sum_{j=1}^d u_j m_j - \frac{1}{2} \sum_{j=1}^d \sigma_j^2 u_j^2\right) \\ &= \prod_{j=1}^d \exp\left(iu_j m_j - \frac{1}{2} \sigma_j^2 u_j^2\right) \\ &= \prod_{j=1}^d \Phi_{X_j}(u_j). \end{aligned}$$

□

Exemple : Soit $X = (X_1, X_2)^T \sim \mathcal{N}(0, V)$. Alors la variable aléatoire $Z = X_1 + X_2 \sim \mathcal{N}(0, u^T V u)$. En effet il suffit juste de remarquer que $Z = u^T X$.

Corollaire 4.2.3. Soit X un vecteur gaussien, alors pour tout $i \neq j$, X_i et X_j sont indépendantes si et seulement si

$$\text{Cov}(X_i, X_j) = 0.$$

Démonstration. La première implication est claire. Pour tout i, j tel que $i \neq j$, on note $\tilde{X} = (X_i, X_j)^T$ qui est un vecteur gaussien de matrice de covariance

$$\text{Var}(\tilde{X}) = \begin{pmatrix} \text{Var}(X_i) & \text{Cov}(X_i, X_j) \\ \text{Cov}(X_i, X_j) & \text{Var}(X_j) \end{pmatrix}$$

Si $\text{Cov}(X_i, X_j) = 0$ la matrice est alors diagonale et X_i, X_j sont donc indépendants. \square

Attention! Encore une fois, ceci est vrai car X est un vecteur gaussien! Il n'y a pas d'implication générale!

Contre-exemple : Soit $X \sim \mathcal{N}(0, 1)$ et Y une variable aléatoire suivant une loi de Rademacher de paramètre $p \in (0, 1)$, i.e

$$\mathbb{P}[Y = 1] = p \quad \text{et} \quad \mathbb{P}[Y = -1] = 1 - p.$$

Dans ce qui suit, on prend $p = 1/2$, on considère que X et Y sont indépendants, et on considère les variables aléatoires X et $Z = XY$. A noter que $\mathbb{E}[Y] = 0$, et donc $\mathbb{E}[Z] = 0$. De plus

$$\text{Cov}(X, Z) = \mathbb{E}[(X - \mathbb{E}[X])(Z - \mathbb{E}[Z])] = \mathbb{E}[XZ] = \mathbb{E}[X^2Y] = \mathbb{E}[X^2] \mathbb{E}[Y] = 0.$$

Cependant, on a

$$\begin{aligned} \mathbb{P}[X \leq 1, Z \leq -2] &= \mathbb{P}[X \leq 1, XY \leq -2] = \mathbb{P}[X \leq 1, X \leq -2, Y = 1] + \underbrace{\mathbb{P}[X \leq 1, -X \leq -2, Y = -1]}_{=0} \\ &= \frac{1}{2} \mathbb{P}[X \leq -2] \end{aligned}$$

et par symétrie de la loi normale

$$\begin{aligned} \mathbb{P}[X \leq 1] \mathbb{P}[Z \leq -2] &= \mathbb{P}[X \leq 1] (\mathbb{P}[X \leq -2, Y = 1] + \mathbb{P}[-X \leq -2, Y = -1]) \\ &= \mathbb{P}[X \leq 1] \left(\frac{1}{2} \mathbb{P}[X \leq -2] + \frac{1}{2} \mathbb{P}[-X \leq -2] \right) \\ &= \underbrace{\mathbb{P}[X \leq 1]}_{\neq 1/2} \mathbb{P}[X \leq -2] \\ &\neq \mathbb{P}[X \leq 1, Z \leq -2]. \end{aligned}$$

Les variables X et Z ne sont donc pas indépendantes. On obtient donc un exemple de variables aléatoires avec une covariance nulle mais pas indépendante. On peut aussi en conclure que le

vecteur $(X, Z)^T$ n'est pas gaussien. Or pour tout $x \in \mathbb{R}$, par symétrie de la loi normale,

$$\mathbb{P}[Z \leq x] = \mathbb{P}[X \leq x, Y = 1] + \mathbb{P}[-X \leq x, Y = -1] = \frac{1}{2} (\mathbb{P}[X \leq x] + \mathbb{P}[X \geq -x]) = \mathbb{P}[X \leq x].$$

Z suit donc la même loi que X , i.e Z suit une loi normale. On a donc un exemple de vecteur $((X, Z)^T)$ dont les coordonnées suivent des loi normales mais qui n'est pas gaussien.

4.3 Théorème de Cochran et applications

Théorème 4.3.1 (Théorème de Cochran). Soit X un vecteur aléatoire de \mathbb{R}^n suivant une loi $\mathcal{N}(m, \sigma^2 I_n)$ avec $\sigma^2 > 0$. Soit $\mathbb{R}^n = E_1 \oplus E_2 \oplus \dots \oplus E_p$ une décomposition de \mathbb{R}^n en somme directe de p sous-espaces vectoriels orthogonaux de dimensions respectives d_1, \dots, d_p . Soit P_k le projecteur orthogonal sur E_k et $Y_k = P_k X$ la projection orthogonale de X sur E_k .

1. Les projections Y_1, \dots, Y_p sont des vecteurs gaussiens indépendants et $Y_k \sim \mathcal{N}(P_k m, \sigma^2 P_k)$.
2. Les variables aléatoires $\|Y_1 - P_1 m\|^2, \dots, \|Y_p - P_p m\|^2$ sont indépendantes et

$$\frac{\|Y_k - P_k m\|^2}{\sigma^2} \sim \chi_{d_k}^2$$

Remarque : A noter que pour tout $k \neq k'$, on a également indépendance entre Y_k et $\|Y_{k'} - P_{k'} m\|^2$.

Démonstration. On note

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_p \end{pmatrix} = \begin{pmatrix} P_1 \\ \vdots \\ P_p \end{pmatrix} X = AX$$

On a donc $Y \sim \mathcal{N}(Am, \sigma^2 AA^T)$. De plus comme les P_k sont des projecteurs orthogonaux, on a $P_k^2 = P_k$ et $P_k P_{k'} = 0$ si $k \neq k'$. On a donc $AA^T = \text{Diag}(P_1, \dots, P_k)$ qui est diagonale par block et on a donc que les Y_k sont indépendants, et en particulier les variables aléatoires $\|Y_k - P_k m\|^2$ sont indépendantes. De plus

$$Y_k \sim \mathcal{N}(P_k m, \sigma^2 P_k).$$

Soit $B_k = \{e_{k,1}, \dots, e_{k,d_k}\}$ une base orthonormée de E_k , on a alors

$$Y_k - P_k m = \sum_{j=1}^{d_k} \langle X - m, e_{k,j} \rangle e_{k,j},$$

et en complétant la base B_k en une base orthonormée de \mathbb{R}^n , on a $Y_k - P_k m$ et sa matrice de cova-

riance qui s'écrivent

$$(Y_k - P_k m)_{B_k} = \begin{pmatrix} \langle X - m, e_{k,1} \rangle \\ \vdots \\ \langle X - m, e_{k,d_k} \rangle \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \text{et} \quad \text{Var} (Y_k - P_k m)_{B_k} = \begin{pmatrix} \sigma^2 I_{d_k} & 0 \\ 0 & 0 \end{pmatrix}$$

et on a donc que les $\langle X - m, e_{k,j} \rangle$ sont indépendants et pour $j = 1, \dots, d_k$,

$$\langle X - m, e_{k,j} \rangle \sim \mathcal{N}(0, \sigma^2)$$

et donc

$$\frac{1}{\sigma^2} \|Y_k - P_k m\|^2 = \sum_{j=1}^{d_k} \frac{\langle X - m, e_{k,j} \rangle^2}{\sigma^2} \sim \chi_{d_k}^2.$$

□

Corollaire 4.3.1. Soient X_1, \dots, X_n des variables aléatoires indépendantes et identiquement distribuées avec $X_1 \sim \mathcal{N}(m, \sigma^2)$. On a alors

1. $\frac{n-1}{\sigma^2} S_n^2 \sim \chi_{n-1}^2$.
2. \bar{X}_n et S_n^2 sont indépendantes.

Démonstration. On note $X = (X_1, \dots, X_n)^T \sim \mathcal{N}(me, \sigma^2 I_n)$ avec $e = (1, \dots, 1)^T \in \mathbb{R}^n$ et $u = \frac{e}{\|e\|} = n^{-1/2} (1, \dots, 1)$. On note $E = \text{Vect}\{u\}$ et

$$\mathbb{R}^n = E \oplus E^\perp.$$

On note P_E la projection orthogonale sur E , et on a

$$P_E X = \langle X, u \rangle u = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i u = \bar{X}_n e.$$

De plus, $P_E(me) = me$ et $P_{E^\perp}(me) = 0$. De plus,

$$P_{E^\perp}(X) = (I_n - P_E)(X) = X - \bar{X}_n e = \begin{pmatrix} X_1 - \bar{X}_n \\ \vdots \\ X_n - \bar{X}_n \end{pmatrix}.$$

Enfin,

$$\|P_E(X) - P_E(me)\|^2 = n(\bar{X}_n - m)^2$$

et

$$\|P_{E^\perp}(X) - P_{E^\perp}(me)\|^2 = \sum_{i=1}^n (X_i - m)^2 = (n-1)S_n^2.$$

On a donc, comme $\dim(E) = 1$, en appliquant le théorème de Cochran,

1. $n(\bar{X}_n - m)^2$ et $(n-1)S_n^2$ sont indépendantes, i.e \bar{X}_n et S_n^2 sont indépendantes.
2. On a

$$\frac{n(\bar{X}_n - m)^2}{\sigma^2} \sim \chi_1^2 \quad \text{et} \quad \frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2.$$

□

Proposition 4.3.1. Soient Z, Z_p, Z_q trois variables aléatoires telles que

1. $Z = Z_p + Z_q$
2. $Z_p \sim \chi_p^2$ et $Z_q \sim \chi_q^2$
3. Z_p et Z_q sont indépendantes

alors $Z \sim \chi_{p+q}^2$.

Démonstration. Soit $Z'_p = \sum_{i=1}^p X_i^2$ et $Z'_q = \sum_{i=p+1}^{p+q} X_i^2$ où les X_i sont indépendants et $X_i \sim \mathcal{N}(0, 1)$.

On a Z'_p et Z'_q qui sont indépendants et pour tout t ,

$$\Phi_{Z'_p + Z'_q}(t) = \Phi_{Z'_p}(t)\Phi_{Z'_q}(t) = \Phi_{Z_p}(t)\Phi_{Z_q}(t) = \Phi_Z(t).$$

et donc Z et $Z'_p + Z'_q$ ont la même loi. Or

$$Z'_p + Z'_q = \sum_{i=1}^{p+q} X_i^2 \sim \chi_{p+q}^2.$$

□

Chapitre 5

Tests

Un test statistique est une procédure mathématique qui permet, à l'aide d'un modèle probabiliste, de rejeter ou non une hypothèse avec un certain risque connu. En pratique, on cherche à expliquer les différences observées entre une hypothèse posée a priori et le résultat obtenu à partir des données. Le test permet alors de décider si ces différences sont seulement dues à la fluctuation d'échantillonnage ou bien si elles sont significatives. Cette dernière possibilité impliquerait alors que la ou les hypothèses posées a priori ne sont pas fondées.

5.1 Généralités

5.1.1 Généralités

Le principe d'un test consiste à

1. Poser une hypothèse H_0 et une hypothèse alternative H_1
2. A l'aide du modèle probabiliste posé, définir une zone de rejet de l'hypothèse nulle H_0 , i.e une zone où les résultats expérimentaux seraient trop éloignés du résultat théorique pour être crédibles.
3. Regarder si le résultat expérimental est dans cette zone de rejet.
4. Si le résultat est dans la zone de rejet, on rejette l'hypothèse nulle H_0 .
5. Si le résultat n'est pas dans la zone de rejet, on ne rejette pas H_0 .

Attention! La sémantique est importante, on dit qu'on ne rejette pas H_0 et non pas que l'on accepte H_0 . Même si la différence semble faible entre ces deux expressions, on verra par la suite pourquoi c'est important.

Exemple : On considère l'exemple du lancer de pièce. On jette 10 fois une pièce et on obtient 9 "Face". Un test statistique nous permettra de dire avec un certain risque si cette proportion de "Face" est seulement due à la fluctuation d'échantillonnage ou bien si la pièce est truquée.

5.1.2 Etapes d'un test statistique

On dispose de n données x_1, \dots, x_n qui sont des mesures de variables qualitatives ou quantitatives.

