

Chapitre II Estimations

1) Echantillon

- On dispose d'une population Ω et d'une variable aléatoire X sur Ω . (associé à une probabilité P).
- Un échantillon de taille n est la donnée de n éléments de Ω choisis aléatoirement.
- l'échantillon est dit non-exhaustif lorsque le choix / tirage des éléments se fait avec remise.
- Si la taille de Ω est grande par rapport à la taille de l'échantillon, on peut considérer que le tirage des éléments est non exhaustif.
- Soit $E = (\omega_1, \dots, \omega_n)$ un échantillon de taille n .
On applique X à chaque individu ω_i pour obtenir un n -uplet de nombres réels:
 $(X(\omega_1), \dots, X(\omega_n))$

On peut voir les choses de la façon suivante. ⁽²⁾

On part d'un n-uplet (X_1, \dots, X_n) de v.a.i qui suivent la même loi X .

Le n-uplet (X_1, \dots, X_n) est appelé échantillon de X .

Soit $E = (\omega_1, \dots, \omega_n)$ un échantillon de Ω

Posons $x_1 = X_1(\omega_1), \dots, x_n = X_n(\omega_n)$.

On dit que (x_1, \dots, x_n) est une réalisation de (X_1, \dots, X_n) .

Définition On appelle **statistique** de (X_1, \dots, X_n) une variable aléatoire $\varphi(X_1, \dots, X_n)$, où $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$.

Exemple. $\frac{X_1 + \dots + X_n}{n}$: moyenne empirique notée \bar{X}

$$s^2 = (X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2$$

$\frac{1}{n} ((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2)$ variance empirique notée $\frac{1}{n} s^2$.

Une estimation de la statistique $\varphi(X_1, \dots, X_n)$ ⁽³⁾
est la valeur prise par elle-même sur un échantillon
 (w_1, \dots, w_n) de taille n .

On pourra parler de la mojeune d'un échantillon,
d'un écart. type d'un échantillon, etc.

objectif. Comparer, par exemple, la mojeune
d'un échantillon avec la mojeune μ_X de la v.a. X

ou inversement. Peut-on avoir des informations sur μ_X
à partir de la mojeune d'un échantillon ?

2) Estimations ponctuelles

Soit X une v.a. de mojeune μ et d'écart. type σ .

Soit la statistique $\bar{X} = \frac{X_1 + \dots + X_n}{n}$ (mojeune empirique)

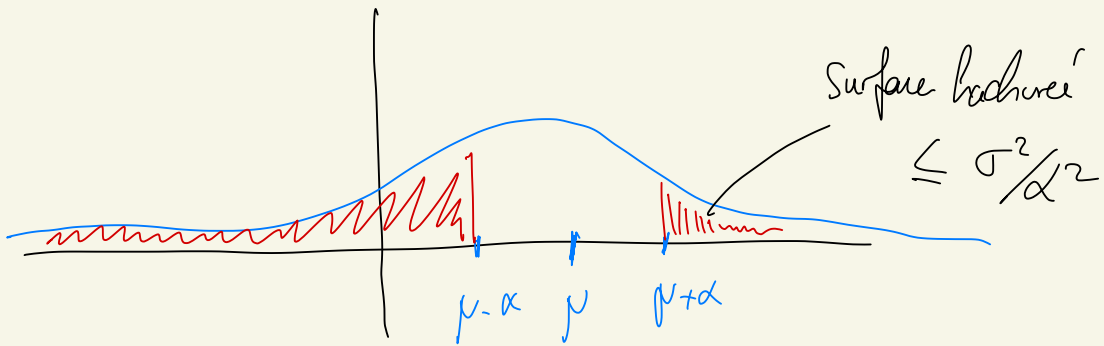
Alors $\mu_{\bar{X}} = \mu$ et $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$ (faits admis)

Inégalité de Bienaymé-Chebyshev (4)

Soit Z une v.a. de moyenne μ et d'écart-type σ .

