

L3 MAPI3
Simulations stochastiques
Partie Statistiques

L'objectif de cette seconde partie du cours est multiple :

- Renforcer le cours de statistique du deuxième semestre par une étude approfondie de quelques questions statistiques.
- L'accent sera mis sur l'utilisation des théorèmes limites (pour la théorie) et sur l'illustration par simulation (pour les TP).
- Nous présenterons aussi une introduction aux tests non-paramétriques, plus délicats du point de vue de la théorie.

Si vous souhaitez approfondir les notions abordées ici, notamment pour un niveau Master, nous vous recommandons la lecture de l'ouvrage suivant.

- Vincent Rivoirard et Gilles Stoltz, *Statistique en action*, Vuibert, seconde édition, 2012.

- 1) Premiers pas en statistiques : le problème du sondage
- 2) Test de Kolmogorov Smirnov

Chapitre I

Premiers pas en statistiques : le problème du sondage

- **Imaginons une situation politique en 2022** : Didier Deschamps se présente aux élections présidentielles en Belgique, et nous sommes le jour du second tour, en fin d'après-midi. Peu de bulletins ont été dépouillés, mais l'institut Iflop a réalisé un sondage à la sortie d'un stade afin d'estimer la proportion des électeurs qui ont voté pour Didier Deschamps.
- Après avoir interrogé 1225 personnes, l'institut de sondage a compté 637 électeurs ayant voté pour Didier Deschamps. Ainsi, 52% des personnes interrogées ont voté pour Didier Deschamps.
- Vous êtes journaliste sur une chaîne de télévision, et vous devez parler des élections. Quel pronostic annoncerez-vous à l'antenne ?
- Notons :
 - $N \approx 10^7$ le nombre total d'électeurs
 - M le nombre total de ceux qui ont voté pour Didier Deschamps parmi les N électeurs.
 - n le nombre de personnes interrogées à la sortie du stade

On souhaite savoir si la proportion $\theta = M/N$ est plus grande que $1/2$.

- **Mais comment se faire une idée de la valeur de θ (qui est inconnue) ?** Une approche très naturelle consiste à dire que θ doit être "proche" de la proportion $\hat{\theta}_n$ observée sur l'échantillon des n personnes interrogées.

$$\hat{\theta}_n \simeq \theta$$

- Plus formellement, en posant $X_i = 1$ si la i -ème personne interrogée a voté pour Didier Deschamps et $X_i = 0$ si elle a voté pour l'autre candidat, la v.a. $\hat{\theta}_n$ est définie par

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

- On dit dans ce cas que la proportion empirique $\hat{\theta}_n$ est un *estimateur* de la *proportion inconnue* θ . Le mot *estimateur* signifie que $\hat{\theta}_n$ est construit à partir de l'échantillon (X_1, \dots, X_n) . *En particulier, cette quantité est une variable aléatoire.* C'est une fonction de l'échantillon
- **Question** : comment mesurer la proximité de $\hat{\theta}_n$ et de θ ? Ces deux valeurs sont-elles d'ailleurs *toujours* proches?
- Non, pas si on choisit un mauvais échantillon de n personnes (par exemple ceux qui se rappellent du 10 juillet 2018). Mais si cet échantillon est choisi "au hasard", alors on pressent que $\hat{\theta}_n$ sera "souvent" proche de θ . Mais que signifie *souvent* précisément? *Une modélisation mathématique est nécessaire.*
- Dans cette partie, nous considérons la modélisation la plus simple : *le modèle de Bernoulli*. On suppose ainsi que parmi la population de N personnes, **on interroge n personnes "complètement au hasard" et "avec remise"**.
Population totale : N , sous population interrogée : n .
- Plus formellement : en numérotant chaque individu de la population par $k \in \{1, \dots, N\}$, le choix d'un échantillon correspond à tirer aléatoirement $K_1, \dots, K_n \stackrel{i.i.d.}{\sim} \mathcal{U}(\{1, \dots, N\})$.
- Pour chaque individu $k \in \{1, \dots, N\}$ de la population, on pose $x_k = 1$ s'il a voté pour DD, et $x_k = 0$ sinon. A l'issue du sondage, on n'observe pas tous les x_k mais seulement les valeurs $x_{K_1}, x_{K_2}, \dots, x_{K_n}$ des n personnes interrogées. On pose $X_i = x_{K_i}$ pour tout $i \in \{1, \dots, n\}$. On pose également $M = |\{k \in \{1, \dots, N\} : x_k = 1\}| = \sum_{k=1}^N x_k$: *le nombre total d'électeurs de DD*.
- Dans ce cadre, la proportion d'électeurs dans toute la population qui ont voté pour DD est donnée par $\theta = M/N$, et la proportion d'électeurs dans l'échantillon qui ont voté pour DD est égale à $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

Proposition 1. — Dans ce modèle de sondage avec remise (dit modèle de Bernoulli), la suite de variables aléatoires X_1, \dots, X_n est i.i.d., de loi $\mathcal{B}(\theta)$.
 — Conséquence importante : la variable aléatoire $\sum_{i=1}^n X_i$ suit la loi binomiale $\mathcal{B}(n, \theta)$.

Attention : on remet à chaque fois la personne choisie dans la population totale sinon le paramètre θ changerait à chaque fois.

- En pratique, on n'interroge jamais deux fois le même individu.
- On peut prendre en compte cette remarque en supposant que le tirage des n personnes dans $\{1, \dots, N\}$ est maintenant "**sans remise**".
- Plus précisément, le sondage consiste à choisir aléatoirement un sous-ensemble $\widehat{S} \subset \{1, \dots, N\}$ uniformément parmi tous les sous-ensembles de $\{1, \dots, N\}$ de cardinal n , au sens suivant : pour tout sous-ensemble $S \subset \{1, \dots, N\}$ de cardinal n ,

$$\mathbb{P}_\theta(\widehat{S} = S) = \frac{1}{\binom{N}{n}} .$$

- On estime alors θ par la proportion $\widehat{\theta}_n = \frac{1}{n} \sum_{i \in \widehat{S}} X_i$ d'individus de l'échantillon \widehat{S} qui ont voté pour DD.

Proposition 2. Considérons le modèle de sondage sans remise ci-dessus, et supposons pour simplifier que $n \leq \min\{M, N - M\}$ (hypothèse très raisonnable en pratique).

Alors, la variable aléatoire $\sum_{i \in \widehat{S}} x_i$ suit *une loi hypergéométrique* : pour tout $k \in \{0, 1, \dots, n\}$,

$$\mathbb{P}\left(\sum_{i \in \widehat{S}} X_i = k\right) = \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}} .$$

Or, on peut montrer que si M et $N - M$ sont grands et $M/N \approx \theta$, alors la loi hypergéométrique ci-dessus est proche de la loi binomiale $\mathcal{B}(n, \theta)$, au sens où :

Proposition 3 (Convergence de la loi hypergéométrique vers la loi binomiale).

Soit $0 \leq k \leq n$ deux entiers fixés. Alors,

$$\frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}} \xrightarrow[\substack{M, N-M \rightarrow +\infty \\ M/N \rightarrow \theta}]{\quad} \binom{n}{k} \theta^k (1 - \theta)^{n-k} .$$

Concrètement, cette proposition signifie que lorsque M et $N - M$ sont grands, un tirage sans remise est très proche d'un tirage avec remise puisqu'on retrouve la loi binomiale. La similarité entre tirage avec ou sans remise n'est pas surprenante car si M et $N - M$ sont grands, on s'attend, même avec remise, à tirer avec grande probabilité un échantillon de n personnes qui sont deux à deux distinctes.

- **Dans toute la suite du chapitre**, on suppose qu'on est dans le *modèle de Bernoulli*, c'est-à-dire que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$.
- En particulier, la somme $\sum_{i=1}^n X_i$ suit la loi binomiale $\mathcal{B}(n, \theta)$, ce qui nous permettra de décrire plus facilement la variable aléatoire

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

- Le modèle de Bernoulli dépend d'un paramètre θ ; on notera $\mathbb{P}_\theta()$, $\mathbf{E}_\theta()$ et $\text{Var}_\theta()$ pour spécifier que ces quantités sont calculées avec le paramètre θ .

Une fois posé le modèle mathématique, le statisticien a pour objectif de retrouver le(s) paramètre(s) du modèle, ou de répondre à des questions qui s'y rattachent.

Dans la suite du chapitre, nous décrirons 3 tâches qu'un statisticien peut être amené à réaliser :

1. Estimation (de la proportion θ).
2. Construction d'un intervalle de confiance (pour la proportion θ).
3. Tests d'hypothèses (du type : est-ce que $\theta \leq 1/2$ ou $\theta > 1/2$?).

L'estimateur

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

est très naturel puisqu'il correspond à la proportion d'individus de l'échantillon ayant voté pour DD.

En fait, $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ peut aussi être interprété comme un cas particulier de deux méthodes plus générales :

Proposition 4. *Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, l'estimateur $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ coïncide avec :*

- l'estimateur obtenu par la méthode des moments ;
- l'estimateur du maximum de vraisemblance.