Introduire un modèle probabiliste : On peut par exemple faire l'hypothèse que les x_i sont des réalisations de n variables aléatoires indépendantes et identiquement distribuées X_1, \dots, X_n .

Dans le cas du lancer de pièce, on notera $x_i = 1$ si le i -ème lancer donne "Pile" et 0 sinon. On considère alors les x_i comme étant les réalisations de v.a i.i.d X_i suivant une loi de Bernoulli de paramètre θ .

Les hypothèses nulle et alternative : Disposant d'un modèle probabiliste, on souhaite vérifier une hypothèse sur ce modèle. On appelle cette hypothèse l'hypothèse nulle H_0 , et on introduit sa négation (hypothèse alternative) H_1 .

Dans le cas du lancer de pièce on veut savoir si la pièce est truquée ou non, i.e si $\theta = 0.5$. On pose alors les hypothèses :

$$H_0 : "\theta = 0.5" \quad \text{contre} \quad H_1 : "\theta \neq 0.5"$$

On peut également se demander si la probabilité d'obtenir "Pile" n'est pas en dessous de 0.5 et on pose alors

$$H_0 : "\theta < 0.5" \quad \text{contre} \quad H_1 : "\theta \geq 0.5"$$

La statistique de test et sa loi sous H_0 : Une fois que l'on a introduit le modèle probabiliste et les hypothèses, il faut pouvoir disposer d'une variable aléatoire appelée statistique de test, que l'on notera Z et :

- On connaît sa loi sous H_0 .
- Sa loi n'est pas la même sous H_1 .

Bien évidemment, Z est définie en fonction des variables aléatoires X_1, \dots, X_n , i.e $Z = T(X_1, \dots, X_n)$ et si on note $z_{obs} = T(x_1, \dots, x_n)$, z_{obs} doit être calculable.

Zone de rejet : Connaissant la loi de Z sous H_0 , on peut déterminer ses valeurs les plus extrêmes, i.e déterminer une zone ZR où $\mathbb{P}[Z \in ZR] \leq \alpha$ si H_0 est vraie. Le seuil α est choisi librement. Dans le cas où on teste :

1. $H_0 : "\theta = \theta_0"$, les valeurs "extrêmes" sont les valeurs les plus éloignées de la valeur attendue.
2. $H_0 : "\theta < \theta_0"$, les valeurs "extrêmes" sont les grandes valeurs.
3. $H_0 : "\theta > \theta_0"$, les valeurs "extrêmes" sont les petites valeurs.

Conclusion : On conclut le test à partir de la valeur observée z_{obs} :

- Si z_{obs} appartient à la zone de rejet, on rejette H_0 au risque α .
- Si z_{obs} n'appartient pas à la zone de rejet, on ne rejette pas H_0 au risque α .

Si on reprend l'exemple du lancer de pièce, on souhaite tester si la pièce est truquée ou non au risque de 5%, i.e on teste

$$H_0 : "\theta = 0.5" \quad \text{contre} \quad H_1 : "\theta \neq 0.5"$$

Sous H_0 , X_i suit une loi de Bernoulli de paramètre 0.5 et $Z = \sum_{i=1}^{10} X_i$ suit une loi binomiale de paramètre $n = 10$, $\theta = 0.5$ sous H_0 , ce qui est faux sous H_1 . On a donc notre statistique de test. Comme les quantiles d'ordres $\alpha/2$ et $1 - \alpha/2$ de la loi binomiale de paramètres 10, 0.5 sont 2 et 8, on a la zone de rejet

$$ZR_{H_0} = \{0, 1\} \cup \{9, 10\}$$

Nous avons observé 9 "Face", donc $z_{obs} \in ZR_{H_0}$ et on rejette donc H_0 au risque de 5%, i.e au risque de 5%, on peut dire que la pièce est truquée.

p-value : Plus le seuil α est petit, plus le test est fiable. Comme il est inconcevable de faire les tests pour toutes les valeurs α possibles, on va s'intéresser à la p -value qui est définie comme

$$p - \text{value} = \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \}.$$

On verra par la suite comment calculer la p -value pour les différents tests. La p -value peut donc être vue comme étant, sous H_0 , la probabilité d'avoir eu un "aussi mauvais résultat" que z_{obs} .

- Si $p - \text{value} \geq \alpha$, alors on ne rejette pas H_0
- Si $p - \text{value} < \alpha$, alors on rejette H_0

Le risque α peut donc être vu comme étant un seuil subjectif où on estime qu'un évènement ayant une probabilité inférieure à ce seuil de se réaliser ne devrait pas se réaliser.

5.2 Tests sur la moyenne et la variance

On considère maintenant que les données réelles x_1, \dots, x_n sont des réalisations de variables aléatoires X_1, \dots, X_n avec $X_1 \sim \mathcal{N}(\mu, \sigma^2)$ et μ, σ^2 sont inconnus.

5.2.1 Test de conformité d'une moyenne

On veut savoir avec un risque α si la moyenne théorique μ est différente ou non d'une valeur μ_0 donnée.

Construction de la statistique de test : On dispose de l'estimateur \bar{X}_n de la moyenne μ qui vérifie

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \sim \mathcal{N}(0, 1).$$

Comme la variance σ^2 est inconnue, on utilise alors

$$\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim T_{n-1},$$

où T_{n-1} suit une loi de Student à $n - 1$ degrés de liberté. Enfin, sous H_0 , comme $\mu = \mu_0$, on a alors

$$\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \sim T_{n-1}.$$

ce qui est bien évidemment faux sous H_1 . En effet, sous H_1 ,

$$\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} = \sqrt{n} \frac{\bar{X}_n - \mu}{S_n} + \sqrt{n} \frac{\mu - \mu_0}{S_n}$$

et ne suit donc pas une loi de Student à $n - 1$ degrés de liberté.

Le test : On teste au risque α l'hypothèse nulle $H_0 : \mu = \mu_0$ contre l'hypothèse alternative $H_1 : \mu \neq \mu_0$, on utilise la statistique de test

$$Z = \sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \sim T_{n-1}, \quad \text{sous } H_0.$$

On définit la zone de rejet par

$$ZR = \{|Z| > t_{n-1, 1-\alpha/2}\}$$

où $t_{n-1, 1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi de Student à $n - 1$ degrés de liberté. On calcule alors $z_{obs} = \sqrt{n} \frac{\bar{x}_n - \mu_0}{s_n}$.

Si $|z_{obs}| > t_{n-1, 1-\alpha/2}$, on rejette H_0 , sinon, on ne rejette pas H_0 .

Remarque 1 : Notons que faire ce test revient donc à vérifier que μ_0 appartient à la réalisation de l'intervalle de confiance de niveau $1 - \alpha$ obtenu à partir des données observées. Si μ_0 appartient à cet intervalle, alors on ne rejette pas H_0 et inversement. En effet, on a

$$\begin{aligned} \text{"On ne rejette pas } H_0 \text{"} &\Leftrightarrow |z_{obs}| \leq t_{n-1, 1-\alpha/2} \\ &\Leftrightarrow -t_{n-1, 1-\alpha/2} \leq z_{obs} \leq t_{n-1, 1-\alpha/2} \\ &\Leftrightarrow t_{n-1, 1-\alpha/2} \leq \sqrt{n} \frac{\bar{x}_n - \mu_0}{s_n} \leq t_{n-1, 1-\alpha/2} \\ &\Leftrightarrow -t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \leq \bar{x}_n - \mu_0 \leq t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \\ &\Leftrightarrow -\bar{x}_n - t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \leq -\mu_0 \leq -\bar{x}_n + t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \\ &\Leftrightarrow \bar{x}_n - t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \leq \mu_0 \leq \bar{x}_n + t_{n-1, 1-\alpha/2} \frac{s_n}{\sqrt{n}} \\ &\Leftrightarrow \mu_0 \in IC_{1-\alpha}(\mu). \end{aligned}$$

Remarque 2 : En réalité, pour trouver la zone de rejet, comme \bar{X}_n doit être proche de μ_0 , on cherche c_α tel que sous H_0

$$\mathbb{P}_{\mu=\mu_0} [\text{"On rejette } H_0 \text{"}] = \mathbb{P}_{\mu=\mu_0} [|\bar{X}_n - \mu_0| \geq c_\alpha] = \alpha.$$

On obtient alors

$$\begin{aligned}\alpha &= \mathbb{P}_{\mu=\mu_0} [|\bar{X}_n - \mu_0| \geq c_\alpha] \\ &= \mathbb{P}_{\mu=\mu_0} \left[\sqrt{n} \frac{|\bar{X}_n - \mu_0|}{S_n} \geq \sqrt{n} \frac{c_\alpha}{S_n} \right] \\ &= \mathbb{P}_{\mu=\mu_0} \left[\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \leq -\sqrt{n} \frac{c_\alpha}{S_n} \right] + \mathbb{P}_{\mu=\mu_0} \left[\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \geq \sqrt{n} \frac{c_\alpha}{S_n} \right].\end{aligned}$$

Or si on note F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté, on obtient par symétrie de la loi de Student,

$$\begin{aligned}\alpha &= F_T \left(-\sqrt{n} \frac{c_\alpha}{S_n} \right) + 1 - F_T \left(\sqrt{n} \frac{c_\alpha}{S_n} \right) \\ &= 2 - 2F_T \left(\sqrt{n} \frac{c_\alpha}{S_n} \right)\end{aligned}$$

et donc

$$F_T \left(\sqrt{n} \frac{c_\alpha}{S_n} \right) = 1 - \alpha/2.$$

Ainsi, en appliquant F_T^{-1} ,

$$t_{n-1, 1-\alpha/2} = \sqrt{n} \frac{c_\alpha}{S_n},$$

et donc

$$c_\alpha = t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}}.$$

On retrouve alors

$$ZR = \left\{ |\bar{X}_n - \mu_0| \geq t_{n-1, 1-\alpha/2} \frac{S_n}{\sqrt{n}} \right\} = \left\{ \left| \sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \right| \geq t_{n-1, 1-\alpha/2} \right\}.$$

Remarque 3 : Soit $T \sim T_{n-1}$, alors on peut calculer la p -value comme

$$p\text{-value} = \mathbb{P} [|T| > |z_{obs}|].$$

En effet, on a

$$\begin{aligned}p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \text{"} \} \\ &= \inf \{ \alpha \in (0, 1), |z_{obs}| \geq t_{n-1, 1-\alpha/2} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(|z_{obs}|) \geq 1 - \alpha/2 \} \\ &= \inf \{ \alpha \in (0, 1), \alpha \geq 2 - 2F_T(|z_{obs}|) \} \\ &= 2 - 2F_T(|z_{obs}|) \\ &= 2\mathbb{P} [T \geq |z_{obs}|]\end{aligned}$$

où $T \sim T_{n-1}$. Par symétrie de la loi de Student,

$$p\text{-value} = 2\mathbb{P}[T \geq |z_{obs}|] = \mathbb{P}[T \geq |z_{obs}|] + \mathbb{P}[-T \geq |z_{obs}|] = \mathbb{P}[|T| \geq |z_{obs}|].$$

5.2.2 Test d'inégalité $\mu \leq \mu_0$

On veut tester au risque α l'hypothèse nulle $H_0 : "\mu \leq \mu_0"$ contre l'hypothèse alternative $H_1 : "\mu > \mu_0"$. Sous H_0 , il existe $\mu' \leq \mu_0$ tel que

$$Z(\mu') = \sqrt{n} \frac{\bar{X}_n - \mu'}{S_n} \sim T_{n-1}.$$

On définit la zone de rejet par

$$ZR = \{Z(\mu_0) > t_{n-1,1-\alpha}\}$$

où $t_{n-1,1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de Student à $n - 1$ degrés de liberté. On calcule alors $z_{obs} = \sqrt{n} \frac{\bar{x}_n - \mu_0}{s_n}$.

Si $z_{obs} > t_{n-1,1-\alpha}$, on rejette H_0 , sinon, on ne rejette pas H_0 .

