

Estimation par intervalle

1. Définition d'une région de confiance

Soit $\alpha \in]0, 1[$ un *niveau de risque* fixé par le statisticien.

Définition 1. Une région de confiance de θ de niveau de confiance $1 - \alpha$ est un ensemble (dépendant de l'observation mais pas du paramètre inconnu θ), $C(X) \subset \Theta$, telle que

$$\forall \theta \in \Theta, \quad \mathbb{P}_\theta(\theta \in C(X)) \geq 1 - \alpha$$

On dit alors qu'on a une région par excès. Dans le cas où on a égalité on parle de niveau exactement égal à $1 - \alpha$.

Lorsqu'on a $X = (X_1, \dots, X_n)$, on parle de *région de confiance asymptotique* de niveau $1 - \alpha$, si

$$\forall \theta \in \Theta, \quad \lim_{n \rightarrow +\infty} \mathbb{P}_\theta(\theta \in C(X_1, \dots, X_n)) \geq 1 - \alpha$$

Les valeurs usuelles de α sont 1%, 5% ou 10%. Dans le cas unidimensionnel, la plupart du temps, une région de confiance s'écrit sous la forme d'un intervalle (unilatère ou bilatère). Un intervalle de confiance de niveau de confiance 95% a une probabilité au moins égale à 0,95 de contenir la vraie valeur inconnue θ . Par passage au complémentaire, le niveau de risque α correspondant à une majoration de la probabilité que la vraie valeur du paramètre θ ne soit pas dans $C(X)$. A niveau de confiance fixé, une région de confiance est d'autant meilleure qu'elle est de taille petite. Avant d'aller plus loin, rappelons la notion de quantile d'une loi de probabilité.

Définition 2. Soit $\alpha \in]0, 1[$. On appelle quantile d'ordre α d'une loi de probabilité \mathbb{P} , la quantité

$$z_\alpha = \inf \{x, \mathbb{P}([-\infty, x]) \geq \alpha\}.$$

Par exemple pour la loi $\mathcal{N}(0, 1)$, le quantile d'ordre 97,5% est 1.96, et celui d'ordre 95% est 1.645.

2. Construction de régions de confiance

Une première méthode consiste à appliquer l'inégalité de Bienaymé-Tchebychev. Rappelons que si X est une variable aléatoire ayant un moment d'ordre 2, alors

$$\forall \varepsilon > 0, \mathbb{P}(|X - \mathbb{E}(X)| > \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

Appliquons cette inégalité dans le cas de variables aléatoires indépendantes X_1, \dots, X_n identiquement distribuées de loi de Bernoulli $\mathcal{B}(\theta)$, où l'on souhaite estimer θ à l'aide de \bar{X}_n . On a

$$\forall \varepsilon > 0, \mathbb{P}(|\bar{X}_n - \theta| > \varepsilon) \leq \frac{\theta(1 - \theta)}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}.$$

On obtient ainsi une région de confiance de niveau $1 - \alpha$ en considérant

$$\left[\bar{X}_n - \frac{1}{2\sqrt{n\alpha}}, \bar{X}_n + \frac{1}{2\sqrt{n\alpha}} \right].$$

Pour $\alpha = 5\%$ et $n = 100$, la précision de l'intervalle est 0.22. Il faut noter que la majoration obtenue par l'application de l'inégalité de Bienaymé-Tchebychev n'est pas très précise.

On obtient un meilleur résultat en utilisant l'inégalité de Hoeffding (*Ouvrad, Probabilité 2, page 132*).

Proposition 3. Soit (X_1, \dots, X_n) une suite de variables aléatoires indépendantes telles que pour tout i , $a_i \leq X_i \leq b_i$ p.s., alors pour tout $\varepsilon > 0$, en posant $S_n = \sum_{i=1}^n X_i$,

$$\mathbb{P}(S_n - \mathbb{E}(S_n) \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right), \quad \mathbb{P}(S_n - \mathbb{E}(S_n) \leq -\varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

et

$$\mathbb{P}(|S_n - \mathbb{E}(S_n)| \geq \varepsilon) \leq 2 \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