- Pour se faire une première idée de la qualité de l'estimation de θ par $\hat{\theta}_n$, on peut déjà remarquer que, en moyenne, l'estimateur $\hat{\theta}_n$ est correct :

$$\mathbf{E}_\theta(\hat{\theta}_n) = \mathbf{E}_\theta\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}_\theta(X_i) = \theta .$$

- On dit dans ce cas que **l'estimateur est sans biais**
- Ce comportement en moyenne n'est pas forcément rassurant : de grosses erreurs positives pourraient compenser de grosses erreurs négatives.
- En fait, **la loi forte des grands nombres** montre que, puisque les v.a. X_i sont i.i.d., intégrables et d'espérance θ , l'estimateur $\hat{\theta}_n$ est asymptotiquement correct avec probabilité 1 :

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \theta .$$

- Cela signifie concrètement qu'il existe un événement A tel que $\mathbb{P}(A) = 1$ et :

$$\forall \omega \in A, \forall \varepsilon > 0, \exists n_0(\omega, \varepsilon), \forall n \geq n_0(\omega, \varepsilon), |\hat{\theta}_n(\omega) - \theta| \leq \varepsilon .$$

- C'est déjà mieux qu'un simple contrôle de l'espérance !
- Subsiste quand même un problème : le nombre $n_0(\omega, \varepsilon)$ d'observations nécessaires pour approcher θ avec une précision ε dépend du résultat $\omega \in \Omega$ de l'expérience aléatoire.
- En pratique, on voudrait savoir ce que vaut n_0 (pour choisir n afin d'assurer une précision ε), mais on ne connaît pas ω ...
- Une façon d'éviter la dépendance en ω est d'utiliser **la loi faible des grands nombres** (convergence en probabilité) qui implique que : pour tout $\varepsilon > 0$,

$$\mathbb{P}_\theta(|\hat{\theta}_n - \theta| > \varepsilon) \xrightarrow[n \rightarrow +\infty]{} 0 .$$

- Par conséquent :

$$\forall \varepsilon > 0, \forall \alpha \in (0, 1), \exists n_0(\varepsilon, \alpha), \forall n \geq n_0(\varepsilon, \alpha), \mathbb{P}_\theta(|\hat{\theta}_n - \theta| > \varepsilon) \leq \alpha .$$

8 CHAPITRE I. PREMIERS PAS EN STATISTIQUES : LE PROBLÈME DU SONDAGE

- Ce résultat est préférable car il indique qu'on peut contrôler la probabilité que $\widehat{\theta}_n$ soit distant de θ de plus de ε dès que le nombre d'observations n est choisi assez grand (sans dépendance en ω).
- Malheureusement, ce résultat est encore très qualitatif car la dépendance de $n_0(\varepsilon, \alpha)$ en ε et α n'est pas quantifiée. Dans la suite, on explique comment sont reliées les quantités n , ε et α .

Remarquons déjà que la variance de $\widehat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ décroît quand le nombre n d'observations augmente :

Proposition 5 (Variance de la moyenne empirique $\widehat{\theta}_n$).

Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, l'estimateur $\widehat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ vérifie :

$$\text{Var}_\theta(\widehat{\theta}_n) = \frac{\theta(1-\theta)}{n} .$$

La proposition suivante permet de montrer qu'avec grande probabilité, la distance entre $\widehat{\theta}_n$ et θ est au plus de l'ordre de $1/\sqrt{n}$:

Proposition 6 (Contrôle avec grande probabilité de la qualité de $\widehat{\theta}_n$).

Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, l'estimateur $\widehat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ vérifie : pour tout $\varepsilon > 0$,

$$\mathbb{P}_\theta(|\widehat{\theta}_n - \theta| \leq \varepsilon) \geq 1 - \frac{1}{4n\varepsilon^2} . \tag{I.1}$$

En particulier, si on fixe $\alpha \in (0, 1)$ et qu'on choisit $\varepsilon = 1/\sqrt{4n\alpha}$, on a :

$$\mathbb{P}_\theta\left(|\widehat{\theta}_n - \theta| \leq \frac{1}{\sqrt{4n\alpha}}\right) \geq 1 - \alpha . \tag{I.2}$$

Preuve :

Le résultat découle de l'inégalité de Bienaymé-Tchebychev.

Rappels :

Inégalité de Markov : Si X est une v.a. positive ou nulle, alors : $\forall x > 0$, $\mathbb{P}(X \geq x) \leq \mathbf{E}(X)/x$.

Inégalité de Bienaymé-Tchebychev : Si X est une v.a. de carré intégrable, alors : $\forall x > 0$, $\mathbb{P}(|X - \mathbf{E}(X)| \geq x) \leq \text{Var}(X)/x^2$.

On applique l'inégalité de Bienaymé-Tchebychev à la v.a. $\widehat{\theta}_n$ qui est bien de carré intégrable (car à valeurs dans $[0, 1]$) et on obtient :

$$\mathbb{P}_\theta(|\widehat{\theta}_n - \theta| > \varepsilon) \leq \mathbb{P}_\theta(|\widehat{\theta}_n - \theta| \geq \varepsilon) \leq \frac{\text{Var}_\theta(\widehat{\theta}_n)}{\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}$$

car $\text{Var}_\theta(\widehat{\theta}_n) = \theta(1 - \theta)/n \leq 1/(4n)$ (étudier la fonction $x \mapsto x(1 - x)$ sur $[0, 1]$). On en déduit l'inégalité (I.1) en passant au complémentaire. Quant à l'inégalité (I.2), elle s'obtient en remarquant que $1/(4n\varepsilon^2) = \alpha$.

La proposition ci-dessus signifie qu'avec grande probabilité ($\geq 1 - \alpha$), l'estimateur $\widehat{\theta}_n$ est à distance au plus $1/\sqrt{4n\alpha}$ du paramètre inconnu θ . Cela permet de construire un intervalle aléatoire qui contient la vraie valeur de θ avec probabilité supérieure ou égale à $1 - \alpha$ (la preuve est une simple réécriture de (I.2)) :

Proposition 7 (Intervalle de confiance pour la proportion θ).

Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, pour tout $\alpha \in (0, 1)$, l'intervalle aléatoire

$$\widehat{I}_{n,\alpha} = \left[\widehat{\theta}_n - \frac{1}{\sqrt{4n\alpha}}; \widehat{\theta}_n + \frac{1}{\sqrt{4n\alpha}} \right]$$

contient la valeur inconnue θ avec probabilité au moins $1 - \alpha$, au sens où

$$\mathbb{P}_\theta(\widehat{I}_{n,\alpha} \ni \theta) \geq 1 - \alpha.$$

Remarque 1. *Attention, c'est l'intervalle $\widehat{I}_{n,\alpha}$ qui est aléatoire et non le paramètre θ .*

C'est pourquoi la notation $\mathbb{P}_\theta(\widehat{I}_{n,\alpha} \ni \theta)$ est peut-être moins ambiguë que $\mathbb{P}_\theta(\theta \in \widehat{I}_{n,\alpha})$.

Remarque 2 (Influence du nombre n d'observations).

La largeur de l'intervalle de confiance est proportionnelle à $1/\sqrt{n}$, donc plus le nombre n d'observations est grand, et plus notre intervalle de confiance est précis.

C'est logique ! Par ailleurs, le fait qu'on ait $1/\sqrt{n}$ plutôt que $1/n$ implique que pour réduire l'erreur d'un facteur 10 (c-à-d gagner un chiffre significatif), il ne faut pas prendre 10 fois plus, mais 100 fois plus d'observations. Nous n'avons pas le choix, c'est la théorie qui l'impose.

— En passant à l'événement complémentaire, on obtient

$$\mathbb{P}_\theta(\widehat{I}_{n,\alpha} \not\ni \theta) \leq \alpha$$

on dit que l'intervalle de confiance $\widehat{I}_{n,\alpha}$ est de niveau α .

- Cela correspond à un niveau de risque, puisqu'il s'agit de la probabilité que l'intervalle de confiance soit incorrect.
- Concrètement, cela signifie que si on répétait l'expérience aléatoire un très grand nombre de fois, alors l'intervalle de confiance obtenu serait incorrect moins de $100 \times \alpha\%$ (environ) des cas.
- En pratique, on n'observe qu'un seul échantillon et on ne calcule qu'un seul intervalle de confiance, mais on espère que l'échantillon observé fait parti des plus de $100 \times (1 - \alpha)\%$ échantillons favorables. Ce phénomène est illustré plus tard par une figure.
- Le contrôle de la probabilité $\mathbb{P}_\theta(\widehat{I}_{n,\alpha} \not\supset \theta) \leq \alpha$ est vrai quel que soit le nombre n d'observations. Si le contrôle n'était vrai qu'asymptotiquement (pour $n \rightarrow +\infty$), c'est-à-dire, si on avait seulement prouvé l'inégalité

$$\limsup_{n \rightarrow +\infty} \mathbb{P}_\theta(\widehat{I}_{n,\alpha} \not\supset \theta) \leq \alpha ,$$

alors on dirait que $\widehat{I}_{n,\alpha}$ serait un intervalle de confiance de niveau *asymptotique* α .

- Construire des intervalles de confiance *asymptotiques* peut paraître moins fort que construire des intervalles dont on contrôle le niveau pour tout $n \in \mathbf{N}^*$, mais cela a quand même un intérêt.
- En effet, comme le montre la proposition suivante, l'intervalle de confiance $\widehat{I}_{n,\alpha}$ construit avec l'inégalité de Bienaymé-Tchebychev est trop large quand n est grand.

Proposition 8 (Intervalle de confiance trop large).

Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, pour tout $\alpha \in (0, 1)$,

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta \left(\left[\widehat{\theta}_n - \frac{1}{\sqrt{4n\alpha}} ; \widehat{\theta}_n + \frac{1}{\sqrt{4n\alpha}} \right] \ni \theta \right) \geq 1 - \exp\left(-\frac{1}{2\alpha}\right)$$

avec $\exp(-1/(2\alpha)) \leq \alpha$ pour tout $\alpha \in (0, 1)$ et $\exp(-1/(2\alpha)) = o(\alpha)$ quand $\alpha \rightarrow 0$.