Remarque 1 : En réalité, pour trouver notre zone de rejet, comme μ peut être aussi petit que l'on veut sous H_0 , on doit rejeter H_0 si \bar{X}_n est trop grand, i.e on cherche c_α tel que sous H_0

$$\sup_{\mu \leq \mu_0} \mathbb{P}_\mu [\text{"On rejette } H_0"] = \sup_{\mu \leq \mu_0} \mathbb{P}_\mu [\bar{X}_n \geq c_\alpha] = \alpha$$

Or, on a

$$\begin{aligned} \alpha &= \sup_{\mu \leq \mu_0} \mathbb{P}_\mu [\bar{X}_n \geq c_\alpha] \\ &= \sup_{\mu \leq \mu_0} \mathbb{P}_\mu [\bar{X}_n - \mu \geq c_\alpha - \mu] \\ &= \sup_{\mu \leq \mu_0} \mathbb{P}_\mu \left[\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \geq \sqrt{n} \frac{c_\alpha - \mu}{S_n} \right]. \end{aligned}$$

Comme $\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \sim T_{n-1}$, si on note F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté, on obtient

$$\begin{aligned} \alpha &= \sup_{\mu \leq \mu_0} \left(1 - F_T \left(\sqrt{n} \frac{c_\alpha - \mu}{S_n} \right) \right) \\ &= 1 - \inf_{\mu \leq \mu_0} F_T \left(\sqrt{n} \frac{c_\alpha - \mu}{S_n} \right) \end{aligned}$$

Comme la fonction F_T est strictement croissante, on obtient

$$1 - \alpha = F_T \left(\sqrt{n} \frac{c_\alpha - \mu_0}{S_n} \right)$$

et donc, en appliquant F_T^{-1} ,

$$t_{n-1,1-\alpha} = \sqrt{n} \frac{c_\alpha - \mu_0}{S_n}$$

et on obtient alors

$$c_\alpha = \mu_0 + t_{n-1,1-\alpha} \frac{S_n}{\sqrt{n}}.$$

Donc, sous H_0

$$\sup_{\mu \leq \mu_0} \mathbb{P} [\bar{X}_n \geq c_\alpha] = \mathbb{P}_{\mu=\mu_0} \left[\bar{X}_n \geq \mu_0 + t_{n-1,1-\alpha} \frac{S_n}{\sqrt{n}} \right] = \mathbb{P}_{\mu=\mu_0} \left[\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \geq t_{n-1,1-\alpha} \right] = \alpha,$$

et on retrouve ainsi la zone de rejet.

Remarque 2 : On a

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \geq t_{n-1,1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \geq 1 - \alpha \} \\ &= \inf \{ \alpha \in (0, 1), \alpha \geq 1 - F_T(z_{obs}) \} \\ &= 1 - F_T(z_{obs}) \\ &= \mathbb{P} [T \geq z_{obs}], \end{aligned}$$

où $T \sim T_{n-1}$.

5.2.3 Test d'inégalité $\mu \geq \mu_0$

On veut tester au risque α l'hypothèse nulle $H_0 : \mu \geq \mu_0$ contre l'hypothèse alternative $H_1 : \mu < \mu_0$. Sous H_0 , il existe $\mu' \geq \mu_0$ tel que

$$Z(\mu') = \sqrt{n} \frac{\bar{X}_n - \mu'}{S_n} \sim T_{n-1}.$$

On définit la zone de rejet par

$$ZR = \{ Z(\mu_0) < t_{n-1,\alpha} \}$$

où $t_{n-1,\alpha}$ est le quantile d'ordre α de la loi de Student à $n - 1$ degrés de libertés. On calcule alors $z_{obs} = \sqrt{n} \frac{\bar{X}_n - \mu_0}{s_n}$. Si $z_{obs} < t_{n-1,\alpha}$, on rejette H_0 , sinon, on ne rejette pas H_0 .

Remarque 1 : En réalité, pour trouver notre zone de rejet, comme μ peut être aussi grand que l'on veut sous H_0 , on doit rejeter H_0 si \bar{X}_n est trop petit, i.e on cherche c_α tel que

$$\sup_{\mu \geq \mu_0} \mathbb{P}_\mu [\text{"On rejette } H_0 \}] = \sup_{\mu \geq \mu_0} \mathbb{P}_\mu [\bar{X}_n \leq c_\alpha] = \alpha$$

Or, si on note F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté, on a

$$\begin{aligned}\alpha &= \sup_{\mu \geq \mu_0} \mathbb{P}_\mu [\bar{X}_n \leq c_\alpha] \\ &= \sup_{\mu \geq \mu_0} \mathbb{P}_\mu [\bar{X}_n - \mu \leq c_\alpha - \mu] \\ &= \sup_{\mu \geq \mu_0} \mathbb{P}_\mu \left[\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \leq \sqrt{n} \frac{c_\alpha - \mu}{S_n} \right] \\ &= \sup_{\mu \geq \mu_0} F_T \left(\sqrt{n} \frac{c_\alpha - \mu}{S_n} \right)\end{aligned}$$

Comme la fonction F_T est strictement croissante, on obtient

$$\alpha = F_T \left(\sqrt{n} \frac{c_\alpha - \mu_0}{S_n} \right)$$

et donc, en appliquant F_T^{-1} ,

$$t_{n-1, \alpha} = \sqrt{n} \frac{c_\alpha - \mu_0}{S_n}.$$

Comme $t_{n-1, \alpha} = -t_{n-1, 1-\alpha}$, on obtient alors

$$c_\alpha = \mu_0 - t_{n-1, 1-\alpha} \frac{S_n}{\sqrt{n}}$$

D'où, sous H_0

$$\sup_{\mu \geq \mu_0} \mathbb{P}_\mu [\bar{X}_n \leq c_\alpha] = \mathbb{P}_{\mu=\mu_0} \left[\bar{X}_n \leq \mu_0 - t_{n-1, 1-\alpha} \frac{S_n}{\sqrt{n}} \right] = \mathbb{P}_{\mu=\mu_0} \left[\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \leq -t_{n-1, 1-\alpha} \right] = \alpha,$$

et on retrouve ainsi la zone de rejet.

Remarque 2 : On a

$$\begin{aligned}p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \leq -t_{n-1, 1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \leq \alpha \} \\ &= F_T(z_{obs}) \\ &= \mathbb{P}[T \leq z_{obs}],\end{aligned}$$

où $T \sim T_{n-1}$.

5.2.4 Application :

Les données suivantes sont extraites du cours de B. Portier à l'INSA de Rouen. A la suite d'un traitement (régime alimentaire) sur une variété de porcs, on prélève un (petit) échantillon de 5

porcs et on les pèse. On obtient les poids suivants (en kg)

83 81 84 80 85

On suppose que les x_i sont des réalisations de variables aléatoires i.i.d suivant une loi normale de paramètres μ, σ^2 . On sait que le poids moyen de cette variété de porcs est de 87,6 kg et on souhaite savoir si le régime alimentaire a eu un impact ou non sur ce poids moyen, i.e on teste au risque de 5%, $H_0 : \mu = 87,6$ contre $H_1 : \mu \neq 87,6$. On a la statistique de test

$$Z = \sqrt{5} \frac{\bar{X}_5 - 87.6}{S_n} \sim T_4 \quad \text{sous } H_0.$$

On a donc la zone de rejet

$$ZR = \{|Z| \geq t_{4,0.975}\}.$$

Ici, $t_{4,0.975} = 2.78$. De plus, on a

$$z_{obs} = \sqrt{5} \frac{\bar{x}_5 - 87.6}{s_5} = \sqrt{5} \frac{82.6 - 87.5}{2.074} = -5.28.$$

Comme $|z_{obs}| > 2.78$, on rejette H_0 .

Un biologiste affirme que ce régime permet d'augmenter le poids des porcs, et que ces résultats ne sont dus qu'à la fluctuation d'échantillonnage. Pour le contredire, on teste alors, au risque de 5%, $H_0 : \mu \geq 87.6$ contre $H_1 : \mu < 87.6$. Sous H_0 , il existe $\mu' \geq 87.6$ tel que

$$\sqrt{5} \frac{\bar{X}_5 - \mu'}{S_5} \sim T_4.$$

On a la zone de rejet définie par

$$ZR = \left\{ \sqrt{5} \frac{\bar{X}_5 - 87.6}{S_5} < -t_{4,0.95} \right\}$$

Ici, $t_{4,0.95} = 2.13$ et

$$\sqrt{5} \frac{\bar{x}_5 - 87.6}{s_5} = \sqrt{5} \frac{82.6 - 87.5}{2.074} = -5.28.$$

Comme $-5.28 < -2.13$, on rejette H_0 .

5.2.5 Test de conformité d'une variance

On souhaite maintenant tester, au risque α , l'hypothèse nulle $H_0 : \sigma^2 = \sigma_0^2$ contre l'hypothèse alternative $H_1 : \sigma^2 \neq \sigma_0^2$.

Construction de la statistique de test : On rappelle que l'on a

$$\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Comme sous H_0 $\sigma^2 = \sigma_0^2$, on a alors la statistique de test

$$\frac{(n-1)S_n^2}{\sigma_0^2} \sim \chi_{n-1}^2 \quad \text{sous } H_0$$

ce qui est bien évidemment faux sous H_1 . En effet, sous H_1 ,

$$\frac{n-1}{\sigma_0^2} S_n^2 = \frac{\sigma^2}{\sigma_0^2} \frac{n-1}{\sigma^2} S_n^2$$

et ne suit donc pas une loi du Chi-deux à $n-1$ degrés de liberté.

Le test : On teste au risque α , $H_0 : \sigma^2 = \sigma_0^2$ contre $\sigma^2 \neq \sigma_0^2$. On a la statistique de test

$$Z = \frac{n-1}{\sigma_0^2} S_n^2 \sim \chi_{n-1}^2 \quad \text{sous } H_0$$

On définit alors la zone de rejet comme

$$ZR = \{Z, Z < k_{\alpha/2}\} \cup \{Z, Z > k_{1-\alpha/2}\},$$

où $k_{\alpha/2}$ et $k_{1-\alpha/2}$ sont les quantiles d'ordre $\alpha/2$ et $1-\alpha/2$ de la loi du chi deux à $n-1$ degrés de liberté. On calcule z_{obs} , et si z_{obs} est dans la zone de rejet, on rejette H_0 et inversement.

Remarque 1 : Notons que faire ce test revient donc à vérifier que σ_0^2 appartient à l'intervalle de confiance obtenu à partir des données observées. Si σ_0^2 appartient à cet intervalle, alors on ne rejette pas H_0 et inversement. En effet, on a

$$\begin{aligned} \text{"On ne rejette pas } H_0 \text{"} &\Leftrightarrow k_{\alpha/2} \leq |z_{obs}| \leq k_{1-\alpha/2} \\ &\Leftrightarrow k_{\alpha/2} \leq \frac{(n-1)S_n^2}{\sigma_0^2} \leq k_{1-\alpha/2} \\ &\Leftrightarrow \frac{k_{\alpha/2}}{(n-1)S_n^2} \leq \frac{1}{\sigma_0^2} \leq \frac{k_{1-\alpha/2}}{(n-1)S_n^2} \\ &\Leftrightarrow \frac{(n-1)S_n^2}{k_{\alpha/2}} \geq \sigma_0^2 \geq \frac{(n-1)S_n^2}{k_{1-\alpha/2}} \\ &\Leftrightarrow \sigma_0^2 \in IC_{1-\alpha}(\sigma^2). \end{aligned}$$

Remarque 2 : Soit $Z \sim \chi_{n-1}^2$, et F_Z sa fonction de répartition, alors

$$\begin{aligned}
 p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \text{"} \} \\
 &= \inf \{ \alpha \in (0, 1), k_{\alpha/2} \geq z_{obs} \text{ ou } z_{obs} \geq k_{1-\alpha/2} \} \\
 &= \min (\inf \{ \alpha \in (0, 1), k_{\alpha/2} \geq z_{obs} \}, \inf \{ \alpha \in (0, 1), k_{1-\alpha/2} \leq z_{obs} \}) \\
 &= \min (\inf \{ \alpha \in (0, 1), \alpha/2 \geq F(z_{obs}) \}, \inf \{ \alpha \in (0, 1), 1 - \alpha/2 \leq F(z_{obs}) \}) \\
 &= \min (2F(z_{obs}), 2 - 2F(z_{obs})) \\
 &= \min (2\mathbb{P}[Z \leq z_{obs}], 2\mathbb{P}[Z \geq z_{obs}]).
 \end{aligned}$$

Exemple : Reprenons l'exemple des porcs. Le même biologiste pense que si les résultats ne sont pas concluants (i.e ils ne vont pas dans son sens), c'est dû à une variance qui serait, selon lui, égale à 25. Sans s'attarder sur le bien fondé de son raisonnement, on teste au risque de 5% l'hypothèse nulle $H_0 : \sigma^2 = 25$ contre l'hypothèse alternative $H_1 : \sigma^2 \neq 25$. On a la statistique de test

$$Z = \frac{4}{25} S_n^2 \sim \chi_4^2 \text{ sous } H_0$$

On a la zone de rejet

$$ZR = \{Z, Z \leq 0.48 \text{ ou } Z \geq 11.14\}.$$

De plus, on a

$$z_{obs} = \frac{4}{25} s_n^2 = 0.29.$$

Comme $z_{obs} < 0.48$, on rejette H_0 .