Appliquons cette inégalité à l'exemple précédent. On a :

$$\forall \varepsilon > 0, \mathbb{P}(|\bar{X}_n - \theta| \geq \varepsilon) \leq \mathbb{P}(|S_n - n\theta| \geq n\varepsilon) \leq 2 \exp(-2n\varepsilon^2).$$

On obtient ainsi l'intervalle de confiance de niveau $1 - \alpha$ suivant :

$$\left[\bar{X}_n - \sqrt{\frac{1}{2n} \ln(2/\alpha)}, \bar{X}_n + \sqrt{\frac{1}{2n} \ln(2/\alpha)} \right].$$

Pour $\alpha = 5\%$ et $n = 100$, la précision de l'intervalle est 0.14 et 0.23 avec la première méthode.

Il peut s'avérer plus pratique de chercher un intervalle de confiance asymptotique.

Supposons que nous cherchions un intervalle de confiance pour un paramètre θ à partir d'un échantillon de taille n de loi \mathbb{P}_θ . Lorsqu'on dispose de suffisamment de données et pour les modèles les plus classiques, le théorème central limite s'avère être un très bon outil, pour obtenir un intervalle de confiance asymptotique. Par exemple si on souhaite estimer la moyenne d'une variable aléatoire dont on connaît la variance $\sigma^2 = 1$. On prend un n -échantillon (X_1, \dots, X_n) . L'application du TCL donne :

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

On obtient alors l'intervalle de confiance asymptotique de niveau α suivant

$$\left[\bar{X}_n - \frac{q_{1-\alpha/2}}{\sqrt{n}}, \bar{X}_n + \frac{q_{1-\alpha/2}}{\sqrt{n}} \right]$$

où $q_{1-\alpha/2}$ est le quantile d'ordre $1 - \alpha/2$ de la loi $\mathcal{N}(0, 1)$.

Ce n'est pas toujours aussi évident. Si on part d'une variable aléatoire de Bernoulli dont on veut estimer le paramètre θ . En considérant l'estimateur du maximum de vraisemblance \bar{X}_n , le TCL donne :

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

Ici la loi limite dépend de θ ce qui est gênant pour construire un intervalle de confiance. Dans ce cas, on peut surmonter ce problème, en remarquant que $\theta(1 - \theta) \leq 0.25$. On obtient donc un intervalle de confiance asymptotique :

$$\left[\bar{X}_n - \frac{q_{1-\alpha/2}}{2\sqrt{n}}, \bar{X}_n + \frac{q_{1-\alpha/2}}{2\sqrt{n}} \right].$$

Dans le cas où on considère (X_1, \dots, X_n) un échantillon de loi de Poisson de paramètre $\theta > 0$ à estimer, le TLC donne :

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

Des outils plus élaborés doivent être utilisés pour construire un intervalle de confiance si on ne connaît pas de majorant de θ . Le *lemme de Slutsky* permet de surmonter certaines difficultés comme le montre l'exemple suivant.

Si on reprend l'exemple, en utilisant les propriétés de forte consistance d'estimateur, on obtient :

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

mais aussi avec

$$\overline{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2$$

qui est un estimateur fortement consistant de la variance

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

3. Exemples classiques d'estimation par intervalle

Soit X une variable aléatoire de loi normale $\mathcal{N}(m, \sigma^2)$.

3.a. Estimation de la moyenne quand la variance est connue.

Théorème 4.

Lorsque σ^2 est connu un intervalle de confiance au niveau $1 - \alpha$ de m est

$$\left[\overline{X}_n - q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}, \overline{X}_n + q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

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

DÉMONSTRATION. On sait que $(\overline{X}_n - m)/\sigma$ suit la loi $\mathcal{N}(0, 1)$. Par conséquent on a

$$\frac{|\overline{X}_n - m|}{\sigma} \in [-q_{1-\alpha/2}, q_{1-\alpha/2}] \iff m \in \left[\overline{X}_n - q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}, \overline{X}_n + q_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

□

Exemple : puisque $q_{0,975} = 1.96$ l'intervalle de confiance de m au niveau 95% est

$$\left[\overline{X}_n - 1.96 \frac{\sigma}{\sqrt{n}}, \overline{X}_n + 1.96 \frac{\sigma}{\sqrt{n}} \right]$$

ce qui signifie que sur un grand nombre d'expériences cet intervalle contiendra effectivement m dans 95% des cas en moyenne.