On a :

$$\begin{aligned} & \mathbb{P}_\theta \left(\left[\widehat{\theta}_n - \frac{1}{\sqrt{4n\alpha}} ; \widehat{\theta}_n + \frac{1}{\sqrt{4n\alpha}} \right] \ni \theta \right) \\ &= \mathbb{P}_\theta \left(\left| \sqrt{n} \frac{\widehat{\theta}_n - \theta}{\sqrt{\theta(1-\theta)}} \right| \leq \frac{1}{\sqrt{4\theta(1-\theta)\alpha}} \right) \xrightarrow{n \rightarrow +\infty} \int_{-\frac{1}{\sqrt{4\theta(1-\theta)\alpha}}}{\frac{1}{\sqrt{4\theta(1-\theta)\alpha}}} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \end{aligned} \quad (\text{I.3})$$

d'après le Théorème Central Limite (TCL). En posant $\bar{\Phi}(u) = \int_u^{+\infty} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx$, l'intégrale ci-dessus se réécrit $1 - 2\bar{\Phi}(u_\alpha)$ avec $u_\alpha = 1/\sqrt{4\theta(1-\theta)\alpha}$.

Inégalités utiles :

$$\forall u > 0, \quad \bar{\Phi}(u) \leq \frac{e^{-u^2/2}}{u\sqrt{2\pi}} \quad (\text{I.4})$$

$$\forall u \geq 1, \quad \bar{\Phi}(u) \leq \frac{e^{-u^2/2}}{2}. \quad (\text{I.5})$$

Preuve de la première inégalité : on pose $\varphi(x) = e^{-x^2/2}/\sqrt{2\pi}$ et on remarque que $\varphi'(x) = -x\varphi(x)$, si bien que

$$\varphi(u) = \int_u^{+\infty} x\varphi(x)dx \geq u \int_u^{+\infty} \varphi(x)dx = u\bar{\Phi}(u).$$

Ensuite, (I.5) découle¹ de (I.4) en remarquant que $u\sqrt{2\pi} \geq \sqrt{2\pi} \geq 2$ car $u \geq 1$.

Nous sommes maintenant en mesure de majorer l'intégrale apparaissant à droite de (I.3). Nous avons dit qu'il s'agissait de $1 - 2\bar{\Phi}(u_\alpha)$. Or, d'après l'inégalité (I.5) ci-dessus et la définition de $u_\alpha = 1/\sqrt{4\theta(1-\theta)\alpha}$,

$$2\bar{\Phi}(u_\alpha) \leq e^{-u_\alpha^2/2} = \exp\left(-\frac{1}{8\theta(1-\theta)\alpha}\right) \leq \exp(-1/(2\alpha)),$$

où la dernière inégalité découle du fait que la fonction $x \mapsto \exp(-1/x)$ est croissante sur \mathbf{R}_+^* et que $8\theta(1-\theta)\alpha \leq 2\alpha$ (rappelons que $\theta(1-\theta) \leq 1/4$).

- Par ailleurs, pour montrer l'inégalité $\exp(-1/(2\alpha)) \leq \alpha$, il suffit de poser $x = 1/\alpha$ et de vérifier que $xe^{-x/2} \leq 1$ pour tout $x > 0$, ce qui est vrai car $\frac{d}{dx}(xe^{-x/2}) = (1-x/2)e^{-x/2}$, si bien que $x \mapsto xe^{-x/2}$ atteint son maximum en $x = 2$ et ce maximum vaut $2/e \leq 1$.
- Enfin, la comparaison $\exp(-1/(2\alpha)) \underset{\alpha \rightarrow 0}{=} o(\alpha)$ découle du fait que $\lim_{x \rightarrow +\infty} xe^{-x/2} = 0$.
- On déduit de la proposition précédente que, quand le nombre n d'observations est grand, l'intervalle de confiance $\hat{I}_{n,\alpha}$ construit avec l'inégalité de Bienaymé-Tchebychev a un niveau de risque bien plus faible que α .
- Par conséquent, si on est prêt à accepter un niveau de risque asymptotique vraiment égal à α , on peut prendre un intervalle moins large, dicté par le Théorème de la Limite Centrale et le lemme de Slutsky.

1. En fait, l'inégalité (I.5) est vraie pour tout $u \geq 0$ mais la preuve est plus délicate.

Proposition 9 (Intervalle de confiance de niveau asymptotique α).

Supposons que les v.a. X_1, \dots, X_n sont i.i.d. de loi $\mathcal{B}(\theta)$ de paramètre θ inconnu. Alors, pour tout $\alpha \in (0, 1)$, si on note $q_{1-\alpha/2}$ le quantile d'ordre $1 - \alpha/2$ de la loi $\mathcal{N}(0, 1)$ (cf. figure I.1), on a :

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta \left(\left[\hat{\theta}_n - q_{1-\alpha/2} \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} ; \hat{\theta}_n + q_{1-\alpha/2} \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \right] \ni \theta \right) = 1 - \alpha .$$

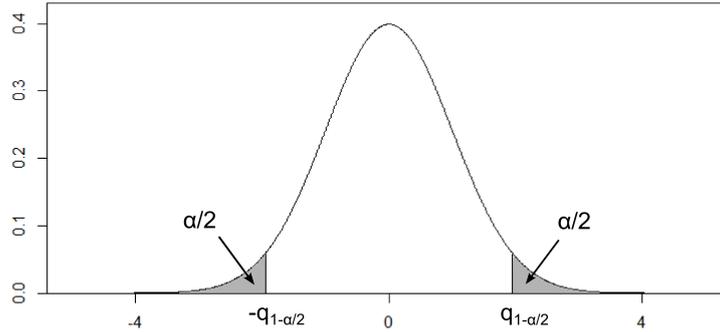


FIGURE I.1 – Représentation du quantile $q_{1-\alpha/2}$ d'ordre $1 - \alpha/2$ de la loi $\mathcal{N}(0, 1)$. Par symétrie, toute v.a. $X \sim \mathcal{N}(0, 1)$ vérifie $\mathbb{P}(|X| \leq q_{1-\alpha/2}) = 1 - (\alpha/2 + \alpha/2) = 1 - \alpha$.

Par le TCL et le lemme de Slutsky (déjà fait dans la première partie du cours). Rappelons à très grands traits la preuve :

$$\begin{aligned} & \mathbb{P}_\theta \left(\left[\hat{\theta}_n - q_{1-\alpha/2} \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} ; \hat{\theta}_n + q_{1-\alpha/2} \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \right] \ni \theta \right) \\ &= \mathbb{P}_\theta \left(\left| \hat{\theta}_n - \theta \right| \leq q_{1-\alpha/2} \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} \right) \\ &= \mathbb{P}_\theta \left(\left| \frac{\sqrt{\theta(1-\theta)}}{\sqrt{\hat{\theta}_n(1-\hat{\theta}_n)}} \sqrt{n} \frac{\hat{\theta}_n - \theta}{\sqrt{\theta(1-\theta)}} \right| \leq q_{1-\alpha/2} \right) \end{aligned} \quad (\text{I.6})$$

Or,

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

Par conséquent, d'après le lemme de Slutsky, la probabilité à droite de l'inégalité (I.6) tend vers la probabilité $\int_{-q_{1-\alpha/2}}^{q_{1-\alpha/2}} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx = 1 - \alpha$ (par définition du quantile $q_{1-\alpha/2}$).

- L'intervalle de confiance de la proposition 9 ci-dessus est plus fin que l'intervalle de confiance qu'on avait construit avec l'inégalité de Bienaymé-Tchebychev (car $\widehat{\theta}_n(1 - \widehat{\theta}_n) \leq 1/4$ et $q_{1-\alpha/2} \leq 1/\sqrt{\alpha}$ pour les mêmes raisons que dans la preuve de la proposition 8).
- Ainsi, si on est dans le régime $n \rightarrow +\infty$, on a une meilleure précision sur θ . On verra également en TD comment raffiner l'intervalle de confiance tout en préservant un niveau α garanti pour tout $n \in \mathbf{N}^*$ (avec l'inégalité de Hoeffding).
- En pratique, en sciences expérimentales ou en ingénierie, il est fréquent de prendre un niveau $\alpha = 5\% = 0.05$, si bien que puisque $q_{97.5\%} \leq 1.96$ (ces valeurs sont tabulées), on en déduit un intervalle de confiance de niveau asymptotique 5% très standard :

$$\left[\widehat{\theta}_n - 1.96 \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} ; \widehat{\theta}_n + 1.96 \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} \right]$$

En figure I.2, nous avons simulé 80 échantillons X_1, \dots, X_n de loi de Bernoulli de paramètre $p = 1/2$ et avons représenté les 80 intervalles de confiance obtenus. Cette figure permet de bien illustrer la notion de *niveau*.

Tests

- Jusqu'à maintenant, nous avons expliqué comment estimer une proportion θ inconnue, et comment fournir des marges d'erreurs $\pm \varepsilon$ pour un niveau de risque α donné.
- Dans ce qui suit, l'objectif ultime n'est plus d'estimer θ , mais de répondre à une question sur θ avec deux réponses au choix.
- Dans l'exemple des élections présidentielles décrit en introduction, la question la plus cruciale est de déterminer, à partir de l'échantillon de sondage, si $\theta \leq 1/2$ ou $\theta > 1/2$.
- On dit dans ce cas qu'on cherche à réaliser un test d'hypothèses ; deux hypothèses sont en jeu :

$$" \theta \in [0, 1/2]" \quad \text{et} \quad " \theta \in (1/2, 1]" .$$
- Ces deux hypothèses sont dites *composites* car les ensembles de valeurs possibles $[0, 1/2]$ ou $(1/2, 1]$ ne sont pas réduits à des singletons.