5.3 Tests de comparaison de deux moyennes

5.3.1 Introduction

Présentation du problème : On considère deux jeux de données

- x_1, \dots, x_p qui sont p mesures d'une variable aléatoire X .
- y_1, \dots, y_q qui sont q mesures d'une variable aléatoire Y .

On souhaite étudier les moyennes théoriques des variables X et Y , moyennes que l'on peut respectivement estimer par \bar{x}_p et \bar{y}_q . Bien évidemment, il est peu probable que ces deux estimations soient égales, et l'objectif est donc de savoir si cette différence est seulement due à la fluctuation d'échantillonnage, ou si elle est réellement significative.

Le modèle probabiliste : On fait les hypothèses suivantes :

- les données x_1, \dots, x_p sont les réalisations de variables aléatoires X_1, \dots, X_p indépendantes et de même loi $\mathcal{N}(\mu_1, \sigma_1^2)$
- les données y_1, \dots, y_q sont les réalisations de variables aléatoires Y_1, \dots, Y_q indépendantes et de même loi $\mathcal{N}(\mu_2, \sigma_2^2)$

— les deux échantillons sont indépendants et de même variance, i.e $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

Exemple : Cet exemple est issu du cours de B. Portier à l'INSA de Rouen. On souhaite comparer les productions laitières de deux races bovines. On choisit donc 50 bêtes de chaque race, i.e $p = q = 50$, et on mesure pour chaque bête la production annuelle totale de lait (en tonnes). On obtient pour ces deux échantillons les rendements moyens suivants :

$$\bar{x}_{50} = 4.860 \quad \text{et} \quad \bar{y}_{50} = 4.690$$

L'objectif sera donc de savoir si la différence entre ces deux moyennes est due à la fluctuation d'échantillonnage ou si il y a une réelle différence de rendement entre les deux races.

5.3.2 Test d'égalité

On veut tester au risque α l'hypothèse nulle $H_0 : \mu_1 = \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 \neq \mu_2$.

Construction de la statistique de test : On dispose des estimateurs \bar{X}_p et \bar{Y}_q et on s'intéresse plus particulièrement à la variable aléatoire $\bar{X}_p - \bar{Y}_q$.

Proposition 5.3.1. *On a*

$$\bar{X}_p - \bar{Y}_q \sim \mathcal{N} \left(\mu_1 - \mu_2, \frac{\sigma^2}{p} + \frac{\sigma^2}{q} \right).$$

Démonstration. On peut voir $\bar{X}_p - \bar{Y}_q$ comme une combinaison linéaire de variables aléatoires indépendantes suivant des lois normales, i.e

$$\bar{X}_p - \bar{Y}_q = \frac{1}{p} \sum_{i=1}^p X_i - \frac{1}{q} \sum_{i=1}^q Y_i,$$

et suit donc une loi normale avec

$$\mathbb{E} [\bar{X}_p - \bar{Y}_q] = \frac{1}{p} \sum_{i=1}^p \mathbb{E} [X_i] - \frac{1}{q} \sum_{i=1}^q \mathbb{E} [Y_i] = \mu_1 - \mu_2$$

et comme \bar{X}_n et \bar{Y}_n sont indépendantes,

$$\mathbb{V} [\bar{X}_p - \bar{Y}_q] = \mathbb{V} [\bar{X}_p] + \mathbb{V} [\bar{Y}_q] = \frac{\sigma^2}{p} + \frac{\sigma^2}{q}.$$

□

On obtient ainsi

$$\frac{\bar{X}_p - \mu_1 + \mu_2 - \bar{Y}_q}{\sigma \sqrt{\frac{1}{p} + \frac{1}{q}}} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \mu_1 + \mu_2 - \bar{Y}_q}{\sigma} \sim \mathcal{N}(0, 1).$$

Malheureusement, on ne connaît pas nécessairement σ^2 et donc pas tous les paramètres de la loi de $\bar{X}_p - \bar{Y}_q$. Cependant, on peut le remplacer par un estimateur. La proposition suivante donne cet estimateur ainsi que sa loi.

Proposition 5.3.2. *On considère S^2 l'estimateur de σ^2 défini par*

$$S^2 = \frac{1}{p+q-2} \left(\sum_{i=1}^p (X_i - \bar{X}_p)^2 + \sum_{i=1}^q (Y_i - \bar{Y}_q)^2 \right) = \frac{1}{p+q-2} ((p-1)S_X^2 + (q-1)S_Y^2),$$

avec $S_X^2 = \frac{1}{p-1} \sum_{i=1}^p (X_i - \bar{X}_p)^2$ et $S_Y^2 = \frac{1}{q-1} \sum_{i=1}^q (Y_i - \bar{Y}_q)^2$.

1. On a

$$\frac{(p+q-2)S^2}{\sigma^2} \sim \chi_{p+q-2}^2.$$

2. Les variables $\bar{X}_p - \bar{Y}_q$ et S^2 sont indépendantes.

Démonstration. Comme

$$\frac{(p-1)S_X^2}{\sigma^2} \sim \chi_{p-1}^2 \quad \text{et} \quad \frac{(q-1)S_Y^2}{\sigma^2} \sim \chi_{q-1}^2,$$

en appliquant la proposition 4.3.1, on obtient

$$\frac{p+q-2}{\sigma^2} S^2 = \frac{(p-1)S_X^2}{\sigma^2} + \frac{(q-1)S_Y^2}{\sigma^2} \sim \chi_{p+q-2}^2.$$

De plus comme S_X^2, \bar{X}_q, S_Y^2 et \bar{Y}_q sont indépendants, S^2 et $\bar{X}_p - \bar{Y}_q$ sont indépendants. \square

Version 2. On note

$$Z = \begin{pmatrix} X_1 - \mu_1 \\ \vdots \\ X_p - \mu_1 \\ Y_1 - \mu_2 \\ \vdots \\ Y_q - \mu_2 \end{pmatrix} \sim \mathcal{N}(0, \sigma^2 I_{p+q}).$$

L'objectif est donc d'utiliser le théorème de Cochran et donc de trouver une décomposition en sous-espaces orthogonaux adéquate. On note $\{e_1, \dots, e_{p+q}\}$ la base canonique de \mathbb{R}^{p+q} , et on considère les sous-espaces vectoriels

$$V = \text{vect} \{e_1, \dots, e_p\} \quad \text{et} \quad W = \text{vect} \{e_{p+1}, \dots, e_{p+q}\}.$$

A noter que l'on a $\mathbb{R}^{p+q} = V \oplus_{\perp} W$ et que les projecteurs orthogonaux sur les sous-espaces sont

définis pour tout h par

$$P_V : h \mapsto \sum_{i=1}^p \langle e_i, h \rangle e_i \quad \text{et} \quad P_W : h \mapsto \sum_{i=p+1}^{p+q} \langle e_i, h \rangle e_i.$$

On considère également $v = \frac{1}{\sqrt{p}} \left(\mathbf{1}_p^T, 0_q^T \right)^T$ avec $\mathbf{1}_p^T = \underbrace{(1, \dots, 1)}_p$ et $0_q^T = \underbrace{(0, \dots, 0)}_q$, ainsi que $w = \frac{1}{\sqrt{q}} (0_p, \mathbf{1}_q)^T$. On note $V_1 = \text{vect}\{v\}$ et $W_1 = \text{vect}\{w\}$ et on a donc V_1 (resp. W_1) qui est un sous-espace vectoriel de V (resp. W) et on note $V_{1,\perp}$ (resp. $W_{2,\perp}$) le sous-espace vectoriel tel que $V = V_1 \oplus_{\perp} V_{1,\perp}$ (resp. $W = W_1 \oplus_{\perp} W_{1,\perp}$). A noter que les projections orthogonales sont donc données par

$$P_{V_1} = vv^T : h \mapsto (v^T h) v \quad P_{W_1} = ww^T : h \mapsto (w^T h) w$$

On obtient donc la décomposition suivante

$$\mathbb{R}^{p+q} = V_1 \oplus_{\perp} V_{1,\perp} \oplus_{\perp} W_1 \oplus_{\perp} W_{1,\perp}.$$

Il ne reste donc plus qu'à calculer les projections. On a facilement

$$P_{V_1}(Z) = (\bar{X}_p - \mu_1) (1, \dots, 1, 0, \dots, 0)^T \quad \text{et} \quad P_{W_1} = (\bar{Y}_q - \mu_q) (0, \dots, 0, 1, \dots, 1)^T.$$

De plus, on a

$$\begin{aligned} P_{V_{1,\perp}}(Z) &= P_V(Z) - P_{V_1}(Z) = ((X_1 - \mu_1), \dots, (X_p - \mu_1), 0, \dots, 0)^T - (\bar{X}_p - \mu_1) (1, \dots, 1, 0, \dots, 0)^T \\ &= ((X_1 - \bar{X}_p), \dots, (X_p - \bar{X}_p), 0, \dots, 0)^T \end{aligned}$$

De la même façon, on a

$$P_{W_{1,\perp}}(Z) = (0, \dots, 0, Y_1 - \bar{Y}_q, \dots, Y_q - \bar{Y}_q)^T.$$

Ainsi, en appliquant le théorème de Cochran,

- $\bar{X}_n, \bar{Y}_n, P_{V_{1,\perp}}(Z), P_{W_{1,\perp}}(Z)$ sont indépendants.
- On a

$$\begin{aligned} \frac{1}{\sigma^2} \|P_{V_{1,\perp}}\|^2 &= \frac{1}{\sigma^2} \sum_{i=1}^p (X_i - \bar{X}_p)^2 = \frac{p-1}{\sigma^2} S_X^2 \sim \chi_{p-1}^2 \\ \frac{1}{\sigma^2} \|P_{W_{1,\perp}}\|^2 &= \frac{1}{\sigma^2} \sum_{i=1}^q (Y_i - \bar{Y}_q)^2 = \frac{q-1}{\sigma^2} S_Y^2 \sim \chi_{q-1}^2 \end{aligned}$$

On a donc indépendance de S^2, \bar{X}_p et \bar{Y}_q . Comme $P_{V_{1,\perp}}(Z)$ et $P_{W_{1,\perp}}(Z)$ sont indépendants,

en appliquant la proposition 4.3.1, on obtient

$$\frac{1}{\sigma^2} ((p-1)S_X^2 + (q-1)S_Y^2) \sim \chi_{p+q-2}^2$$

□

Version 3. On considère U le sous-espace vectoriel tel que

$$\mathbb{R}^{p+q} = V_1 \oplus_{\perp} W_1 \oplus_{\perp} U.$$

La projection orthogonale de Z sur U est donnée par

$$P_U(Z) = Z - P_{V_1}(Z) - P_{W_1}(Z) = ((X_1 - \bar{X}_p), \dots, (X_p - \bar{X}_p), (Y_1 - \bar{Y}_q), \dots, (Y_q - \bar{Y}_q)).$$

Ainsi

$$\|P_U(Z)\|^2 = \sum_{i=1}^p (X_i - \bar{X}_p)^2 + \sum_{i=1}^q (Y_i - \bar{Y}_q)^2 = (p+q-2)S^2.$$

En appliquant le théorème de Cochran, on obtient donc

- \bar{X}_p, \bar{Y}_q et S^2 sont indépendants.
- $\frac{1}{\sigma^2} \left(\sum_{i=1}^p (X_i - \mu_1)^2 + \sum_{i=1}^q (Y_i - \mu_2)^2 \right) \sim \chi_{p+q-2}^2$.

□

Le corollaire suivant est une application directe de la proposition précédente.

Corollaire 5.3.1. On a

$$\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{(\bar{X}_p - \bar{Y}_q) - (\mu_1 - \mu_2)}{S} \sim T_{p+q-2}$$

où T_{p+q-2} est une loi de Student à $p+q-2$ degrés de liberté.

Démonstration. On a

$$\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{(\bar{X}_p - \bar{Y}_q) - (\mu_1 - \mu_2)}{S} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{(\bar{X}_p - \bar{Y}_q) - (\mu_1 - \mu_2)}{\sigma} \frac{1}{\sqrt{\frac{(p+q-2)S^2}{\sigma^2(p+q-2)}}} \sim T_{p+q-2}$$

car S^2 et $\bar{X}_p - \bar{Y}_q$ sont indépendants.

□

On obtient donc la statistique de test

$$\frac{\sqrt{pq}(\bar{X}_p - \bar{Y}_q)}{S\sqrt{p+q}} \sim T_{p+q-2} \quad \text{sous } H_0,$$

ce qui est évidemment faux sous H_1 car alors $\mu_1 - \mu_2 \neq 0$.