3.b. Estimation de la moyenne quand la variance est inconnue.

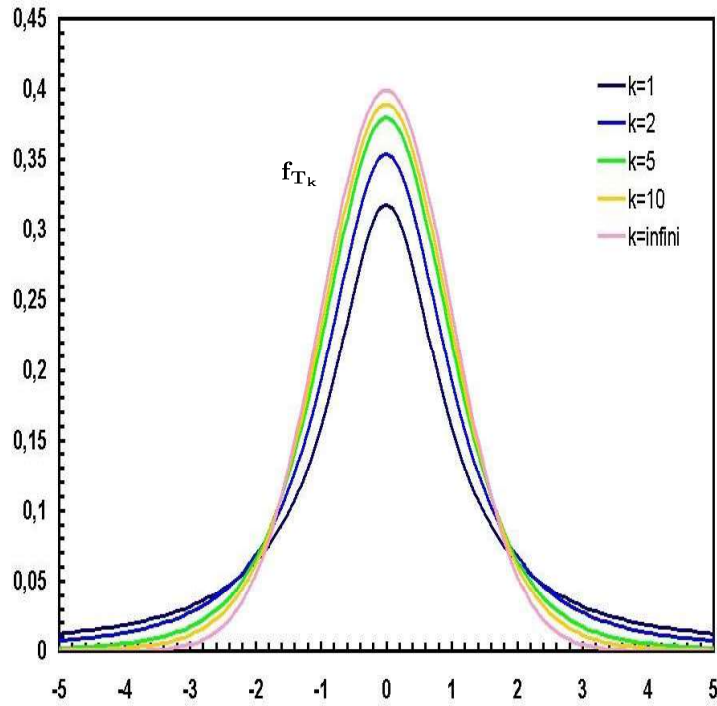
Définition 5. Soient X et Y deux variables aléatoires indépendantes suivant respectivement la loi normale centrée réduite et la loi du χ^2 à n degrés de liberté. La variable aléatoire

$$T = \frac{X}{\sqrt{Y/n}}$$

suit la loi de Student à n degrés de liberté. La densité de cette loi est donnée par :

$$f_T(u) = \frac{1}{\sqrt{n\pi}} \frac{\Gamma((n+1)/2)}{\Gamma(n/2)} \frac{1}{(1 + u^2/n)^{(n+1)/2}}$$

Cette variable n'a pas d'espérance pour $n = 1$ et pas de variance pour $n \leq 2$. Sinon on a $\mathbb{E}(T) = 0$ et $\text{Var}(T) = n/(n-2)$.



Théorème 6.

Lorsque σ^2 est inconnu un intervalle de confiance au niveau $1 - \alpha$ de m est

$$\left[\bar{X}_n - t_{n-1, 1-\alpha/2} \frac{\sqrt{\widehat{S}_n^2}}{\sqrt{n}}, \bar{X}_n + t_{n-1, 1-\alpha/2} \frac{\sqrt{\widehat{S}_n^2}}{\sqrt{n}} \right]$$

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.

Cela provient du résultat précédent et de l'estimation de σ^2 par \widehat{S}_n^2 .

Exemple : pour $n = 10$, avec un niveau de confiance de 95% et un intervalle symétrique on obtient l'intervalle

$$\left[\bar{X}_n - 2,26 \frac{\sqrt{\widehat{S}_n^2}}{\sqrt{n}}, \bar{X}_n + 2,26 \frac{\sqrt{\widehat{S}_n^2}}{\sqrt{n}} \right]$$

L'intervalle de confiance est plus grand que celui obtenu lorsqu'on connaît la variance.

3.c. Estimation de la variance quand la moyenne est connue.