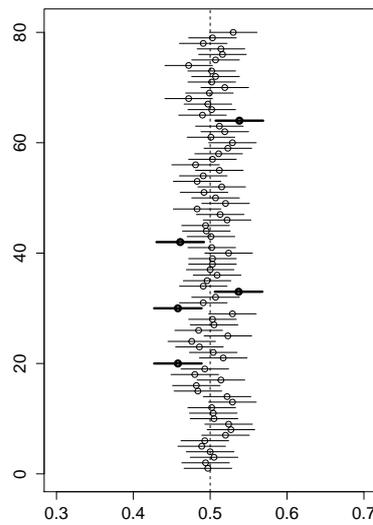


FIGURE I.2 – Visualisation des intervalles de confiance pour θ obtenus avec 80 échantillons (aléatoires) de taille $n = 1000$. Ces intervalles sont construits à partir du Théorème Central Limite et du lemme de Slutsky pour un niveau $\alpha = 0.05$. La vraie valeur de θ est $\theta = 0.5$. Sur les 80 intervalles de confiance obtenus, 5 ne contiennent pas la valeur de θ (en gras). Ainsi, une proportion de $5/80 = 6.25\%$ des échantillons ont donné un mauvais intervalle de confiance. On s'attendait à une proportion de l'ordre de 5% car on a choisi un niveau de risque asymptotique $\alpha = 0.05$.

- Un autre exemple de test de deux hypothèses composites faisant intervenir une proportion θ inconnue peut être tiré du domaine biomédical.
- Ainsi, dans le cadre d'un essai clinique visant à mettre sur le marché (ou non) un nouveau médicament, on peut vouloir tester si la proportion θ d'individus atteints par une certaine pathologie qui seraient guéris par ce nouveau médicament est supérieure au taux de guérison $\theta_0 = 0.7$ observé de longue date pour un médicament qui a déjà fait ses preuves.
- La proportion θ est relative à toute la population de personnes atteintes par cette pathologie. Pour des raisons bioéthiques évidentes, on ne peut pas mettre le médicament tout de suite le marché et attendre d'observer l'effet du nouveau médicament sur toute cette population.
- Un essai clinique est donc réalisé² sur un échantillon de taille n ; en d'autres termes, on effectue un "sondage". A partir du taux de guérison $\hat{\theta}_n$ observé sur l'échantillon, on cherche à identifier laquelle des deux hypothèses composites suivantes est correcte :

2. Le processus d'autorisation de mise sur le marché d'un médicament est un peu plus compliqué ; il y a en fait plusieurs phases d'essais cliniques. Nous avons simplifié volontairement.

" $\theta \leq \theta_0$ " (le nouveau médicament est moins efficace)
 versus " $\theta > \theta_0$ " (le nouveau médicament est plus efficace)

Comme nous allons le voir tout à l'heure, l'ordre dans lequel on écrit les deux hypothèses n'est pas du tout anodin (les conséquences statistiques sont importantes).

- Remarquons tout d'abord qu'il est vain de vouloir identifier correctement et à coup sûr (avec probabilité 1) laquelle des deux hypothèses

$$"\theta \leq \theta_0" \quad \text{et} \quad "\theta > \theta_0"$$

est vraie.

- En effet, à partir d'un échantillon de taille n , la connaissance qu'on a sur θ n'est pas d'une précision infinie : typiquement, elle se résume à un intervalle de confiance centré en $\hat{\theta}_n$ et de largeur de l'ordre de $1/\sqrt{n}$.
- Ainsi, si la proportion inconnue θ se situe à distance bien inférieure à $1/\sqrt{n}$ de θ_0 , il paraît difficile de pouvoir conclure de façon raisonnable entre l'une des deux hypothèses.
- Cet argument très informel peut être réécrit plus rigoureusement (à l'aide de concepts et d'outils de la *théorie de la décision*) pour montrer qu'en effet, cet objectif est trop ambitieux. On reviendra sur cette question en fin de chapitre.

- Commençons par décrire en détail un cadre volontairement plus simple : les deux hypothèses ne sont plus *composites* mais *simples*, c'est-à-dire de la forme

$$H_0 : "\theta = \theta_0" \quad \text{et} \quad H_1 : "\theta = \theta_1"$$

avec $\theta_0 \neq \theta_1$.

- On se place toujours dans le modèle de Bernoulli, c'est-à-dire que le statisticien dispose d'un échantillon $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \mathcal{B}(\theta)$ de paramètre θ inconnu.
- Dans ce cadre très simple, on peut introduire sans ambiguïté le vocabulaire usuel des tests.

- Un *test de l'hypothèse H_0 contre l'hypothèse H_1* est une fonction mesurable $\varphi_n : \mathcal{X}^n \rightarrow \{0, 1\}$ qui à un vecteur d'observations $(x_1, \dots, x_n) \in \mathcal{X}^n$ associe une décision $\varphi_n(x) \in \{0, 1\}$. Concrètement, lorsque $\varphi_n(X_1, \dots, X_n) = 0$, on dira qu'on conserve l'hypothèse H_0 , et lorsque $\varphi_n(X_1, \dots, X_n) = 1$, on dira qu'on rejette l'hypothèse H_0 et qu'on accepte l'hypothèse H_1 .
- Le *risque de première espèce* d'un test de deux hypothèses simples est la **probabilité de rejeter H_0 alors que H_0 est vraie**, ce qu'on peut noter

$$\mathbb{P}_{\theta_0}(\varphi_n(X_1, \dots, X_n) = 1)$$

- Le *risque de seconde espèce* d'un test de deux hypothèses simples est la **probabilité de conserver H_0 alors que H_1 est vraie**, ce qu'on peut noter

$$\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0)$$

- La *puissance* d'un test de deux hypothèses simples est reliée au risque de seconde espèce : c'est la probabilité d'accepter H_1 quand H_1 est vraie : $\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 1) = 1 - \mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0)$. La puissance quantifie la capacité du test à détecter l'hypothèse H_1 quand elle est vraie.
- On dit qu'un test de deux hypothèses simples est *de niveau α* ssi son risque de première espèce est inférieur ou égal à α , c-à-d $\mathbb{P}_{\theta_0}(\varphi_n(X_1, \dots, X_n) = 1) \leq \alpha$.
- Lorsqu'on a seulement un contrôle asymptotique

$$\limsup_{n \rightarrow +\infty} \mathbb{P}_{\theta_0}(\varphi_n(X_1, \dots, X_n) = 1) \leq \alpha ,$$

on dit que la suite de tests $(\varphi_n)_{n \in \mathbf{N}^*}$ est *de niveau asymptotique α* .

Résumons par un tableau :

	$\varphi_n = 0$	$\varphi_n = 1$
$\theta = \theta_0$	OK	erreur de 1ère espèce
$\theta = \theta_1$	erreur de 2nde espèce	OK

Supposons $\theta_0 < \theta_1$. Comment à partir de l'observation de

$$X_1, \dots, X_n \sim \mathcal{B}(\theta)$$

peut-on savoir si $H_0 : \theta = \theta_0$ ou $H_1 : \theta = \theta_1$ est vraie ? Une idée naive serait de comparer $\widehat{\theta}_n$ à $(\theta_0 + \theta_1)/2$.

- Si $\widehat{\theta}_n \leq \frac{\theta_0 + \theta_1}{2}$ on choisit H_0
- Si $\widehat{\theta}_n > \frac{\theta_0 + \theta_1}{2}$ on choisit H_1

et donc

$$\varphi_n(X_1, \dots, X_n) = \begin{cases} 0 & \text{si } \hat{\theta}_n \leq \frac{\theta_0 + \theta_1}{2} \\ 1 & \text{si } \hat{\theta}_n > \frac{\theta_0 + \theta_1}{2} \end{cases}$$

Problème : imaginons que θ_1 est proche de θ_0 par exemple $\theta_1 = \theta_0 + \frac{c}{\sqrt{n}}$ (pour n suffisamment grand et c petit) on a que le risque de première espèce :

$$\mathbb{P}_{\theta_0}(\hat{\theta}_n > \frac{\theta_0 + \theta_1}{2}) = \mathbb{P}_{\theta_0}(\hat{\theta}_n > \theta_0 + \frac{c}{2\sqrt{n}}) \tag{I.7}$$

$$= \mathbb{P}_{\theta_0} \left(\sqrt{n} \frac{\hat{\theta} - \theta_0}{\sqrt{\theta_0(1-\theta_0)}} > \frac{c}{2\sqrt{\theta_0(1-\theta_0)}} \right) \tag{I.8}$$

$$\simeq \int_{\frac{c}{2\sqrt{\theta_0(1-\theta_0)}}}^{\infty} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \quad (TCL) \tag{I.9}$$

Si c est petit la quantité $\frac{c}{2\sqrt{\theta_0(1-\theta_0)}}$ est proche de 0 et

$$\int_{\frac{c}{2\sqrt{\theta_0(1-\theta_0)}}}^{\infty} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \simeq \frac{1}{2}$$

Pour les mêmes raisons le risque de seconde espèce sera proche de 1/2

- Conséquences : si θ_0 et θ_1 sont proches, il faut choisir entre risque de 1ère espèce petit OU risque de 2nde espèce petit.
- Dissymétrie entre H_0 et H_1 : puisqu'il est parfois impossible de contrôler à la fois le risque de 1ère et celui de 2nde espèce, on décide de se concentrer sur l'erreur la plus grave.

Exemples :

- Test d'un nouveau médicament : on cherche à savoir si le taux de guérison θ d'un nouveau médicament est plus grand ou plus petit que celui d'un médicament utilisé depuis plusieurs décennies.
- On suppose que $\theta_0 < \tilde{\theta} < \theta_1$ correspond à l'efficacité du médicament sachant que $\tilde{\theta}$ correspond à l'efficacité de l'ancien médicament.

	$\varphi_n = 0$	$\varphi_n = 1$
$\theta = \theta_0$	OK	erreur de 1ère espèce : on choisit le nouveau médicament alors qu'il est moins efficace
$\theta = \theta_1$	erreur de 2nde espèce : on garde l'ancien médicament alors que le nouveau est plus efficace	OK

Dans cet exemple, l'erreur de première espèce est plus grave que celle de seconde espèce.