Le test :

— On teste au risque α

$$H_0 : " \mu_1 = \mu_2 " \quad \text{contre} \quad H_1 : " \mu_1 \neq \mu_2 " .$$

— On a la statistique de test

$$Z = \frac{\sqrt{pq} (\bar{X}_p - \bar{Y}_q)}{S\sqrt{p+q}} \sim T_{p+q-2} \quad \text{sous } H_0.$$

— On a alors la zone de rejet :

$$ZR = \{Z, |Z| > t_{p+q-2, 1-\alpha/2}\}$$

où $t_{p+q-2, 1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi de Student à $p + q - 2$ degrés de liberté.

— Si $|z_{obs}| > t_{p+q-2, 1-\alpha/2}$, on rejette H_0 et inversement.

Remarque 1 : Cela revient à vérifier que 0 appartient à l'intervalle de confiance de niveau $1 - \alpha$ de $\mu = \mu_1 - \mu_2$. En effet

$$\begin{aligned} |z_{obs}| \leq t_{p+q-2, 1-\alpha/2} &\Leftrightarrow -t_{p+q-2, 1-\alpha/2} \leq \frac{\sqrt{pq} (\bar{x}_p - \bar{y}_q)}{s\sqrt{p+q}} \leq t_{p+q-2, 1-\alpha/2} \\ &\Leftrightarrow \bar{x}_p - \bar{y}_q - \frac{st_{p+q-2, 1-\alpha/2}\sqrt{p+q}}{\sqrt{pq}} \leq 0 \leq \bar{x}_p - \bar{y}_q + \frac{st_{p+q-2, 1-\alpha/2}\sqrt{p+q}}{\sqrt{pq}} \\ &\Leftrightarrow 0 \in IC_{1-\alpha}(\mu_1 - \mu_2). \end{aligned}$$

Remarque 2 : En réalité, pour trouver la zone de rejet, comme $\bar{X}_p - \bar{Y}_q$ doit être proche de 0, on cherche c_α tel que

$$\mathbb{P}_{\mu_1=\mu_2} ["\text{On rejette } H_0"] = \mathbb{P}_{\mu_1=\mu_2} [|\bar{X}_p - \bar{Y}_q| \geq c_\alpha] = \alpha.$$

On obtient donc

$$\begin{aligned} \alpha &= \mathbb{P}_{\mu_1=\mu_2} [|\bar{X}_p - \bar{Y}_q| \geq c_\alpha] \\ &= \mathbb{P}_{\mu_1=\mu_2} \left[\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{|\bar{X}_p - \bar{Y}_q|}{S} \geq \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S} \right], \end{aligned}$$

et on obtient donc, en notant F_T la fonction de répartition de la loi de Student à $p + q - 2$ degrés de liberté, sous H_0 ,

$$\alpha = 2 - 2F_T \left(\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S} \right).$$

En appliquant F_T^{-1} , on obtient

$$t_{p+q-2, 1-\alpha/2} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S},$$

et donc

$$c_\alpha = t_{p+q-2, 1-\alpha/2} \frac{\sqrt{p+q}}{\sqrt{pq}} S.$$

On retrouve alors la zone de rejet, i.e

$$ZR = \left\{ |\bar{X}_p - \bar{Y}_q| \geq t_{p+q-2, 1-\alpha/2} \frac{\sqrt{p+q}}{\sqrt{pq}} S \right\} = \left\{ \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{|\bar{X}_p - \bar{Y}_q|}{S} \geq t_{p+q-2, 1-\alpha/2} \right\}.$$

Remarque 3 : Soit $T \sim T_{p+q-2}$, on a alors

$$p\text{-value} = \mathbb{P}[|T| > |z_{obs}|].$$

En effet, on note F_T la fonction de répartition de T ,

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), |z_{obs}| \geq t_{p+q-2, 1-\alpha/2} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(|z_{obs}|) \geq 1 - \alpha/2 \} \\ &= 2 - 2F_T(|z_{obs}|) \\ &= \mathbb{P}[|T| \geq |z_{obs}|], \end{aligned}$$

la dernière égalité résultant de la symétrie de la loi de Student.

Exemple : les Iris de Fisher On considère un jeu de données "classique". L'objectif est de reconnaître le type d'Iris à partir de la longueur de ses sépales. On considère ici deux espèces, et pour chaque espèce on dispose de 50 individus. Plus précisément, on note x_1, \dots, x_{50} les longueurs des sépales des iris de la variété *Virginica* et y_1, \dots, y_{50} celle de la variété *Versicolor*. On obtient les résultats suivants

$$\begin{array}{ll} \bar{x}_{50} = 5.94 & s_X = 0.52 \\ \bar{y}_{50} = 6.59 & s_Y = 0.64 \end{array}$$

On teste au risque 5% l'hypothèse nulle $H_0 : \mu_1 = \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 \neq \mu_2$. On a la statistique de test

$$Z = \frac{\sqrt{50 \times 50} (\bar{X}_{50} - \bar{Y}_{50})}{S \sqrt{50 + 50}} \sim T_{98} \text{ sous } H_0.$$

Ici $z_{obs} = -5.63$. On a la zone de rejet

$$ZR = \{|Z| > t_{98, 0.975}\}.$$

Ici $t_{98,0.975} = 1.98$. Donc $|z_{obs}| > t_{98,0.975}$ et on rejette donc H_0 , i.e on peut donc penser que la longueur moyenne des sépales des Iris dépend de la variété.

5.3.3 Test d'inégalité $\mu_1 \leq \mu_2$

On veut tester au risque α l'hypothèse nulle $H_0 : \mu_1 \leq \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 > \mu_2$. Sous H_0 , il existe $\mu' \leq 0$ tel que

$$Z(\mu') = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q - \mu'}{S} \sim T_{p+q-2}.$$

On définit alors la zone de rejet par

$$ZR = \{Z(0) > t_{p+q-2,1-\alpha}\},$$

où $t_{p+q-2,1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de Student à $p + q - 2$ degrés de liberté. On calcule $z_{obs} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{x}_p - \bar{y}_q}{s}$. Si $z_{obs} > t_{p+q-2,1-\alpha}$, on rejette H_0 et inversement.

Remarque 1 : En réalité, pour trouver la zone de rejet, comme $\mu = \mu_1 - \mu_2$ peut être aussi petit que l'on veut sous H_0 , on doit rejeter H_0 si $\bar{X}_p - \bar{Y}_q$ est trop grand, i.e on cherche c_α tel que

$$\sup_{\mu \leq 0} \mathbb{P}_\mu [\text{"On rejette } H_0\text{"}] = \sup_{\mu \leq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q \geq c_\alpha] = \alpha.$$

Or si on note F_T la fonction de répartition de la loi de Student à $p + q - 2$ degrés de liberté, on a

$$\begin{aligned} \alpha &= \sup_{\mu \leq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q \geq c_\alpha] \\ &= \sup_{\mu \leq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q - \mu \geq c_\alpha - \mu] \\ &= \sup_{\mu \leq 0} \mathbb{P}_\mu \left[\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q - \mu}{S} \geq \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha - \mu}{S} \right] \\ &= \sup_{\mu \leq 0} \left(1 - F_T \left(\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha - \mu}{S} \right) \right) \end{aligned}$$

Comme la fonction F_T est strictement croissante, on obtient

$$\alpha = 1 - F_T \left(\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S} \right)$$

et donc, en appliquant F_T^{-1} ,

$$t_{p+q-2,1-\alpha} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S},$$

et donc

$$c_\alpha = \frac{\sqrt{p+q}}{\sqrt{pq}} t_{p+q-2,1-\alpha} S$$

On retrouve ainsi la zone de rejet

$$ZR = \left\{ \bar{X}_p - \bar{Y}_q \geq t_{p+q-2, 1-\alpha} \frac{\sqrt{p+q}}{\sqrt{pq}} S \right\} = \left\{ \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q}{S} \geq t_{p+q-2, 1-\alpha} \right\}.$$

Remarque 2 : On a

$$p - \text{value} = \mathbb{P} [T \geq z_{obs}],$$

où $T \sim T_{p+q-2}$. En effet,

$$\begin{aligned} p - \text{value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \geq t_{p+q-2, 1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \geq 1 - \alpha \} \\ &= 1 - F_T(z_{obs}). \end{aligned}$$

5.3.4 Test d'inégalité $\mu_1 \geq \mu_2$

On veut tester au risque α l'hypothèse nulle $H_0 : \mu_1 \geq \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 < \mu_2$. Sous H_0 , il existe $\mu' \geq 0$ tel que

$$Z(\mu') = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q - \mu'}{S} \sim T_{p+q-2}.$$

On définit alors la zone de rejet par

$$ZR = \{ Z(0) < -t_{p+q-2, 1-\alpha} \},$$

où $t_{p+q-2, 1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de Student à $p + q - 2$ degrés de liberté. On calcule $z_{obs} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{x}_p - \bar{y}_q}{s}$. Si $z_{obs} < -t_{p+q-2, 1-\alpha}$, on rejette H_0 et inversement.

Remarque 1 : En réalité, pour trouver la zone de rejet, comme $\mu = \mu_1 - \mu_2$ peut être aussi grand que l'on veut sous H_0 , on doit rejeter H_0 si $\bar{X}_p - \bar{Y}_q$ est trop petit, i.e on cherche c_α tel que

$$\sup_{\mu \geq 0} \mathbb{P}_\mu [\text{"On rejette } H_0 \}] = \sup_{\mu \geq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q \leq c_\alpha] = \alpha.$$

Or si on note F_T la fonction de répartition de la loi de Student à $p + q - 2$ degrés de liberté, on a

$$\begin{aligned}\alpha &= \sup_{\mu \geq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q \leq c_\alpha] \\ &= \sup_{\mu \geq 0} \mathbb{P}_\mu [\bar{X}_p - \bar{Y}_q - \mu \leq c_\alpha - \mu] \\ &= \sup_{\mu \geq 0} \mathbb{P}_\mu \left[\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q - \mu}{S} \leq \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha - \mu}{S} \right] \\ &= \sup_{\mu \geq 0} F_T \left(\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha - \mu}{S} \right)\end{aligned}$$

Comme la fonction F_T est strictement croissante, on obtient

$$\alpha = F_T \left(\frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S} \right)$$

et donc, en appliquant F_T^{-1} ,

$$t_{p+q-2,\alpha} = \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{c_\alpha}{S},$$

et donc, comme $t_{p+q-2,\alpha} = -t_{p+q-2,1-\alpha}$,

$$c_\alpha = -\frac{\sqrt{p+q}}{\sqrt{pq}} t_{p+q-2,1-\alpha} S.$$

On retrouve ainsi la zone de rejet,

$$ZR = \left\{ \bar{X}_p - \bar{Y}_q \leq -t_{p+q-2,1-\alpha} \frac{\sqrt{p+q}}{\sqrt{pq}} S \right\} = \left\{ \frac{\sqrt{pq}}{\sqrt{p+q}} \frac{\bar{X}_p - \bar{Y}_q}{S} \leq t_{p+q-2,\alpha} \right\}.$$

Remarque 2 : On a

$$p\text{-value} = \mathbb{P}[T \leq z_{obs}],$$

où $T \sim T_{p+q-2}$. En effet,

$$\begin{aligned}p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \leq t_{p+q-2,\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \leq \alpha \} \\ &= F_T(z_{obs}).\end{aligned}$$

5.4 Test de Fischer et test de Shapiro-Wilk

On rappelle que pour faire les tests de Student, on a du faire les hypothèses suivantes :

- les données x_1, \dots, x_p sont les réalisations de variables aléatoires X_1, \dots, X_p indépendantes

et de même loi $\mathcal{N}(\mu_1, \sigma_1^2)$

- les données y_1, \dots, y_q sont les réalisations de variables aléatoires Y_1, \dots, Y_q indépendantes et de même loi $\mathcal{N}(\mu_2, \sigma_2^2)$
- les deux échantillons sont indépendants et de même variance, i.e $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

L'indépendance des variables est souvent dépendante du protocole expérimental. L'objectif de cette section et des prochaines est de donner des tests permettant de vérifier chacune des autres hypothèses, i.e que les deux variables X et Y suivent des lois normales, et qu'elles ont la même variance.

5.4.1 Test de Fischer

On veut tester au risque α l'hypothèse nulle $H_0 : \sigma_1^2 = \sigma_2^2$ contre l'hypothèse alternative $H_1 : \sigma_1^2 \neq \sigma_2^2$.

Le cadre probabiliste :

- les données x_1, \dots, x_p sont les réalisations de variables aléatoires X_1, \dots, X_p indépendantes et de même loi $\mathcal{N}(\mu_1, \sigma_1^2)$
- les données y_1, \dots, y_q sont les réalisations de variables aléatoires Y_1, \dots, Y_q indépendantes et de même loi $\mathcal{N}(\mu_2, \sigma_2^2)$
- les deux échantillons sont indépendants

Construction de la statistique de test : On dispose des estimateurs de σ_1^2 et σ_2^2 définis par

$$S_X^2 = \frac{1}{p-1} \sum_{i=1}^p (X_i - \bar{X}_p)^2 \quad \text{et} \quad S_Y^2 = \frac{1}{q-1} \sum_{i=1}^q (Y_i - \bar{Y}_q)^2.$$

La construction du test de Fisher repose sur la définition suivante.