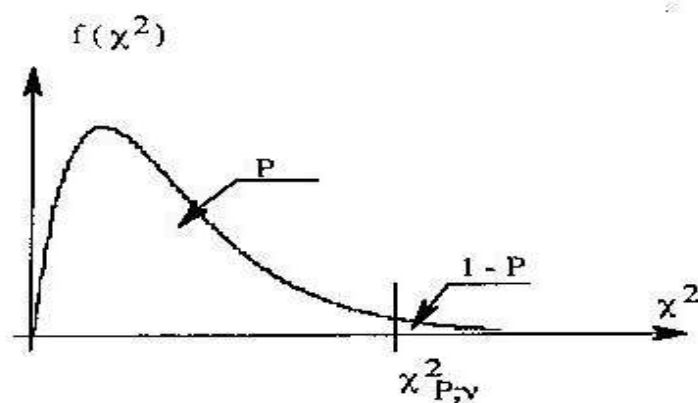
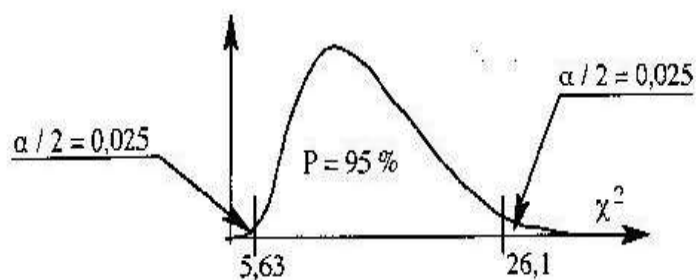
Théorème 7.

Lorsque m est connu un intervalle de confiance au niveau $1 - \alpha$ de σ^2 est

$$\left[\frac{1}{u_2} \sum_{k=1}^n (X_k - m)^2, \frac{1}{u_1} \sum_{k=1}^n (X_k - m)^2 \right]$$

où u_1 et u_2 sont les quantiles d'ordre $\alpha/2$ et $1 - \alpha/2$ de la loi du χ^2 à n degrés de liberté.

Exemples d'intervalles bilatères et unilatères pour la loi du χ^2 :



DÉMONSTRATION. Si X_k suit une loi $\mathcal{N}(m, \sigma^2)$ alors $(X_k - m)/\sigma$ suit une loi $\mathcal{N}(0, 1)$ et par conséquent $\sum_{k=1}^n \left(\frac{X_k - m}{\sigma}\right)^2$ suit une loi du χ^2 à n degrés de liberté. On définit alors u_1 et u_2 tels que

$$\mathbb{P}(\chi^2 \leq u_1) = \frac{\alpha}{2} \quad \text{et} \quad \mathbb{P}(\chi^2 \geq u_2) = \frac{\alpha}{2},$$

et donc on a

$$\mathbb{P} \left(u_1 \leq \frac{1}{\sigma^2} \sum_{k=1}^n (X_k - m)^2 \leq u_2 \right) = 1 - \alpha.$$

D'où on en déduit le résultat. □

3.d. Estimation de la variance quand la moyenne est inconnue.

Théorème 8.

Lorsque m est inconnu un intervalle de confiance au niveau $1 - \alpha$ de σ^2 est

$$\left[\frac{1}{u_2} \sum_{k=1}^n (X_k - \bar{X}_n)^2, \frac{1}{u_1} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \right]$$

où u_1 et u_2 sont les quantiles d'ordre $\alpha/2$ et $1 - \alpha/2$ de la loi du χ^2 à $n - 1$ degrés de liberté.

DÉMONSTRATION. On estime m par \bar{X}_n , puis

$$\sum_{k=1}^n \left(\frac{X_k - \bar{X}_n}{\sigma} \right)^2$$

suit une loi du χ^2 à $n - 1$ degrés de liberté. Ensuite on procède comme dans la preuve précédente. □

Lorsqu'on s'intéresse à l'écart-type on prend les racines carrées des bornes des intervalles obtenus pour la variance.

4. Comparaison de moyennes et de variances

Soient $(X_1, X_2, \dots, X_{n_1})$ un échantillon d'une population suivant la loi normale $\mathcal{N}(m_1, \sigma_1^2)$ et $(Y_1, Y_2, \dots, Y_{n_2})$ un échantillon d'une population suivant la loi normale $\mathcal{N}(m_2, \sigma_2^2)$; ces deux échantillons sont supposés indépendants. Nous souhaitons comparer les moyennes, m_1 et m_2 , et les variances, σ_1^2 et σ_2^2 , à l'aide de ces échantillons. Pour cela nous allons construire des intervalles de confiance pour $m_1 - m_2$ et pour σ_1^2 et σ_2^2 .