- Contrôle de qualité : on cherche à savoir si la proportion θ de produits défectueux dans un entrepôt dépasse ou non la norme autorisée par un contrat de vente (0.1% de produits défectueux). Deux rumeurs courent : $\theta = 0.07\%$ et $\theta = 0.15\%$. Si un expert mandaté par l'industriel effectue le test :

$$H_0 : \theta = 0.07$$

contre $H_1 : \theta = 0.15$

alors, du point de vue de l'industriel, c'est l'erreur de première espèce qui est la plus grave (car conclure " $\theta = 0.15$ " si la réalité est $\theta = 0.07$ le conduirait à revoir pour rien son processus de fabrication, d'où des coûts importants inutiles).

- NB : du point de vue du consommateur, c'est l'erreur dans l'autre sens qui est la plus grave. En pratique, on fera toujours en sorte que l'erreur la plus grave corresponde à l'erreur de 1ère espèce. Ainsi, le point de vue du consommateur conduirait à choisir les hypothèses de test suivantes :

$$H_0 : \theta = 0.15$$

contre $H_1 : \theta = 0.07$

Dans ce deuxième scénario, c'est bien l'erreur de 1ère espèce qui est la plus grave pour le consommateur, puisqu'elle correspond à ne pas sanctionner l'industriel alors que les produits sont défectueux.

Conclusion : dans tous les cas, on choisira les hypothèses H_0 et H_1 de sorte que **l'erreur la plus grave corresponde à l'erreur de 1ère espèce** (quand c'est possible).

Revenons au cas. de deux hypothèse simples

$$H_0 : \theta = \theta_0 \quad vs \quad H_1 : \theta = \theta_1$$

En accord avec le fait que l'on considère l'erreur de première espèce comme la plus grave, on choisit $\alpha \in (0, 1)$ tel que l'on contrôle le risque de première espèce par

$$\mathbb{P}_{\theta_0}(\varphi_n(X_1, \dots, X_n) = 1) \leq \alpha$$

ce qui va entraîner une contrainte sur φ_n

Il faut également prendre en compte le fait que l'on veut avoir un risque de deuxième espèce petit i.e

$$\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0) \text{ le plus petit possible}$$

En particulier $\varphi_n \equiv 0$ n'est pas un bon choix on a $\mathbb{P}_{\theta_0}(\varphi_n(X_1, \dots, X_n) = 1) = 0$ mais $\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0) = 1$

De manière concrète on va choisir judicieusement un nombre t_α tel que

$$\varphi_n = \mathbf{1}_{\{\hat{\theta}_n > t_\alpha\}}$$

On choisira la règle de décision suivante

$$\begin{aligned} \text{"} H_0 \text{"} & \text{ si } \hat{\theta}_n > t_\alpha \\ \text{"} H_1 \text{"} & \text{ si } \hat{\theta}_n \leq t_\alpha \end{aligned}$$

Reste à choisir t_α . Posons

$$t_\alpha = \theta_0 + \frac{1}{\sqrt{4n\alpha}}$$

On a d'après B-T

$$\mathbb{P}(\hat{\theta}_n > t_\alpha) \leq \mathbb{P}(|\hat{\theta}_n - \theta_0| > \frac{1}{\sqrt{4n\alpha}}) \leq \alpha$$

On considère donc

$$\varphi_n(X_1, \dots, X_n) = \mathbf{1}_{\hat{\theta}_n > \theta_0 + \frac{1}{\sqrt{4n\alpha}}}$$

Quid du risque de seconde espèce ?

$$\begin{aligned} \mathbb{P}_{\theta_1}(\hat{\theta}_n \leq t_\alpha) &= \mathbb{P}_{\theta_1}(\hat{\theta}_n \leq \theta_0 + \frac{1}{\sqrt{4n\alpha}}) \\ &= \mathbb{P}_{\theta_1}(\hat{\theta}_n \leq \theta_1 - (\theta_1 - \theta_0 - \frac{1}{\sqrt{4n\alpha}})) \\ &\stackrel{\leq B-T}{\leq} \frac{\text{Var}_{\theta_1}(\hat{\theta}_n)}{(\theta_1 - \theta_0 - \frac{1}{\sqrt{4n\alpha}})^2} = \frac{\theta_1(1 - \theta_1)}{n(\theta_1 - \theta_0 - \frac{1}{\sqrt{4n\alpha}})^2} \\ &\leq \frac{1}{4n(\theta_1 - \theta_0 - \frac{1}{\sqrt{4n\alpha}})^2} \end{aligned}$$

Supposons $\theta_1 \geq \theta_0 + \frac{1}{\sqrt{4n\alpha}} + \frac{1}{\sqrt{4n\beta}}$ alors

$$\mathbb{P}_{\theta_1}(\hat{\theta}_n \leq t_\alpha) \leq \beta$$

Puissance du test : capacité à détecter H_1 si H_1 est vraie. On a sous les conditions ci-dessus

$$\mathbb{P}_{\theta_1}(\hat{\theta}_n > t_\alpha) = 1 - \mathbb{P}_{\theta_1}(\hat{\theta}_n \leq t_\alpha) \geq 1 - \beta$$

- Attention, l'interprétation du résultat d'un test est très dangereuse : beaucoup de professionnels/étudiants interprètent mal ces résultats.
- Il faut avoir conscience que puisque les résultats d'un test entraînent souvent des décisions importantes (choix d'un médicament, choix d'une politique, choix d'un investissement économique, etc), l'interprétation est cruciale!

— *La statistique est un outil d'aide à la prise de décisions.*

La difficulté d'interprétation vient du fait qu'un raisonnement logique un peu subtil doit être réalisé. Voici le détail :

- Si les données x_1, \dots, x_n observées conduisent au résultat $\varphi_n(x_1, \dots, x_n) = 1$, il y a de bonnes raisons de croire que H_1 est vraie. En effet : si H_0 était vraie, on aurait $\mathbb{P}_0(\varphi_n(X_1, \dots, X_n) = 1) \leq \alpha$, ce qui signifie que la probabilité d'observer un échantillon conduisant à $\varphi_n(X_1, \dots, X_n) = 1$ serait petite. Puisqu'on a observé $\varphi_n(x_1, \dots, x_n) = 1$, par un pseudo-raisonnement par l'absurde, on fait le pari que c'est H_1 qui est vraie. On ne peut pas en être certain (le seul moyen d'être sûr serait de connaître θ , par exemple en faisant un recensement au lieu d'un sondage) ; il s'agit d'un pari au vu des données observées.
- A l'inverse, si les données x_1, \dots, x_n observées conduisent au résultat $\varphi_n(x_1, \dots, x_n) = 0$, la conclusion est beaucoup moins claire. On pourrait en effet être dans l'une des deux situations suivantes :
 - H_0 est vraie et le test a raison ;
 - H_1 est vraie et le test a tort ; malheureusement, la probabilité d'erreur $\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0)$ n'est pas aussi bien contrôlée que celle de l'erreur de première espèce. On ne peut donc pas exclure cette deuxième situation, *sauf si on sait que le risque de deuxième espèce $\mathbb{P}_{\theta_1}(\varphi_n(X_1, \dots, X_n) = 0)$ est lui aussi petit, i.e., si la puissance du test est proche³ de 1.*

En conclusion :

- Si le résultat d'un test est $\varphi_n = 1$, on peut conclure que les données observées sont vraisemblablement en désaccord avec l'hypothèse H_0 et donc conclure "vraisemblablement H_1 " ;
- Si le résultat d'un test est $\varphi_n = 0$, la conclusion la plus prudente est : "les données ne permettent pas d'exclure l'hypothèse H_0 ", autrement dit : "peut-être H_0 , dans le doute". Mais si on sait par ailleurs que le test est puissant, alors on peut faire le pari que H_0 est vraie.

Les considérations ci-dessus expliquent elles aussi comment choisir les hypothèses H_0 et

3. On verra un peu plus tard dans le cours comment étudier la puissance d'un test. Typiquement, il faut avoir à l'idée qu'un test est d'autant plus puissant que le nombre d'observations est grand.

H_1 . En effet, si on veut prouver une hypothèse au travers d'une expérience puis d'un test statistique, il faut mettre cette hypothèse comme H_1 (si on la met comme H_0 et que le résultat du test est $\varphi_n = 0$, ça ne prouve pas forcément grand chose). Le choix des hypothèses est donc le suivant :

H_0 : hypothèse par défaut, de prudence ("hypothèse nulle")

H_1 : hypothèse qui doit à tout prix être démontrée par l'expérience ("hypothèse alternative").

- $H_0 : \theta \leq \theta_0$ $H_1 : \theta > \theta_0$
- Risque de première espèce

$$\sup_{\theta \leq \theta_0} \mathbb{P}[\varphi_n(X_1, \dots, X_n) = 1]$$

- Risque de seconde espèce

$$\sup_{\theta > \theta_0} \mathbb{P}[\varphi_n(X_1, \dots, X_n) = 0]$$

Proposition 10. Soit $X_1, \dots, X_n \sim \mathcal{B}(\theta)$ avec $\theta \in (0, 1)$ inconnu et θ_0 fixé alors

$$\sup_{\theta \leq \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n > t] + \sup_{\theta > \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n \leq t] = 1,$$

pour tout $t \in \mathbb{R}$.