Definition 5.4.1. Soient p, q deux entiers positifs et soit $Z_p \sim \chi_p^2$ et $Z_q \sim \chi_q^2$ deux variables aléatoires indépendantes. Alors

$$\frac{\frac{Z_p}{p}}{\frac{Z_q}{q}} \sim F(p, q)$$

où $F(p, q)$ suit une loi de Fisher à p, q degrés de liberté.

Corollaire 5.4.1. On a

$$\frac{\sigma_2^2 S_X^2}{\sigma_1^2 S_Y^2} \sim F(p-1, q-1).$$

Démonstration. Comme les variables aléatoires X_1, \dots, X_p et Y_1, \dots, Y_q sont indépendantes, les variables S_X^2 et S_Y^2 sont indépendantes. De plus

$$\frac{p-1}{\sigma_1^2} S_X^2 \sim \chi_{p-1}^2 \qquad \frac{q-1}{\sigma_2^2} S_Y^2 \sim \chi_{q-1}^2,$$

ce qui conclut la preuve. □

Comme sous H_0 , $\sigma_1^2 = \sigma_2^2$, on a

$$\frac{S_X^2}{S_Y^2} \sim F(p-1, q-1) \quad \text{sous } H_0,$$

ce qui est bien évidemment faux sous H_1 .

Le test :

- On teste au risque α l'hypothèse nulle $H_0 : \sigma_1^2 = \sigma_2^2$ contre $H_1 : \sigma_1^2 \neq \sigma_2^2$.
- On a la statistique de test

$$Z = \frac{S_X^2}{S_Y^2} \sim F(p-1, q-1) \quad \text{sous } H_0.$$

- On a la zone de rejet

$$ZR = \{Z < f_{p,q,\alpha/2}\} \cup \{Z > f_{p,q,1-\alpha/2}\},$$

où $f_{p,q,\alpha/2}$ et $f_{p,q,1-\alpha/2}$ sont les quantiles d'ordre $\alpha/2$ et $1-\alpha/2$ de la loi de Fisher à $p-1, q-1$ degrés de liberté.

- On calcule z_{obs} . Si z_{obs} appartient à la zone de rejet, on rejette H_0 et inversement.

Remarque 1 : A noter que pour tout $\alpha \in (0, 1)$ et pour tout entiers positifs p, q , on a

$$f_{p,q,\alpha} = \frac{1}{f_{p,q,1-\alpha}}.$$

Remarque 2 : On peut intervertir S_X^2 et S_Y^2 et on a alors la statistique de test

$$Z = \frac{S_Y^2}{S_X^2} \sim F(q-1, p-1) \quad \text{sous } H_0.$$

Remarque 3 : Soit F une variable aléatoire suivant une loi de Fisher de paramètres p, q , on a

$$p\text{-value} = \min \{ \mathbb{P} [F \leq z_{obs}], \mathbb{P} [F \geq z_{obs}] \}.$$

En effet, notons $F_{p,q}$ la fonction de répartition de la loi de Fisher qui est strictement croissante, on a

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), z_{obs} < f_{p,q,\alpha} \quad \text{ou} \quad z_{obs} > f_{p,q,1-\alpha} \} \\ &= \min \{ \inf \{ \alpha \in (0, 1), z_{obs} < f_{p,q,\alpha} \}, \inf \{ \alpha \in (0, 1), z_{obs} > f_{p,q,1-\alpha} \} \} \\ &= \min \{ \inf \{ \alpha \in (0, 1), F_{p,q}(z_{obs}) < \alpha \}, \inf \{ \alpha \in (0, 1), F_{p,q}(z_{obs}) > 1 - \alpha \} \} \\ &= \min \{ F_{p,q}(z_{obs}), 1 - F_{p,q}(z_{obs}) \} \end{aligned}$$

Remarque 4 : On peut rencontrer la statistique de test

$$Z = \frac{\max(S_X^2, S_Y^2)}{\min(S_X^2, S_Y^2)}.$$

Exemple : On reprend l'exemple des iris de Fisher. Rappelons que l'on a $s_X^2 = 0.52^2$ et $s_Y^2 = 0.64^2$. On teste au risque de 5% l'hypothèse nulle $H_0 : \sigma_1^2 = \sigma_2^2$ contre l'hypothèse alternative $H_1 : \sigma_1^2 \neq \sigma_2^2$. On a la statistique de test

$$Z = \frac{S_X^2}{S_Y^2} \sim F(49, 49) \quad \text{sous } H_0$$

On a la zone de rejet

$$Z = \{Z < 0.57\} \cup \{Z > 1.76\}.$$

Ici, $z_{obs} = 0.66$ et on ne rejette donc pas H_0 , i.e on ne rejette pas le fait que les variances puissent être égales.

5.4.2 Test de Shapiro-Wilk

Afin de continuer de vérifier le cadre probabiliste pour les tests de Fisher, il reste entre autre à vérifier que les variables aléatoires X et Y suivent des lois normales. On considère x_1, \dots, x_n qui sont des réalisations des variables aléatoires indépendantes et de même loi X_1, \dots, X_n . Le test de Shapiro-Wilk permet de tester l'hypothèse nulle H_0 : "La variable X est gaussienne" contre l'hypothèse alternative H_1 : "La variable X n'est pas gaussienne". On ne donnera pas ici les détails de ce test, on retiendra seulement la conclusion :

- si la p -value est supérieure à α , alors on ne peut pas rejeter H_0 , i.e on ne peut pas rejeter le fait que les données soient distribuées selon une loi normale.
- sinon, on rejette H_0 .

Exemple : Si on reprend l'exemple des iris, on obtient (pour les espèces Versicolor et Virginica) pour le test de Shapiro-Wilk les p -values suivantes : 0.46 et 0.26. On ne rejette donc pas H_0 , i.e on ne rejette pas le caractère gaussien des données observées.

5.5 Test de Student dans le cas apparié

5.5.1 Introduction

On considère des données $((x_1, y_1), \dots, (x_n, y_n))$ qui sont des réalisations de couples de variables aléatoires indépendants $(X_1, Y_1), \dots, (X_n, Y_n)$ de même loi que (X, Y) . On souhaite comparer les moyennes des variables aléatoires X et Y .

Le cadre probabiliste

- Les variables aléatoires X_i et Y_i ne sont pas (nécessairement) indépendantes.
- Les espérances de X et Y sont données par μ_1 et μ_2 .
- La variable aléatoire $X - Y$ suit une loi normale d'espérance $\mu = \mu_1 - \mu_2$ et de variance σ^2 .

5.5.2 Test d'égalité

On veut tester au risque α l'hypothèse nulle $H_0 : \mu_1 = \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 \neq \mu_2$.

Construction de la statistique de test : On dispose des estimateurs \bar{X}_n et \bar{Y}_n et on s'intéresse à la variable aléatoire $\bar{X}_n - \bar{Y}_n$.

Proposition 5.5.1. *On a*

$$\bar{X}_n - \bar{Y}_n \sim \mathcal{N}\left(\mu_1 - \mu_2, \frac{\sigma^2}{n}\right)$$

Démonstration. On a $\bar{X}_n - \bar{Y}_n = \frac{1}{n} \sum_{i=1}^n (X_i - Y_i)$ qui est une combinaison linéaire de variables aléatoires indépendantes suivant des lois normales et par linéarité

$$\mathbb{E}[\bar{X}_n - \bar{Y}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Y_i] = \mu_1 - \mu_2,$$

et par indépendance

$$\mathbb{V}[\bar{X}_n - \bar{Y}_n] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{V}[X_i - Y_i] = \frac{\sigma^2}{n}.$$

□

On obtient ainsi

$$\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - (\mu_1 - \mu_2)}{\sigma} \sim \mathcal{N}(0, 1).$$

Cependant, on ne connaît pas la variance σ^2 , que l'on va donc remplacer par un estimateur. La proposition suivante nous donne cet estimateur ainsi que sa loi.

Proposition 5.5.2. *On considère l'estimateur S^2 de σ^2 défini par*

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - Y_i - (\bar{X}_n - \bar{Y}_n))^2.$$

Cet estimateur vérifie alors

1. *On a*

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

2. *Les variables $\bar{X}_n - \bar{Y}_n$ et S^2 sont indépendantes.*

Démonstration. Il suffit juste de considérer la variable aléatoire $Z = X - Y$, $\mu = \mu_1 - \mu_2$. En remarquant que

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2,$$

on obtient

1. $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$
2. $\frac{(n-1)S^2}{\sigma^2}$ et \bar{Z}_n sont indépendants.

□

Corollaire 5.5.1. On a

$$\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - (\mu_1 - \mu_2)}{S} \sim T_{n-1}$$

où T_{n-1} suit une loi de Student à $n - 1$ degrés de liberté.

On obtient donc la statistique de test

$$\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \sim T_{n-1} \quad \text{sous } H_0,$$

ce qui est bien évidemment faux sous H_1 .

Le test :

- On teste au risque α l'hypothèse nulle $H_0 : \mu_1 = \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 \neq \mu_2$.
- On a la statistique de test

$$Z = \sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \sim T_{n-1} \quad \text{sous } H_0.$$

- On a la zone de rejet

$$ZR = \{Z, |Z| > t_{n-1, 1-\alpha/2}\}$$

où $t_{n-1, 1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi de Student à $n - 1$ degrés de liberté.

- Si $|z_{obs}| > t_{n-1, 1-\alpha/2}$, on rejette H_0 et inversement.

Remarque 1 : Cela revient à vérifier que 0 appartient à l'intervalle de confiance de niveau $1 - \alpha$ de $\mu = \mu_1 - \mu_2$. En effet

$$\begin{aligned} |z_{obs}| \leq t_{n-1, 1-\alpha/2} &\Leftrightarrow -t_{n-1, 1-\alpha/2} \leq \sqrt{n} \frac{\bar{x}_n - \bar{y}_n}{s} \leq t_{n-1, 1-\alpha/2} \\ &\Leftrightarrow \bar{x}_n - \bar{y}_n - t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}} \leq 0 \leq \bar{x}_n - \bar{y}_n + t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}} \\ &\Leftrightarrow \bar{z}_n - t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}} \leq 0 \leq \bar{z}_n + t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}} \\ &\Leftrightarrow 0 \in IC_{1-\alpha}(\mu). \end{aligned}$$

Remarque 2 : En réalité, pour trouver la zone de rejet, comme $\bar{X}_n - \bar{Y}_n$ doit être proche de 0, on cherche c_α tel que

$$\mathbb{P}_{\mu=0} [\text{"On rejette } H_0"] = \mathbb{P}_{\mu=0} [|\bar{X}_n - \bar{Y}_n| \geq c_\alpha] = \alpha.$$

On obtient donc, par symétrie de la loi de Student

$$\begin{aligned} \alpha &= \mathbb{P}_{\mu=0} [|\bar{X}_n - \bar{Y}_n| \geq c_\alpha] \\ &= \mathbb{P}_{\mu=0} \left[\sqrt{n} \frac{|\bar{X}_n - \bar{Y}_n|}{S} \geq \frac{\sqrt{n}}{S} c_\alpha \right] \\ &= \mathbb{P}_{\mu=0} \left[\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \geq \frac{\sqrt{n}}{S} c_\alpha \right] + \mathbb{P}_{\mu=0} \left[\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \leq -\frac{\sqrt{n}}{S} c_\alpha \right] \\ &= 2\mathbb{P}_{\mu=0} \left[\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \geq \frac{\sqrt{n}}{S} c_\alpha \right]. \end{aligned}$$

On obtient donc, en notant F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté,

$$\alpha = 2 - 2F_T \left(\frac{\sqrt{n}}{S} c_\alpha \right).$$

En appliquant F_T^{-1} , on obtient

$$t_{n-1, 1-\alpha/2} = \frac{\sqrt{n}}{S} c_\alpha$$

et donc

$$c_\alpha = t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}}.$$

On retrouve ainsi la zone de rejet

$$ZR = \left\{ |\bar{X}_n - \bar{Y}_n| \geq t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}} \right\} = \left\{ \sqrt{n} \frac{|\bar{X}_n - \bar{Y}_n|}{S} \geq t_{n-1, 1-\alpha/2} \right\}.$$

Remarque 3 : Soit $T \sim T_{n-1}$, on a alors

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0" \} \\ &= \inf \{ \alpha \in (0, 1), |z_{obs}| \geq t_{n-1, 1-\alpha/2} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(|z_{obs}|) \geq 1 - \alpha/2 \} \\ &= 2 - 2F_T(|z_{obs}|) \\ &= \mathbb{P}[|T| \geq |z_{obs}|]. \end{aligned}$$

Exemple : Cet exemple est tiré du cours de B. Portier à l'INSA de Rouen. On s'intéresse à un échantillon de 30 matières fécales. On soumet cet échantillon à deux méthodes différentes de spectrométrie pour étudier leur teneur en lutécium radioactif. Les mesures relatives aux deux méthodes ne

peuvent donc bien évidemment pas être indépendantes. On souhaite savoir si les résultats obtenus avec ces deux méthodes sont équivalents, i.e si les moyennes sont les mêmes. On obtient

$$\bar{x}_{30} = 120.83 \quad \text{et} \quad \bar{y}_{30} = 119.33.$$

En effectuant le test pour les données appariées, on obtient une p -value égale à 0.66 et au risque de 5%, on ne rejette donc pas H_0 . A noter que si on avait effectué le test de Student pour des variables indépendantes, on aurait obtenue un p -value égale à 0.0031 et on aurait alors rejeter H_0 . Cette différence est due à la forte corrélation entre les données. En effet le coefficient de corrélation entre les données est de 0.982.