4.a. Intervalle de confiance de la différence de deux moyenne.

On pose

$$\bar{X} = \frac{1}{n_1} \sum_{i=1}^{n_1} X_i, \quad \bar{Y} = \frac{1}{n_2} \sum_{i=1}^{n_2} Y_i, \quad D = \bar{X} - \bar{Y}$$

Proposition 9. *L'estimateur D est un estimateur sans biais $m_1 - m_2$.*

DÉMONSTRATION. Par définition \bar{X} suit une loi $\mathcal{N}(m_1, \sigma_1^2/n_1)$ et \bar{Y} suit une loi $\mathcal{N}(m_2, \sigma_2^2/n_2)$ et par conséquent D suit la loi $\mathcal{N}(m_1 - m_2, \sigma_1^2/n_1 + \sigma_2^2/n_2)$, d'où le résultat. □

Théorème 10.

Si σ_1 et σ_2 sont connues, un intervalle de confiance de $m_1 - m_2$ au niveau $1 - \alpha$ est

$$\left[D - q_{1-\alpha/2} \sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}, D + q_{1-\alpha/2} \sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2} \right]$$

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

DÉMONSTRATION. Un intervalle de confiance de $m_1 - m_2$ au niveau $1 - \alpha$ est $[D - a, D + b]$ si

$$\begin{aligned} 1 - \alpha &= \mathbb{P}(D - a \leq m_1 - m_2 \leq D + b) = \mathbb{P}(-b \leq D - (m_1 - m_2) \leq a) \\ &= \mathbb{P}\left(-\frac{b}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} \leq \frac{D - (m_1 - m_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} \leq \frac{a}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}\right) \\ &= \frac{1}{\sqrt{2\pi}} \int_{u_1}^{u_2} \exp(-x^2/2) dx \end{aligned}$$

On en déduit l'intervalle annoncé. \square

En général les variances ne sont pas connues. Il peut alors se présenter deux cas.

Posons :

$$\widehat{S_{n_1, X}^2} = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_i - \bar{X})^2 \quad \widehat{S_{n_2, Y}^2} = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2$$

et

$$\widehat{S_{X, Y}^2} = \frac{1}{n_1 + n_2 - 2} \left(\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2 \right) = \frac{(n_1 - 1)\widehat{S_{n_1, X}^2} + (n_2 - 1)\widehat{S_{n_2, Y}^2}}{n_1 + n_2 - 2}$$

Théorème 11.

Si les variances σ_1 et σ_2 sont inconnues mais égales, un intervalle de confiance de $m_1 - m_2$ au niveau $1 - \alpha$ est

$$\left[D - t_{n_1+n_2-2, 1-\alpha/2} \sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}, D + t_{n_1+n_2-2, 1-\alpha/2} \sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2} \right]$$

où $t_{n_1+n_2-2, 1-\alpha/2}$ représente le fractile d'ordre $1 - \alpha/2$ de la loi de Student à $n_1 + n_2 - 2$ degrés de liberté.

On aurait pu remplacer σ_1^2 et σ_2^2 par $\widehat{S_{n_1, X}^2}$ et $\widehat{S_{n_2, Y}^2}$, mais en général on préfère prendre un estimateur basé sur la réunion des deux échantillons. On a déjà vu que :

$$Z_1 = \sum_{i=1}^{n_1} \left(\frac{X_i - \bar{X}}{\sigma_1} \right)^2 \quad \text{et} \quad Z_2 = \sum_{i=1}^{n_2} \left(\frac{Y_i - \bar{Y}}{\sigma_2} \right)^2$$

suivent respectivement une loi du χ^2 à $n_1 - 1$ et $n_2 - 1$ degrés de liberté. Par conséquent, comme Z_1 et Z_2 sont indépendantes, $Z_1 + Z_2$ suit une loi du χ^2 à $n_1 + n_2 - 2$ degrés de liberté. On obtient alors, en posant $\sigma_1^2 = \sigma_2^2 = \sigma^2$,

$$\mathbb{E}(Z_1 + Z_2) = n_1 + n_2 - 2 = \frac{1}{\sigma^2} \mathbb{E} \left[\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2 \right]$$

c'est-à-dire que $\widehat{S_{X, Y}^2}$ est un estimateur sans biais de σ^2 .