Le test est donc $\varphi_n = \mathbf{1}_{\widehat{\theta}_n > t}$

- L'application $\theta \rightarrow \mathbb{P}_\theta[\widehat{\theta}_n > t]$ est croissante et continue !
- Soit $\theta_0 \leq \theta_1$ et soit $U_1, \dots, U_n \sim \mathcal{U}([0, 1])$, on pose $X_i = \mathbf{1}_{U_i \leq \theta_0}$ et $Y_i = \mathbf{1}_{U_i \leq \theta_1}$ alors on a $X_i \leq Y_i$ pour tout i de sorte que

$$\left\{ \frac{1}{n} \sum_{i=1}^n X_i \geq t \right\} \subset \left\{ \frac{1}{n} \sum_{i=1}^n Y_i \geq t \right\}$$

et donc

$$\mathbb{P}_{\theta_0}[\widehat{\theta}_n > t] = \mathbb{P}\left[\left\{ \frac{1}{n} \sum_{i=1}^n X_i \geq t \right\}\right] \tag{I.10}$$

$$\leq \mathbb{P}\left[\left\{ \frac{1}{n} \sum_{i=1}^n Y_i \geq t \right\}\right] \tag{I.11}$$

$$= \mathbb{P}_{\theta_1}[\widehat{\theta}_n > t] \tag{I.12}$$

- La continuité en θ se montre en faisant intervenir la loi binomiale car $n\widehat{\theta}_n \sim \mathcal{B}(n, \theta)$.

— Donc on a par continuité et croissance

$$\sup_{\theta \leq \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n > t] = \mathbb{P}_{\theta_0}[\widehat{\theta}_n > t] \text{ (croissance)} \quad (\text{I.13})$$

$$\sup_{\theta > \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n \leq t] = \mathbb{P}_{\theta_0}[\widehat{\theta}_n \leq t] \text{ (croissance + continuité)} \quad (\text{I.14})$$

et la somme vaut naturellement 1 ce qui conclue la preuve.

— On introduit la fonction puissance d'un test φ_n :

$$\theta \rightarrow \mathbb{P}_\theta[\varphi_n = 1]$$

— Puisque $\theta \rightarrow \mathbb{P}_\theta[\widehat{\theta}_n > t]$ est croissante le test $\varphi_n = \mathbf{1}_{\widehat{\theta}_n > t_\alpha}$ satisfait

Proposition 11. Soit $X_1, \dots, X_n \sim \mathcal{B}(\theta)$ avec $\theta \in (0, 1)$ alors le choix de

$$t_\alpha = \theta_0 + \frac{1}{\sqrt{4n\alpha}}$$

permet d'assurer

$$\sup_{\theta \leq \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n > t_\alpha] \leq \alpha$$

Par ailleurs

$$\sup_{\theta > \theta_0} \mathbb{P}_\theta[\widehat{\theta}_n \leq t_\alpha] \leq \beta$$

— On peut remarquer que pour tout $\theta > \theta_0$

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta[\widehat{\theta}_n > t_\alpha] = 1$$

— On dit dans ce cas que la suite $(\mathbf{1}_{\widehat{\theta}_n > t_\alpha})$ est consistente ou alors on dit que le test est consistant.

— On a

$$\mathbb{P}_\theta[\widehat{\theta}_n > t_\alpha] = 1 - \mathbb{P}_\theta[\widehat{\theta}_n \leq t_\alpha] \quad (\text{I.15})$$

$$\geq 1 - \frac{1}{\sqrt{4n \left(\theta - \theta_0 - \frac{1}{\sqrt{4n\alpha}} \right)}} \quad (\text{I.16})$$

vertiges	pas de vertige
108	1055

— Nombre total de patients 1163

- Test $H_0 : \theta \leq 8 \text{ pourcent}$ $H_1 : \theta > 8 \text{ pourcent}$
- On choisit $\alpha = 0,05$

$$\varphi_n = \mathbf{1}_{\hat{\theta}_n > 0,08 + \frac{1}{\sqrt{4n0,05}}}$$

- $\hat{\theta}_n = 0,0929 < 0,08 + \frac{1}{\sqrt{4n0,05}} = 0,1456$
- Donc on conserve H_0 ou plutôt les observations n'excluent pas que H_0 soit vrai
- On prend un risque α de choisir H_0 vraie.

- On rappelle que l'on pose

$$\Phi(v) = \int_v^\infty \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx$$

qui est une fonction décroissante.

- Revenons au test précédent en considérant n grand. Nous avons

$$\mathbb{P}_{\theta_0} \left[\hat{\theta}_n > \theta_0 + \frac{1}{\sqrt{4n\alpha}} \right] = \mathbb{P} \left[\sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\sqrt{\theta_0(1-\theta_0)}} > \frac{1}{\sqrt{4\alpha\theta_0(1-\theta_0)}} \right] \quad (\text{I.17})$$

$$\rightarrow_n \Phi \left(\frac{1}{\sqrt{4\alpha\theta_0(1-\theta_0)}} \right) \quad (\text{I.18})$$

$$\leq \Phi \left(\frac{1}{\sqrt{\alpha}} \right) \quad \text{car } \theta_0(1-\theta_0) \leq 1/4 \quad (\text{I.19})$$

$$\leq \frac{\exp\left(-\frac{1}{2\alpha}\right)}{2} = o(\alpha) \quad (\text{I.20})$$

- On peut donc choisir t_α plus petit. Avec le quantile $q_{1-\alpha}$ d'ordre α on peut choisir

$$t_\alpha = \theta_0 + q_{1-\alpha} \sqrt{\frac{\theta_0(1-\theta_0)}{n}}$$

avec $\Phi(q_{1-\alpha}) = \alpha$

- Dessin de la densité
- On a alors

$$\sup_{\theta \geq \theta_0} \mathbb{P}_\theta[\hat{\theta}_n > t_\alpha] = \mathbb{P}_{\theta_0}[\hat{\theta}_n > t_\alpha] \quad (\text{I.21})$$

$$= \mathbb{P}_{\theta_0} \left[\sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\sqrt{\theta_0(1-\theta_0)}} > \frac{1}{\sqrt{4\alpha\theta_0(1-\theta_0)}} \right] \quad (\text{I.22})$$

$$\rightarrow_n \Phi(q_{1-\alpha}) = \alpha \quad (\text{I.23})$$

— En pratique on calcul

$$T = \sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\sqrt{\theta_0(1 - \theta_0)}}$$

où T est parfois appelée variable de test

— Si $T > q_{1-\alpha}$ on rejette H_0 et on fait le pari que H_1 est vrai

— Si $T \leq q_{1-\alpha}$ on ne rejette pas H_0

— Retour à l'exemple

vertiges	pas de vertige
108	1055

— Nombre total de patients 1163

— Test $H_0 : \theta \leq 8 \text{ pourcent}$ $H_1 : \theta > 8 \text{ pourcent}$

— On calcule

$$T = \sqrt{n} \frac{\hat{\theta}_n - \theta_0}{\sqrt{\theta_0(1 - \theta_0)}} \simeq 1,617$$

— Le quantile pour $\alpha = 0,05$ est

$$q_{1-0,05} = 1,645$$

— Ainsi $1,617 \leq 1,645$ on ne peut pas rejeter H_0 au vu des observations

— La notion de p -valeur est fréquemment utilisée en science expérimentale

— Donnons sa définition dans le cadre d'un test du type

$$\varphi_n = \mathbf{1}_{\hat{\theta}_n > t_\alpha}$$

Definition 1. On se place dans le cadre d'un test

$$H_0 : \theta \leq \theta_0 \quad H_1 : \theta > \theta_0$$

et on considère un test du type $\varphi_n = \mathbf{1}_{\hat{\theta}_n > t_\alpha}$ où

$$\alpha \rightarrow t_\alpha$$

est une fonction décroissante. On appelle p -valeur de cette famille de test associée à l'observation (X_1, \dots, X_n) la variable aléatoire

$$\hat{\alpha}(X_1, \dots, X_n) = \sup\{\alpha, \hat{\theta}_n \geq t_\alpha\}$$

- La p -valeur est la première valeur du risque pour lequel on rejette H_0

$$\hat{\alpha}(X_1, \dots, X_n) = \sup\{\alpha, \hat{\theta}_n \leq t_\alpha\} \quad (\text{I.24})$$

$$= \inf\{\alpha, \hat{\theta}_n > t_\alpha\} \quad (\text{I.25})$$

- Si $\alpha > \hat{\alpha}$ on a $t_\alpha \leq t_{\hat{\alpha}}$ et donc $\hat{\theta}_n > t_\alpha$ et on rejette H_0 .
- On rejette H_0 quand la p -valeur est petite.
- Attention la p -valeur est une variable aléatoire et donc dépend de notre observation. Pour un autre jeu de donnée cette valeur serait différente.
- **Rq** : Le risque α correspond à la probabilité de rejeter H_0 alors que H_0 est vraie mais il est incorrect d'interpréter la p -valeur de la même manière car celle-ci change selon les observations.

- Estimation du paramètre θ dans un modèle de Bernoulli avec un échantillon $X_1, \dots, X_n \sim \mathcal{B}(\theta)$

- Estimateur

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

- C'est un estimateur sans biais

$$\mathbb{E}_\theta[\hat{\theta}_n] = \theta$$

- La variance est donnée par

$$\text{Var}_\theta(\hat{\theta}_n) = \frac{\theta(1-\theta)}{n}$$

- C'est un estimateur asymptotiquement consistant : LFGN

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \theta .$$

- Une fois obtenu une estimation de θ on veut construire un intervalle de confiance de risque α

- Byénaimé Tchebyshev

$$\hat{I}_{n,\alpha} = \left[\hat{\theta}_n - \frac{1}{\sqrt{4n\alpha}}; \hat{\theta}_n + \frac{1}{\sqrt{4n\alpha}} \right]$$

- TCL+ Slutsky

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

- Test d'hypothèse : H_0 vs H_1
- Risque de 1ere espèce : probabilité de rejeter H_0 alors que H_0 est vraie
- Risque de seconde espèce : probabilité de conserver H_0 alors que H_1 est vraie
- Seuil du risque α permet d'obtenir un test du type

$$\varphi_n = \mathbf{1}_{\hat{\theta}_n > t_\alpha}$$

- Pour déterminer t_α on utilise soit B-T

$$t_\alpha = \theta_0 + \frac{1}{\sqrt{4n\alpha}}$$

- Soit TCL

$$t_\alpha = \theta_0 + q_{1-\alpha} \sqrt{\frac{\theta_0(1-\theta_0)}{n}}$$

Soit X_1, \dots, X_n des variables aléatoires réelles i.i.d. de densité $f : \mathbf{R} \rightarrow \mathbf{R}$ (par rapport à la mesure de Lebesgue sur \mathbf{R}) définie par

$$f_\theta(x) = \frac{x^{\theta-1} e^{-x}}{\Gamma(\theta)} \mathbf{1}_{x \geq 0},$$

où $\Gamma(\theta) = \int_0^{+\infty} x^{\theta-1} e^{-x} dx$, et où $\theta \in \mathbf{R}_+^*$ est un paramètre inconnu. On note \mathbf{E}_θ , Var_θ et \mathbb{P}_θ les espérances, variances et probabilités calculées avec cette valeur de θ .