Alors pour tout nombre $\alpha > 0$:

$$\mathbb{P}(|Z - \mu| \geq \alpha) \leq \sigma^2 / \alpha^2$$



Application à $Z = \bar{X} = \frac{X_1 + \dots + X_n}{n}$.

$$\mathbb{P}(|\bar{X} - \mu| \geq \alpha) \leq \frac{\sigma^2}{\alpha^2 n}$$

C'est le théorème des grands nombres.
Si α est petit et n grand, la probabilité que $\frac{X_1 + \dots + X_n}{n}$ soit loin de μ , est petite.

Pour n assez grand, l'inégalité de Bienaymé-Chebyshev ⁽⁵⁾
indique que $\bar{X}(w_1, \dots, w_n) = \frac{X_1(w_1) + \dots + X_n(w_n)}{n}$
est une bonne estimation de μ

De même que la moyenne d'un échantillon
est une bonne estimation de μ .

Qu'en est-il pour σ^2 ?

$$\text{Soit } \bar{S}^2 = \frac{1}{n} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right)$$

Variance empirique.

On a : $\bar{S}^2 = \frac{n-1}{n} \sigma^2$ et $\text{Var}(\bar{S}^2)$ petit
(fait admis) quand n grand.

Par l'inégalité de Bienaymé-Chebyshev, il vient :

$\bar{S}^2(w_1, \dots, w_n)$ est une bonne estimation de $\frac{n-1}{n} \sigma^2$

Or pour n grand, $\frac{n-1}{n} \sigma^2$ proche de σ^2 .

La variance d'un échantillon est une estimation
biaisée de σ^2 .

On pose alors :

(6)

$$\hat{S}^2 = \frac{n}{n-1} \quad \bar{S}^2 = \frac{1}{n-1} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right)$$

Alors
$$\mathbb{V} \hat{S}^2 = \frac{n}{n-1} \quad \mathbb{V} \bar{S}^2 = \sigma^2$$

Par le théorème de Bienaymé-Chebyshev, \hat{S}^2 est un bon estimateur de σ^2 . (= variance).

Conclusion Soit $(\omega_1, \dots, \omega_n)$ un échantillon de Ω .

Soit $x_1 = X_1(\omega_1), \dots, x_n = X_n(\omega_n)$ une réalisation de (X_1, \dots, X_n) .

Soit $\bar{x} = \frac{x_1 + \dots + x_n}{n}$ la moyenne de l'échantillon.

Alors, pour n grand, \bar{x} est une bonne approximation de μ .

et $\frac{1}{n-1} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)$ est une bonne approximation de la variance σ^2 .

3) Intervalle de confiance par une moyenne (7)

Soit X une variable aléatoire de moyenne μ et d'écart-type σ .

préliminaire

On cherche un intervalle contenant μ , aussi petit que possible, avec une probabilité contrôlée.