DÉMONSTRATION. On va remplacer l'écart-type de D , $\sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}$ par $\sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}$. On a

$$\begin{aligned} \alpha &= \mathbb{P}(D - a \leq m_1 - m_2 \leq D + b) = \mathbb{P}(-b \leq D - (m_1 - m_2) \leq a) \\ &= \mathbb{P}\left(-\frac{b}{\sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}} \leq \frac{D - (m_1 - m_2)}{\sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}} \leq \frac{a}{\sqrt{\widehat{S_{X, Y}^2}/n_1 + \widehat{S_{X, Y}^2}/n_2}}\right) \end{aligned}$$

On pose

$$T = \frac{D - (m_1 - m_2)}{\sqrt{\widehat{S_{X,Y}^2}/n_1 + \widehat{S_{X,Y}^2}/n_2}} = \frac{\frac{D - (m_1 - m_2)}{\sqrt{\sigma^2/n_1 + \sigma^2/n_2}}}{\sqrt{\frac{\widehat{S_{X,Y}^2}}{\sigma^2}}}$$

Le numérateur suit une loi normale centrée réduite. Au dénominateur on a $\sqrt{\frac{\widehat{S_{X,Y}^2}}{\sigma^2}} = \frac{U}{n_1 + n_2 - 2}$ où U suit une loi du χ^2 à $n_1 + n_2 - 2$ degrés de liberté. On en déduit que T suit une loi de Student à $n_1 + n_2 - 2$ degrés de liberté, et donc on obtient le résultat. \square

Lorsqu'on σ_1 et σ_2 sont inconnues mais non nécessairement égales, on utilise la méthode approchée suivante.

Théorème 12.

Si les échantillons ont des tailles importantes et égales à $n_1 = n_2 = n > 30$, un intervalle de confiance de $m_1 - m_2$ au niveau $1 - \alpha$ est

$$\left[D - q_{1-\alpha/2} \sqrt{(\widehat{S_{n_1,X}^2} + \widehat{S_{n_2,Y}^2})/n}, D + q_{1-\alpha/2} \sqrt{(\widehat{S_{n_1,X}^2} + \widehat{S_{n_2,Y}^2})/n} \right]$$

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

Il suffit de remarquer que la variable $\frac{D - (m_1 - m_2)}{\sqrt{(\widehat{S_{n_1,X}^2} + \widehat{S_{n_2,Y}^2})/n}}$ suit sensiblement une loi normale centrée réduite.

4.b. Intervalle de confiance du rapport de deux variances.

Théorème 13.

Un intervalle de confiance au niveau $1 - \alpha$ de σ_1^2/σ_2^2 est

$$\left[\frac{1}{f_{n_1-1, n_2-1, 1-\alpha/2}} \frac{\widehat{S_{n_1,X}^2}}{\widehat{S_{n_2,Y}^2}}, \frac{1}{f_{n_1-1, n_2-1, 1-\alpha/2}} \frac{\widehat{S_{n_1,X}^2}}{\widehat{S_{n_2,Y}^2}} \right]$$

où les fractiles sont ceux de loi de Fisher-Snedecor $\mathcal{F}(n_1 - 1, n_2 - 1)$.

Le résultat s'obtient par les mêmes méthodes que pour les théorèmes précédents.

La loi de Fisher-Snedecor peut être obtenue comme le quotient de deux lois du χ^2 indépendantes :

$$F(n_1 - 1, n_2 - 1) = \frac{\chi_{n_1-1}^2/(n_1 - 1)}{\chi_{n_2-1}^2/(n_2 - 1)}$$

5. Estimation d'une proportion

Dans une certaine population, la proportion d'individus ayant une propriété donnée est égale à p . Soit X le nombre d'individus d'un échantillon de taille n ayant la propriété.

5.a. Estimation ponctuelle.

Théorème 14.