1. Montrer que $\Gamma(t+1) = t \cdot \Gamma(t)$ pour tout $t > 0$.
2. Montrer que $\mathbf{E}_\theta(X_1) = \theta$ et $\text{Var}_\theta(X_1) = \theta$.
3. Montrer que $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ est un estimateur sans biais et fortement consistant du paramètre inconnu θ .
4. Soit $\alpha \in]0, 1[$. Construire un intervalle $[A_n, B_n]$ tel que $\mathbb{P}_\theta(A_n \leq \theta \leq B_n) \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha$.
5. On suppose ici que $\theta \in]0, 9]$. Calculer $\text{Var}_\theta(\hat{\theta}_n)$ et en déduire un intervalle de confiance de niveau α pour le paramètre θ , valable pour tout $n \in \mathbf{N}$.

Un de vos amis affirme que moins de la moitié des gens seraient prêts à ramasser votre stylo si vous le laissez tomber dans un ascenseur. De nature optimiste, vous voulez prouver le contraire.

1. Quelle modélisation du problème parmi les deux suggestions suivantes vous paraît la plus appropriée ? Justifier brièvement, et dites à quoi correspondent θ et θ_0 .

$$\left\{ \begin{array}{l} H_0 : \theta \leq \theta_0 \\ H_1 : \theta > \theta_0 \end{array} \right. \quad \text{ou} \quad \left\{ \begin{array}{l} H_0 : \theta > \theta_0 \\ H_1 : \theta \leq \theta_0 \end{array} \right.$$

2. Vous faites l'expérience avec n personnes choisies au hasard (n grand) et calculez la proportion $\hat{\theta}_n$ de personnes ayant ramassé votre stylo. Donner la formule d'un test φ_n de niveau $\alpha = 5\%$ de H_0 contre H_1 (justifier en 8 lignes maximum).

Indication : les quantiles d'ordre 90%, 95% et 97.5% de la loi $\mathcal{N}(0, 1)$ valent respectivement 1.28, 1.64 et 1.96 environ.

3. Vous faites l'expérience avec $n = 100$ personnes et observez les résultats suivants :

a ramassé votre stylo	60	Quelle est votre conclusion ?
n'a pas ramassé votre stylo	40	

Chapitre II

Test non-paramétriques

- On va aborder des test dits "non paramétrique" c'est à dire que l'on va tester le jeu d'hypothèse

$$\begin{cases} H_0 : F = F_{ref} \\ H_1 : F \neq F_{ref} \end{cases}$$

où F_{ref} correspond à une fonction de répartition (f.d.r) de référence et on a un échantillon

$$X_1, \dots, X_n \sim F$$

où F est une f.d.r inconnue et les X_i sont i.i.d

- On parle de test non paramétriques car F est décrite par un nombre infini de paramètres ($F(t), t \in \mathbb{R}$)
- On rappelle que $F(t) = \mathbb{P}[X_1 \leq t] = \mathbb{P}[X_i \leq t]$

- L'objet mathématique qui va nous permettre de faire l'estimation de la fonction F inconnue est la **fonction de répartition empirique**

$$\begin{aligned} F_n : \mathbb{R} &\rightarrow [0, 1] \\ t &\rightarrow \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t} \end{aligned}$$

- La valeur

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t}$$

dépend des valeurs prises par l'échantillon X_1, \dots, X_n qui est aléatoire donc la fonction F_n est une fonction aléatoire.

— Si on définit la mesure

$$\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i(\omega)}$$

alors

$$F_n(t) = \mu_n(] - \infty, t])$$

— C'est donc bien une f.d.r, on l'appelle empirique (comme pour la moyenne empirique) car elle dépend des résultats de l'expérience.

— Dessin

— Fixons $t \in \mathbb{R}$ comme les v.a (X_i) sont i.i.d alors les v.a

$$(\mathbf{1}_{X_i \leq t})$$

sont également i.i.d et on a que pour tout i

$$\mathbb{E}(\mathbf{1}_{X_i \leq t}) = \mathbb{P}[X_i \leq t] = \mathbb{P}[X_1 \leq t]$$

— On peut donc appliquer la LFGN

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t} \xrightarrow[n \rightarrow +\infty]{\text{p.s.}} \mathbb{P}[X_1 \leq t] = F(t)$$

— Les v.a ($\mathbf{1}_{X_i \leq t}$) sont des v.a i.i.d qui suivent la loi $\mathcal{B}(\mathbb{P}[X_1 \leq t]) = \mathcal{B}(F(t))$

— La v.a $nF_n(t)$ compte donc le nombre de succès ainsi

$$nF_n(t) \sim \mathcal{B}(n, F(t))$$

— On a appliqué la LFGN, le TCL nous donne

$$\sqrt{n}(F_n(t) - F(t)) \xrightarrow[n \rightarrow +\infty]{\text{Loi}} \mathcal{N}(0, F(t)(1 - F(t)))$$

$$\frac{\sqrt{n}}{\sqrt{F(t)(1 - F(t))}} \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{1}_{X_i \leq t} - F(t)) \right) \xrightarrow[n \rightarrow +\infty]{\text{Loi}} \mathcal{N}(0, 1)$$

— Jusqu'ici on a considéré des résultats lorsque t était fixé et on a obtenu des convergences lorsque F_n était évalué en t .

— Si on fixe un t_0 alors pour un n_0 grand on peut dire des choses sur la proximité de $F_n(t_0)$ et $F(t_0)$ mais si on change et qu'on regarde un temps t_1 alors on aura un autre n_1 qui n'a aucun lien à priori avec le premier n_0 .

- Ainsi à chaque $t \in \mathbb{R}$ correspond un n_t suffisamment grand tel que $F_{n_t}(t)$ est proche de $F(t)$. Ce résultat n'est pas satisfaisant car on voudrait pouvoir prendre un n suffisamment grand indépendant de t . Une sorte d'uniformité.
- On aimerait établir des résultats du type

$$\|F_n - F\| = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow +\infty]{p.s.} 0$$

- On a le résultat suivant appelé parfois théorème fondamental de la statistique ou théorème de Glivenko Cantelli
- theoreme 1** (Glivenko Cantelli). *Soit (X_n) une suite de v.a.i.i.d de f.d.r F alors*

$$\|F_n - F\| = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow +\infty]{p.s.} 0$$

- On doit donc passer d'une convergence simple à une convergence uniforme
 - Maintenant qu'on a un résultat de convergence presque sûre uniforme de F_n vers F on a donc un résultat de consistance. On voudrait maintenant pouvoir faire le test.
 - Pour cela nous avons le résultat suivant
- theoreme 2** (Kolmogorov). *Soit (X_n) une suite de v.a.i.i.d de f.d.r F CONTINUE alors la loi de la v.a*

$$D_n = \|F_n - F\|_\infty$$

ne dépend pas de F

- Notons qu'il y a une contrainte c'est la continuité de F .
- Nous donnons ici des éléments de preuve des deux résultats ci-dessus.
- Dans le premier résultat on ne suppose pas F forcément continue
- Pour établir ces résultats nous allons utiliser l'inverse généralisée de la f.d.r

$$F^{(-1)}(u) = \inf\{t : F(t) \geq u\}, u \in]0, 1[$$

- On rappelle que si $U \sim \mathcal{U}([0, 1])$ alors

$$X = F^{(-1)}(U)$$

a pour f.d.r la fonction F .

— Soit U_1, \dots, U_n n v.a.i.i.d $\mathcal{U}([0, 1])$ alors

$$X_1, \dots, X_n = F^{(-1)}(U_1), \dots, F^{(-1)}(U_n)$$

sont n v.a.i.i.d de f.d.r la fonction F .

— Ainsi on a

$$\begin{aligned} \|F_n - F\|_\infty = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| &= \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t} - F(t) \right| \\ &\sim \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| \end{aligned} \quad (\text{II.1})$$

— Or on a $\{F^{(-1)}(U_i) \leq t\} = \{U_i \leq F(t)\}$ donc

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| = \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq F(t)} - F(t) \right| \quad (\text{II.2})$$

— Dans le sup ci-dessus on pose $s = F(t)$ on a donc

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| = \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq F(t)} - F(t) \right| \quad (\text{II.3})$$

$$= \sup_{s \in F(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \quad (\text{II.4})$$

— En utilisant que le *sup* sur $F(\mathbb{R})$ est plus petit que le sup sur $[0, 1]$ on a

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| = \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq F(t)} - F(t) \right| \quad (\text{II.5})$$

$$= \sup_{s \in F(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \quad (\text{II.6})$$

$$\leq \sup_{s \in [0, 1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \quad (\text{II.7})$$

— Grâce à la loi forte des grands nombres on sait

$$\forall s \in [0, 1], \exists A_s, t, q, \mathbb{P}(A_s) = 1, \text{ et } \forall \omega \in A_s, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— Pour pouvoir utiliser la phrase

$$\forall s \in [0, 1], \exists A_s, t.q, \mathbb{P}(A_s) = 1, \text{ et } \forall \omega \in A_s, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

on veut trouver un ensemble A indépendant de s c'est à dire obtenir

$$\exists A, t.q \forall s \in [0, 1], \forall \omega \in A, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— Problème déjà évoqué l'ensemble A_s dépend de s . Pourquoi n'a t-on pas le droit de considérer

$$\bigcap_{s \in [0, 1]} A_s \quad \text{car c'est une intersection non dénombrable}$$

— On va se ramener à un ensemble dénombrable suffisamment gros pour pouvoir trouver cet ensemble A

— On va utiliser le fait que \mathbb{Q} est dénombrable on a donc

$$\forall s \in [0, 1] \cap \mathbb{Q}, \exists A_s, t.q, \mathbb{P}(A_s) = 1, \text{ et } \forall \omega \in A_s, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— Ainsi on peut considérer l'ensemble

$$A = \bigcap_{s \in [0, 1] \cap \mathbb{Q}} A_s$$

qui est une intersection dénombrable donc c'est un ensemble mesurable.