5.5.3 Test d'inégalité $\mu_1 \leq \mu_2$

Le test : On teste au risque $\alpha \in (0, 1)$ l'hypothèse nulle $H_0 : \mu_1 \leq \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 > \mu_2$. Sous H_0 , il existe $\mu' \leq 0$ tel que

$$Z(\mu') = \sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - \mu'}{S} \sim T_{n-1}.$$

On définit alors la zone de rejet par

$$ZR = \{Z(0) > t_{n-1, 1-\alpha}\},$$

où $t_{n-1, 1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de Student à $n - 1$ degrés de liberté. On calcule $Z_{obs} = \sqrt{n} \frac{\bar{x}_n - \bar{y}_n}{s}$. Si $Z_{obs} > t_{n-1, 1-\alpha}$, on rejette H_0 et inversement.

Remarque 1 : En réalité, pour trouver la zone de rejet, comme $\mu = \mu_1 - \mu_2$ peut être aussi petit que l'on veut sous H_0 , on doit rejeter H_0 si $\bar{X}_n - \bar{Y}_n$ est trop grand, i.e on cherche c_α tel que

$$\sup_{\mu \leq 0} \mathbb{P}_\mu [\text{"On rejette } H_0\text{"}] = \sup_{\mu \leq 0} \mathbb{P}_\mu [\bar{X}_n - \bar{Y}_n \geq c_\alpha] = \alpha.$$

Or, si on note F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté, on obtient

$$\begin{aligned} \alpha &= \sup_{\mu \leq 0} \mathbb{P}_\mu [\bar{X}_n - \bar{Y}_n \geq c_\alpha] \\ &= \sup_{\mu \leq 0} \mathbb{P}_\mu \left[\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - \mu}{S} \geq \sqrt{n} \frac{c_\alpha - \mu}{S} \right] \\ &= \sup_{\mu \leq 0} \left(1 - F_T \left(\sqrt{n} \frac{c_\alpha - \mu}{S} \right) \right). \end{aligned}$$

Comme la fonction de répartition F_T est strictement croissante, on obtient

$$\alpha = 1 - F_T \left(\frac{c_\alpha}{\sqrt{n}S} \right),$$

et donc en appliquant F_T^{-1} ,

$$t_{n-1,1-\alpha} = \frac{c_\alpha}{\sqrt{n}S},$$

et donc

$$c_\alpha = t_{n-1,1-\alpha} \frac{S}{\sqrt{n}}.$$

On retrouve ainsi la zone de rejet, i.e

$$ZR = \left\{ \bar{X}_n - \bar{Y}_n \geq t_{n-1,1-\alpha} \frac{S}{\sqrt{n}} \right\} = \left\{ \sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \geq t_{n-1,1-\alpha} \right\}.$$

Remarque 2 : On a

$$p\text{-value} = \mathbb{P}[T \geq z_{obs}],$$

où $T \sim T_{n-1}$. En effet;

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \geq t_{n-1,1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \geq 1 - \alpha \} \\ &= 1 - F_T(z_{obs}). \end{aligned}$$

5.5.4 Test d'inégalité $\mu_1 \geq \mu_2$

On veut tester au risque $\alpha \in (0, 1)$ l'hypothèse nulle $H_0 : \mu_1 \geq \mu_2$ contre l'hypothèse alternative $H_1 : \mu_1 < \mu_2$. Sous H_0 , il existe $\mu' \geq 0$ tel que

$$Z(\mu') = \sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - \mu'}{S} \sim T_{n-1}.$$

On définit alors la zone de rejet par

$$ZR = \{ Z(0) < -t_{n-1,1-\alpha} \},$$

où $t_{n-1,1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de Student à $n - 1$ degrés de liberté. On calcul $z_{obs} = \sqrt{n} \frac{\bar{x}_n - \bar{y}_n}{s}$. Si $z_{obs} < -t_{n-1,1-\alpha}$, on rejette H_0 et inversement.

Remarque 1 : En réalité, pour trouver la zone de rejet, comme $\mu = \mu_1 - \mu_2$ peut être aussi grand que l'on veut sous H_0 , on doit rejeter H_0 si $\bar{X}_n - \bar{Y}_n$ est trop petit, i.e on cherche c_α tel que

$$\sup_{\mu \geq 0} \mathbb{P}_\mu [\text{"On rejette } H_0"] = \sup_{\mu \geq 0} \mathbb{P}_\mu [\bar{X}_n - \bar{Y}_n \leq c_\alpha] = \alpha.$$

Or si on note F_T la fonction de répartition de la loi de Student à $n - 1$ degrés de liberté, on obtient

$$\begin{aligned}\alpha &= \sup_{\mu \geq 0} \mathbb{P}_\mu [\bar{X}_n - \bar{Y}_n \leq c_\alpha] \\ &= \sup_{\mu \geq 0} \mathbb{P}_\mu \left[\sqrt{n} \frac{\bar{X}_n - \bar{Y}_n - \mu}{S} \leq \sqrt{n} \frac{c_\alpha - \mu}{S} \right] \\ &= \sup_{\mu \geq 0} F_T \left(\sqrt{n} \frac{c_\alpha - \mu}{S} \right)\end{aligned}$$

Comme la fonction F_T est strictement croissante, on obtient

$$\alpha = F_T \left(\sqrt{n} \frac{c_\alpha}{S} \right),$$

et donc, en appliquant F_T^{-1} ,

$$t_{n-1, \alpha} = \sqrt{n} \frac{c_\alpha}{S}.$$

Ainsi, comme $t_{n-1, \alpha} = -t_{n-1, 1-\alpha}$,

$$c_\alpha = -t_{n-1, 1-\alpha} \frac{S}{\sqrt{n}}.$$

On retrouve ainsi la zone de rejet, i.e

$$ZR = \left\{ \bar{X}_n - \bar{Y}_n \leq -t_{n-1, 1-\alpha} \frac{S}{\sqrt{n}} \right\} = \left\{ \sqrt{n} \frac{\bar{X}_n - \bar{Y}_n}{S} \leq -t_{n-1, 1-\alpha} \right\}.$$

Remarque 2 : On a

$$p\text{-value} = \mathbb{P}[T \leq z_{obs}],$$

où $T \sim T_{n-1}$. En effet,

$$\begin{aligned}p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \leq -t_{n-1, 1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F_T(z_{obs}) \leq \alpha \} \\ &= F_T(z_{obs}).\end{aligned}$$

5.6 Tests du Khi-deux

5.6.1 Test d'indépendance du Khi-deux

Le test d'indépendance du Khi-deux permet de vérifier l'indépendance entre deux variables aléatoires X, Y , i.e on veut tester au risque α l'hypothèse nulle H_0 : "X et Y sont indépendantes" contre l'hypothèse alternative H_1 : "X et Y ne sont pas indépendantes".

Le cadre : On dispose de n données $(x_1, y_1), \dots, (x_n, y_n)$ qui sont des réalisations d'un couple de

variables aléatoires (X, Y) à valeurs dans $\{a_1, \dots, a_p\} \times \{b_1, \dots, b_q\}$.

Construction du test : Le test du Khi-deux repose sur le fait que les variables aléatoires X et Y sont indépendantes si et seulement si pour tout $(i, j) \in \{1, \dots, p\} \times \{1, \dots, q\}$,

$$\mathbb{P}[X = a_i \text{ et } Y = b_j] = \mathbb{P}[X = a_i] \mathbb{P}[Y = b_j].$$

Pour tester si les variables aléatoires X et Y sont indépendantes, il suffit donc de comparer les estimations des probabilités jointes au produit des estimations des probabilités marginales.

Notations : Les réalisations $(x_1, y_1), \dots, (x_n, y_n)$ sont des réalisations des couples de variables aléatoires (X_k, Y_k) , où les (X_k, Y_k) sont indépendants et de même loi que (X, Y) . Pour tout $(i, j) \in \{1, \dots, p\} \times \{1, \dots, q\}$, on note :

- $O_{i,j} = \sum_{k=1}^n \mathbf{1}_{\{X_k=a_i, Y_k=b_j\}}$ qui désigne le nombre de données pour lesquelles $(X, Y) = (a_i, b_j)$.
- $O_{i.} = \sum_{j=1}^q O_{i,j}$ qui désigne le nombre de données pour lesquelles $X = a_i$.
- $O_{.j} = \sum_{i=1}^p O_{i,j}$ qui désigne le nombre de données pour lesquelles $Y = b_j$.
- $E_{i,j} = \frac{O_{i.} \times O_{.j}}{n}$.

Notons que

- $O_{i,j}/n$ est un estimateur de $\mathbb{P}[X = a_i, Y = b_j]$.
- $E_{i,j}/n$ est un estimateur de $\mathbb{P}[X = a_i] \mathbb{P}[Y = b_j]$.

Sous l'hypothèse H_0 , ces deux estimateurs sont censés être proches. On introduit la variable

$$Z = \sum_{i=1}^p \sum_{j=1}^q \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}.$$

Proposition 5.6.1 (Admise). *Si X et Y sont indépendantes, alors,*

$$Z \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_{(p-1)(q-1)}^2.$$

La proposition précédente est fautive sous H_1 .

Le test :

- On teste au risque α l'hypothèse nulle H_0 : "X et Y sont indépendantes" contre l'hypothèse alternative H_1 : "X et Y ne sont pas indépendantes".
- On a la statistique de test

$$Z = \sum_{i=1}^p \sum_{j=1}^q \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_{(p-1)(q-1)}^2 \quad \text{sous } H_0.$$

- On a la zone de rejet

$$ZR = \left\{ Z > k_{(p-1)(q-1), 1-\alpha} \right\}$$

où $k_{(p-1)(q-1), 1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi du Khi-deux à $(p - 1)(q - 1)$ degrés de liberté.

— Si z_{obs} appartient à la zone de rejet, alors on rejette H_0 et inversement.

Exemple : Cet exemple est issu du cours de B. Portier à l'INSA de Rouen. On étudie l'influence du sexe sur la couleur des cheveux d'un groupe d'élèves. On veut savoir si la couleur des cheveux est indépendante du sexe. Pour cela, on dispose des données suivantes :

Sexe	Blond	Roux	Châtain	Brun	Noir de Jais	Total
Garçon	592	119	849	504	36	2100
Fille	544	97	677	451	14	1783
Total	1136	2126	1526	955	50	3883

On désigne par X la variable couleur des cheveux et par Y la variable sexe.

On teste au risque de 5% l'hypothèse H_0 : "X et Y sont indépendantes" contre H_1 : "X et Y ne sont pas indépendantes". On a la statistique de test

$$Z = \sum_{i=1}^5 \sum_{j=1}^2 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_4^2 \quad \text{sous } H_0.$$

On a la zone de rejet

$$ZR = \{Z > 9.50\}.$$

Ici, $z_{obs} = 10.47$ et on rejette donc l'hypothèse H_0 . A noter que l'on a une p -value égale à 0.033 et on n'aurait donc pas rejeté H_0 au risque de 1%.

5.6.2 Test d'adéquation

On dispose de n données x_1, \dots, x_n qui sont les réalisations de variables aléatoires indépendantes X_1, \dots, X_n et de même loi, à valeurs dans $\{a_1, \dots, a_K\}$. On note P la loi inconnue de la variable aléatoire discrète X et on souhaite savoir si $P = P_0$ où P_0 est une loi connue. On veut donc tester au risque α l'hypothèse nulle H_0 : "la loi de X est P_0 " i.e " $P = P_0$ " contre l'hypothèse alternative H_1 : "la loi de X n'est pas P_0 ", i.e " $P \neq P_0$ ".