Un estimateur sans biais et consistant de p est :

$$T = \frac{X}{n}$$

En effet, le nombre X d'individus de l'échantillon ayant la propriété suit la loi binomiale $B(n, p)$. On a :

$$\mathbb{E}(T) = \frac{\mathbb{E}(X)}{n} = p \quad \text{et} \quad \text{Var}(T) = \frac{\text{Var}(X)}{n^2} = \frac{p(1-p)}{n} \rightarrow 0$$

5.b. Estimation par intervalle.

On ne sait pas déterminer exactement un intervalle de confiance. On utilise des solutions approchées, qui fonctionnent lorsqu'on dispose d'échantillon de grande taille. Ainsi, lorsque n est grand ou/et p voisin de 0,5 on peut approcher la loi binomiale par une loi normale.

Rappel : Soit une suite de variables aléatoires Z_n suivant la loi binomiale $B(n, p)$; la suite des variables réduites $Z_n^* = \frac{Z_n - np}{\sqrt{np(1-p)}}$ converge en loi vers la loi normale centrée réduite, et on a :

$$\mathbb{P}(a \leq Z_n^* \leq b) \rightarrow \frac{1}{\sqrt{2\pi}} \int_a^b \exp(-x^2/2) dx$$

Théorème 15.

Un intervalle de confiance approché de p au niveau $1 - \alpha$ est donnée par

$$\left[T - q_{1-\alpha/2} \sqrt{\frac{T(1-T)}{n}}, T + q_{1-\alpha/2} \sqrt{\frac{T(1-T)}{n}} \right]$$

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

DÉMONSTRATION. D'après le rappel,

$$\frac{X - np}{\sqrt{np(1-p)}} = \frac{T - p}{\sqrt{p(1-p)/n}}$$

suit approximativement la loi normale centrée réduite. L'intervalle $[T - a, T + b]$ est un intervalle de confiance de p au niveau α si :

$$\begin{aligned} \alpha &= \mathbb{P}(T - a < p < T + b) = \mathbb{P}\left(-\frac{b}{\sqrt{p(1-p)/n}} < \frac{T - p}{\sqrt{p(1-p)/n}} < \frac{a}{\sqrt{p(1-p)/n}}\right) \\ &\simeq \Phi\left(\frac{a}{\sqrt{p(1-p)/n}}\right) - \Phi\left(-\frac{b}{\sqrt{p(1-p)/n}}\right) \quad \text{où} \quad \Phi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u \exp(-x^2/2) dx \end{aligned}$$

En choisissant un intervalle symétrique on obtient la quantile $q_{1-\alpha/2}$. Comme $\sqrt{p(1-p)}$ n'est pas connu, on obtient alors un intervalle de confiance approché en remplaçant p par l'estimateur T . \square

On peut donner une approximation de l'intervalle de confiance un peu plus précise.

Un intervalle de confiance approché de p au niveau $1 - \alpha$ est donnée par

$$\left[\frac{1}{1 + \frac{q_{1-\alpha/2}^2}{n}} \left(T + \frac{q_{1-\alpha/2}^2}{2n} - \Delta \right), \frac{1}{1 + \frac{q_{1-\alpha/2}^2}{n}} \left(T + \frac{q_{1-\alpha/2}^2}{2n} + \Delta \right) \right]$$

où $\Delta = \frac{q_{1-\alpha/2}}{\sqrt{n}} \sqrt{T(1-T) + \frac{q_{1-\alpha/2}^2}{4n}}$ et $q_{1-\alpha/2}$ représente le fractile d'ordre $1 - \alpha/2$ de la loi normale centrée réduite.

DÉMONSTRATION. On a :

$$-q_{1-\alpha/2} < \frac{T-p}{\sqrt{\frac{p(1-p)}{n}}} < q_{1-\alpha/2} \iff (T-p)^2 < q_{1-\alpha/2}^2 \frac{p(1-p)}{n}$$

tiononpeutcalculer

$$\text{une estimation de la moyenne et de l'écart-type :} \iff p^2 \left(1 + \frac{q_{1-\alpha/2}^2}{n} \right) - 2p \left(T + \frac{q_{1-\alpha/2}^2}{2n} \right) + T^2 < 0$$

Le paramètre p doit donc être compris entre les racines de l'équation du second degré. On vérifie aisément qu'elle a deux racines réelles appartenant à l'intervalle $[0, 1]$. D'où, on obtient l'intervalle de confiance indiqué.