— Comme on a $\forall s \in [0, 1] \cap \mathbb{Q}, \mathbb{P}(A_s) = 1$ on a que $\mathbb{P}(A) = 1$

— Reprenons la phrase

$$\forall s \in [0, 1] \cap \mathbb{Q}, \exists A_s, t.q, \mathbb{P}(A_s) = 1, \text{ et } \forall \omega \in A_s, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

et regardons ce qui se passe pour notre ensemble $A = \bigcap_{s \in [0, 1] \cap \mathbb{Q}} A_s$

— Soit A l'ensemble ci-dessus et soit $s \in [0, 1] \cap \mathbb{Q}$ et $\omega \in A$ alors nécessairement $s \in A_s$ et donc

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— On a donc montré

$$\exists A, t.q \forall s \in [0, 1] \cap \mathbb{Q}, \forall \omega \in A, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— Mais on voudrait la même chose sans se restreindre à \mathbb{Q}

$$\exists A, t.q \forall s \in [0, 1], \forall \omega \in A, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \xrightarrow[n \rightarrow +\infty]{} s$$

— Pour passer de $s \in [0, 1] \cap \mathbb{Q}$ à $s \in [0, 1]$, on va utiliser la densité de \mathbb{Q} . Soit $s \in [0, 1]$, alors il existe deux suites (s_k) et (t_k) telle que

$$s_k \rightarrow s, \text{ et } t_k \rightarrow s$$

et $s_k \leq s$ et $t_k \geq s$. On peut donc approcher s par dessus et par dessous par des éléments de \mathbb{Q} .

— Ainsi $\{U_i \leq s_k\} \subset \{U_i \leq s\} \subset \{U_i \leq t_k\}$ et par suite soit $\omega \in A$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s_k} \leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq t_k}$$

— On considère les limite en n . On doit utiliser la notion de liminf et limsup

$$\begin{aligned} s_k &= \lim_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s_k} \leq \liminf_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \\ &\leq \limsup_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \leq \lim_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq t_k} = t_k \end{aligned}$$

— On reprend la dernière inégalité on a donc

$$s_k \leq \liminf_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \leq \limsup_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} \leq t_k \quad (\text{II.8})$$

— On prend $k \rightarrow \infty$ et on a donc par le théorème des gendarmes

$$s = \liminf_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} = \limsup_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} = s$$

— On a donc montré que

$$\lim_n \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq s} = s$$

- La conclusion se base sur un théorème dû à Dini.
- Ce résultat dit que si (f_n) est une suite de fonctions croissante i.e $t \rightarrow f_n(t)$ est croissante en t et que la suite de fonction converge simplement vers une fonction f qui est continue alors la convergence est uniforme.
- On est exactement dans ce cadre. Soit $\omega \in A$ alors pour tout n la fonction

$$t \rightarrow F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i(\omega) \leq t}$$

est une fonction croissante en t .

- De plus la suite de fonction (F_n) est une suite de fonction croissante qui converge simplement vers t . Or $t \rightarrow t$ est une fonction continue donc le théorème s'applique et
- Attention, on a utilisé à un moment une égalité en loi. Montrer que la convergence presque sûre ci-dessus se transmet malgré l'égalité en loi est assez ardue et nous l'admettrons.

- Revenons au Théorème de Kolmogorov

theoreme 3 (Kolmogorov). *Soit (X_n) une suite de v.a.i.i.d de f.d.r F CONTINUE alors la loi de la v.a*

$$D_n = \|F_n - F\|_\infty$$

ne dépend pas de F

- Nous introduisons

$$G_n(u) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq u}$$

- Nous allons montrer que

$$\|F_n - F\| \sim \|G_n - G\|$$

où $G(u) = u$.

- On reprend

$$\sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| = \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq F(t)} - F(t) \right| \quad (\text{II.9})$$

$$= \sup_{s \in F(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \quad (\text{II.10})$$

$$(\text{II.11})$$

— On remarque que

$$]0, 1[\subset \text{Im}F = F(\mathbb{R}) \subset [0, 1]$$

car F est continue

— Ainsi on en déduit

$$\sup_{s \in]0, 1[} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \leq \sup_{s \in F(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \leq \sup_{s \in [0, 1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right|$$

— Or $G_n(1) = G(1)$ et $G_n(0) = G(0)$ donc

$$\sup_{s \in]0, 1[} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| = \sup_{s \in [0, 1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right|$$

— On a donc

$$\begin{aligned} \|F_n - F\|_\infty &= \sup_{t \in \mathbb{R}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{F^{(-1)}(U_i) \leq t} - F(t) \right| \\ &= \sup_{s \in F(\mathbb{R})} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \\ &= \sup_{s \in [0, 1]} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{U_i \leq s} - s \right| \\ &= \|G_n - G\|_\infty \end{aligned} \tag{II.12}$$

— Le résultat en découle

— On a donc besoin de connaître les quantiles de la loi

$$D_n = \|F_n - F\|_\infty$$

— On définit $d_{n, 1-\alpha}$ tel que

$$\alpha = \mathbb{P}[D_n > d_{n, 1-\alpha}]$$

— Ces valeurs sont connues (des bonnes approximations) tabulées.

— On s'intéresse souvent à D_n pour des valeurs grandes. En particulier il existe W_{max} telle que

$$\sqrt{n}D_n \xrightarrow[n \rightarrow +\infty]{\text{Loi}} W_{max}$$

— En particulier pour tout $x \geq 0$ et n grand

$$\mathbb{P}[\sqrt{n}D_n \leq x] \simeq \mathbb{P}[W_{max} \leq x]$$

— On en déduit que

$$\mathbb{P}[D_n > \frac{x}{\sqrt{n}}] \simeq \mathbb{P}[W_{max} > x]$$

et donc

$$d_{n,1-\alpha} \simeq \frac{w_{1-\alpha}}{\sqrt{n}}$$

— Pour n grand les quantiles de W_{max} donnent ceux de D_n

— On est en mesure de présenter le test de Kolmogorov qui va nous permettre de tester

$$\begin{cases} H_0 : F = F_{ref} \\ H_1 : F \neq F_{ref} \end{cases}$$

— On pose le test

$$\varphi_n(X_1, \dots, X_n) = \mathbf{1}_{D_n > d_{n,1-\alpha}},$$

avec $D_n = \|F_n - F_{ref}\|_\infty$

— Le résultat reste vrai même si F_{ref} n'est pas continue.

— On montrera en T.D que ce test est consistant i.e sous H_1

$$\mathbb{P}[\varphi_n(X_1, \dots, X_n) = 1] \xrightarrow[n \rightarrow +\infty]{} 1$$

Soit (X_1, \dots, X_n) un n -échantillon de fonction de répartition F inconnue. Notre seule connaissance sur F est qu'elle est continue. On note \mathbf{F}_n la fonction de répartition empirique associée à l'échantillon :

$$\forall t \in \mathbf{R}, \quad \mathbf{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t}.$$

On notera F^{-1} l'inverse généralisée de la fonction F .

1. On rappelle que si U_1, \dots, U_n sont des variables aléatoires i.i.d. uniformes sur $[0, 1]$, alors $(F^{-1}(U_1), \dots, F^{-1}(U_n))$ est un n -échantillon de fonction de répartition F . En déduire que la loi de $\|\mathbf{F}_n - F\|_\infty$ est la même quelle que soit F continue.

Dans la suite, pour tout $\alpha \in]0, 1[$, on notera $d_{n,1-\alpha}$ un quantile d'ordre $1 - \alpha$ de cette loi.

2. En cours, nous avons construit un intervalle de confiance centré en la moyenne empirique qui contenait le vrai paramètre θ d'une loi de Bernoulli avec probabilité supérieure ou égale à $1 - \alpha$. En vous inspirant de ce résultat, construisez une région de confiance $\widehat{\mathcal{F}}_n \subset [0, 1]^{\mathbf{R}}$ qui contienne la vraie fonction de répartition F avec probabilité supérieure ou égale à $1 - \alpha$.

1. Application : vous observez les valeurs d'échantillon suivantes :

0.38 0.16 0.04 0.63 0.44 0.27 0.51 0.32 0.06 0.13

2. Représentez graphiquement la fonction $t \mapsto \mathbf{F}_n(t; \omega)$ sur cet exemple. Expliquez brièvement comment on pourrait représenter sur le graphique la région de confiance $\widehat{\mathcal{F}}_n(\omega)$, pour le niveau de risque $\alpha = 5\%$.

On donne quelques valeurs approchées du quantile $d_{10,1-\alpha}$:

α	0.01	0.025	0.05	0.1
$d_{10,1-\alpha}$	0.489	0.445	0.409	0.369

3. Que signifie très concrètement la valeur 5% ?