Notations : On note $p_{0,1}, \dots, p_{0,K} > 0$ les probabilités définissant la loi P_0 . On note également pour tout $k = 1, \dots, K$:

— E_k l'effectif observé pour la modalité a_k , i.e

$$E_k = \sum_{i=1}^n \mathbf{1}_{\{X_i = a_k\}}.$$

— N_k l'effectif théorique pour la modalité a_k sous la loi P_0 , i.e

$$N_k = np_{0,k}.$$

Construction de la statistique de test : La variable E_k/n est donc un estimateur naturel de $p_k = \mathbb{P}[X = a_k]$. Sous l'hypothèse nulle, cet estimateur doit donc converger vers $N_k/n = p_{0,k}$. On cherche alors une statistique de test et on s'intéresse pour cela à

$$Q^2 = \sum_{k=1}^K \frac{(E_k - N_k)^2}{N_k}.$$

Proposition 5.6.2. *Si la loi de X est PO, alors*

$$Q^2 \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_{K-1}^2$$

Démonstration. A noter que le TLC pour chaque E_k nous donne

$$\sqrt{n} \left(\frac{E_k}{n} - p_k \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, p_k(1 - p_k))$$

mais il est inexploitable pour obtenir le résultat. On va plutôt utiliser le théorème de Cochran. Pour cela, on remarque que l'on peut écrire

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} Y_{1,i} \\ \vdots \\ Y_{K,i} \end{pmatrix} = \frac{1}{n} \begin{pmatrix} E_1 \\ \vdots \\ E_K \end{pmatrix}$$

avec $Y_{k,i} = \mathbf{1}_{X_i=k}$. Les vecteurs Y_i sont i.i.d avec

$$\mathbb{E}[Y_1] = \begin{pmatrix} p_1 \\ \vdots \\ p_K \end{pmatrix}$$

et comme si $k \neq k'$, on a $Y_{i,k}Y_{i,k'} = 0$, on obtient

$$\text{Var}[Y_1] = \begin{pmatrix} p_1(1-p_1) & -p_1p_2 & \cdots \\ -p_1p_2 & p_2(1-p_2) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

Le TLC multidimensionnel nous donne donc

$$\sqrt{n} (\bar{Y}_n - \mathbb{E}[Y_1]) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \text{Var}[Y_1])$$

et comme $\text{Var}[Y_1] = \text{diag}(\sqrt{p_1}, \dots, \sqrt{p_K}) \left(I - \sqrt{p} \sqrt{p}^T \right) \text{diag}(\sqrt{p_1}, \dots, \sqrt{p_K})$, on peut donc ré-

écrire le TLC comme

$$\sqrt{n} \begin{pmatrix} \frac{\bar{Y}_{1,n-p_1}}{\sqrt{p_1}} \\ \vdots \\ \frac{\bar{Y}_{K,n-p_1}}{\sqrt{p_K}} \end{pmatrix} = \sqrt{n} \begin{pmatrix} \frac{E_1/n-p_1}{\sqrt{p_1}} \\ \vdots \\ \frac{E_K/n-p_1}{\sqrt{p_K}} \end{pmatrix} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} \left(0, I - \sqrt{p} \sqrt{p}^T \right)$$

Or $\sqrt{p} \sqrt{p}^T$ est la matrice de projection orthogonale sur $V = \text{Vect}\{p\}$, et $I - \sqrt{p} \sqrt{p}^T$ est la projection orthogonale sur V^\perp . Et on a donc, en appliquant le théorème de continuité et le théorème de Cochran

$$\sum_{k=1}^K \frac{(E_k - N_k)^2}{N_k} = n \left\| \begin{pmatrix} \frac{E_1/n-p_1}{\sqrt{p_1}} \\ \vdots \\ \frac{E_K/n-p_1}{\sqrt{p_K}} \end{pmatrix} \right\|^2 \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \left\| \mathcal{N} \left(0, I - \sqrt{p} \sqrt{p}^T \right) \right\|^2 \sim \chi_{K-1}^2.$$

□

La proposition précédente est fautive si l'hypothèse H_0 n'est pas vérifiée.

Le test

- On test au risque α l'hypothèse nulle H_0 : "X suit la loi P_0 " contre H_1 : "X ne suit pas la loi P_0 ".
- On a la statistique de test

$$Q^2 = \sum_{k=1}^K \frac{(E_k - N_k)^2}{N_k} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_{K-1}^2 \quad \text{sous } H_0$$

- On a la zone de rejet

$$ZR = \{Q^2 > k_{1-\alpha}\}$$

où $k_{1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi du Khi-deux à $K - 1$ degrés de liberté.

- Si q_{obs}^2 appartient à la zone de rejet, on rejette H_0 et inversement.

Exemple : Les données suivantes sont issues du cours de B. Portier à l'INSA de Rouen. Un biologiste pense qu'à un âge donné :

- 50% des bébés marchent
- 12% ont une ébauche de marche
- 38% ne marchent pas

On fait une étude sur 80 bébés et on obtient :

- 35 de ces bébés ne marchent pas
- 4 ont une ébauche de marche
- 41 ne marchent pas

On souhaite donc vérifier que le biologiste a tort. On résume les données dans le tableau suivant :

	Marche	Ebauche	Ne marche pas	Total
Effectif observé	35	4	41	80
Effectif Théorique	80×0.5	80×0.12	80×0.38	80
"Distance"	0.625	3.267	3.696	7.588

On teste donc au risque 5% l'hypothèse nulle H_0 : "X suit la loi P_0 " contre H_1 : "X ne suit pas la loi P_0 ". On a la statistique de test

$$Q^2 = \sum_{k=1}^3 \frac{(E_k - N_k)^2}{N_k} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_2^2 \quad \text{sous } H_0$$

On a la zone de rejet

$$ZR = \{Q^2 > 5.99\}.$$

Ici, $q_{obs}^2 = 7.5878$ et on rejette donc H_0 .

5.7 Tests asymptotiques

5.7.1 Introduction

On considère ici des données x_1, \dots, x_n qui sont des réalisations de variables aléatoires X_1, \dots, X_n qui sont indépendantes et de même loi que X . On s'intéresse à l'estimation d'un paramètre θ de la loi de X , et on suppose que l'on a un estimateur $\hat{\theta}_n$ asymptotiquement normal, i.e qu'il existe $\sigma^2 > 0$ tel que

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2).$$

On suppose également que l'on dispose d'un estimateur consistant $\hat{\sigma}_n^2$ de la variance asymptotique σ^2 . A l'aide du théorème de Slutsky, on a alors

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

5.7.2 Test d'égalité

On veut tester au risque α l'hypothèse nulle H_0 : " $\theta = \theta_0$ " contre l'hypothèse alternative H_1 : " $\theta \neq \theta_0$ ". On a la statistique de test

$$Z = \sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, 1) \quad \text{sous } H_0$$

On obtient alors la zone de rejet

$$ZR = \{|z_{obs}| > q_{1-\alpha/2}\},$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite. Si z_{obs} appartient à la zone de rejet, on rejette H_0 et inversement.

Remarque 1 : Cela revient à vérifier que θ_0 appartient à l'intervalle de confiance asymptotique de niveau $1 - \alpha$ de θ . En effet,

$$\begin{aligned} |z_{obs}| \leq q_{1-\alpha/2} &\Leftrightarrow -q_{1-\alpha/2} \leq \sqrt{n} \frac{\theta_n - \theta_0}{\sigma_n} \leq q_{1-\alpha/2} \\ &\Leftrightarrow \theta_n - q_{1-\alpha/2} \frac{\sigma_n}{\sqrt{n}} \leq \theta_0 \leq \theta_n + q_{1-\alpha/2} \frac{\sigma_n}{\sqrt{n}} \\ &\Leftrightarrow \theta_0 \in IC_{1-\alpha}(\theta) \end{aligned}$$

Remarque 2 : En réalité, pour trouver la zone de rejet, comme $\hat{\theta}_n$ doit être proche de θ , on cherche c_α tel que

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\theta=\theta_0} [\text{"On rejette } H_0"] = \lim_{n \rightarrow \infty} \mathbb{P} [|\hat{\theta}_n - \theta_0| \geq c_\alpha] = \alpha.$$

On obtient donc

$$\alpha = \lim_{n \rightarrow \infty} \mathbb{P}_{\theta=\theta_0} [|\hat{\theta}_n - \theta_0| \geq c_\alpha] = \lim_{n \rightarrow \infty} 2\mathbb{P}_{\theta=\theta_0} \left[\sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n} \geq \frac{\sqrt{n}}{\hat{\sigma}_n} c_\alpha \right].$$

Comme

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

on choisit donc

$$\frac{\sqrt{n}}{\hat{\sigma}_n} c_\alpha = q_{1-\alpha/2} \Leftrightarrow c_\alpha = q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}}.$$

On retrouve alors la zone de rejet

$$ZR = \left\{ |\hat{\theta}_n - \theta_0| \geq q_{1-\alpha/2} \frac{\hat{\sigma}_n}{\sqrt{n}} \right\} = \left\{ \sqrt{n} \frac{|\hat{\theta}_n - \theta_0|}{\hat{\sigma}_n} \geq q_{1-\alpha/2} \right\}.$$

Remarque 3 : Soit $Z \sim \mathcal{N}(0, 1)$, on a alors

$$p\text{-value} = \mathbb{P} [|Z| \geq |z_{obs}|].$$

En effet,

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0" \} \\ &= \inf \{ \alpha \in (0, 1), |z_{obs}| \geq q_{1-\alpha/2} \} \\ &= \inf \{ \alpha \in (0, 1), F(|z_{obs}|) \geq 1 - \alpha/2 \} \\ &= 2 - 2F(|z_{obs}|) \\ &= 2\mathbb{P} [Z \geq |z_{obs}|]. \end{aligned}$$

Par symétrie de la loi normale,

$$p\text{-value} = 2\mathbb{P}[Z \geq |z_{obs}|] = \mathbb{P}[Z \geq |z_{obs}|] + \mathbb{P}[-Z \geq |z_{obs}|] = \mathbb{P}[|Z| \geq |z_{obs}|].$$

5.7.3 Test d'inégalité $\theta_0 \geq \theta$

On veut tester au risque α l'hypothèse nulle $H_0 : \theta_0 \geq \theta$ contre l'hypothèse alternative $H_1 : \theta_0 < \theta$. Sous H_0 , il existe $\theta' \leq \theta_0$ tel que

$$Z(\theta') = \sqrt{n} \frac{\hat{\theta}_n - \theta'}{\hat{\sigma}_n} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

On définit alors la zone de rejet

$$ZR = \{Z(\theta_0) > q_{1-\alpha}\}$$

où $q_{1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi normale centrée réduite. On calcule $z_{obs} = \sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n}$, et si $z_{obs} > q_{1-\alpha}$ on rejette H_0 et inversement.

Remarque : On a

$$p\text{-value} = \mathbb{P}[Z \geq z_{obs}],$$

où $Z \sim \mathcal{N}(0, 1)$. En effet,

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \text{"} \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \geq q_{1-\alpha} \} \\ &= \inf \{ \alpha \in (0, 1), F(z_{obs}) \geq 1 - \alpha \} \\ &= 1 - F(z_{obs}). \end{aligned}$$

5.7.4 Test d'inégalité $\theta \geq \theta_0$

On veut tester au risque $\alpha \in (0, 1)$, l'hypothèse nulle $H_0 : \theta \geq \theta_0$ contre l'hypothèse alternative $H_1 : \theta < \theta_0$. Sous H_0 , il existe $\theta' \geq \theta_0$ tel que

$$Z(\theta') = \sqrt{n} \frac{\hat{\theta}_n - \theta'}{\hat{\sigma}_n} \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1).$$

On définit alors la zone de rejet par

$$ZR = \{Z(\theta_0) < -q_{1-\alpha}\},$$

où q_α est le quantile d'ordre α de la loi normale centrée réduite. Si $z_{obs} < -q_{1-\alpha}$, alors on rejette H_0 et inversement.

Remarque : On a

$$p\text{-value} = \mathbb{P}[Z \leq z_{obs}],$$

où Z suit une loi normale centrée réduite. En effet,

$$\begin{aligned} p\text{-value} &= \inf \{ \alpha \in (0, 1), \text{"On rejette } H_0 \} \\ &= \inf \{ \alpha \in (0, 1), z_{obs} \leq q_\alpha \} \\ &= \inf \{ \alpha \in (0, 1), F(z_{obs}) \leq \alpha \} \\ &= F(z_{obs}). \end{aligned}$$