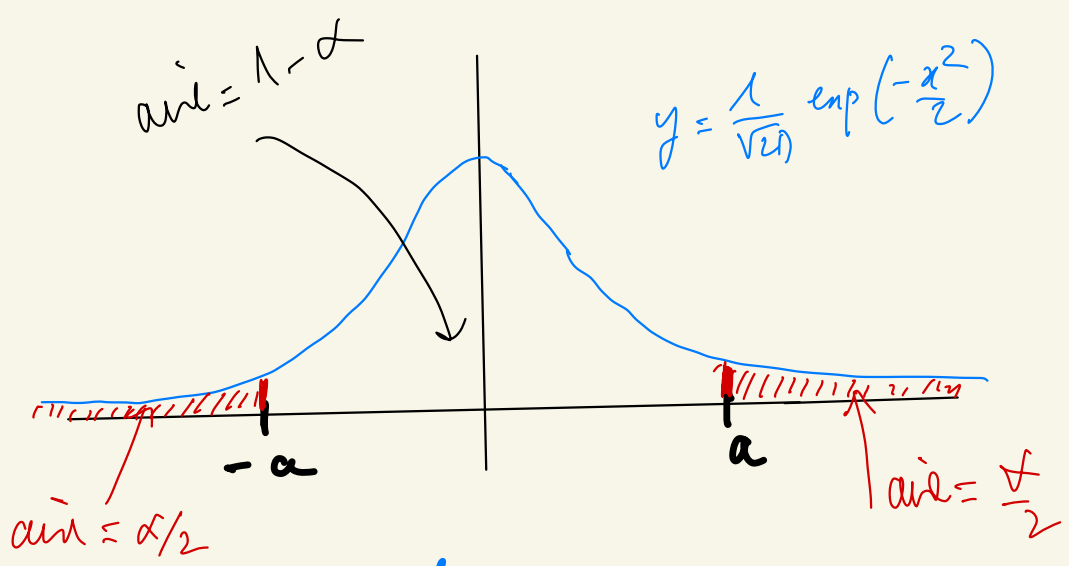
L'inégalité de Bienaymé-Chebichev permet d'en avoir un, mais ce n'est pas le "meilleur".

Pour n grand, on utilisera le théorème limite central... et la courbe de Gauss associée à la loi $\mathcal{N}(0,1)$.

Rapels Probabilités bilatérales de $\mathcal{N}(0,1)$.

Soit $Y \sim \mathcal{N}(0, 1)$. Soit α donné!
On cherche a positif tel que: (erreur).

$$P(|Y| \geq a) = \alpha$$



On pose $a = t_{\alpha/2}$.

$$P(|Y| \geq t_{\alpha/2}) = \alpha$$

$$P(|Y| \leq t_{\alpha/2}) = 1 - \alpha$$

Contexte

- X_1, \dots, X_n n variables indépendantes qui suivent la même loi X .

- μ = moyenne de X
 σ = écart-type de X .

- $\bar{X} = \frac{X_1 + \dots + X_n}{n}$ \rightsquigarrow $\bar{x} = \frac{x_1 + \dots + x_n}{n}$
échantillon

On veut donner un intervalle (petit si possible) contenant μ . Il va apparaître 3 cas :

- $X \sim \mathcal{N}(\mu, \sigma)$ avec σ connu
- $X \sim \mathcal{N}(\mu, \sigma)$ avec σ inconnu
- X quelconque, n grand.

Quand $X \sim \mathcal{N}(\mu, \sigma)$ avec σ connu

X_i : n v.a.i avec $X_i \sim \mathcal{N}(\mu, \sigma)$.

$\bar{X} = \frac{X_1 + \dots + X_n}{n}$: moyenne empirique.

Soit $Y = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$. Alors $Y \sim \mathcal{N}(0, 1)$.

Ainsi $P(|Y| \leq t_{\alpha/2}) = 1 - \alpha$.

(α = erreur)

Il vient ainsi l'intervalle estimé à partir d'un échantillon de taille n:

$|\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}| \leq t_{\alpha/2}$

ou encore:

$\bar{x} - t_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

\bar{x} : moyenne de l'échantillon

μ est dans l'intervalle $[\bar{x} - t_{\alpha/2} \frac{\sigma}{\sqrt{n}} ; \bar{x} + t_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$ avec une probabilité $1 - \alpha$.

longueur de l'intervalle: $\frac{2\sigma}{\sqrt{n}} t_{\alpha/2}$: meilleur que B.C.

Exemple. Le caractère d'une population suit une loi ⁽¹⁰⁾ normale d'écart. type $\sigma = 2$ mais de moyenne μ inconnue. On dispose d'un échantillon de taille 100 sur lequel la moyenne des valeurs de X est 2,5. Donner un intervalle de confiance pour μ avec $\alpha = 0,05$ (5%)

Ici $t_{\alpha/2} = 1,96$.

$$t_{\alpha/2} \times \frac{\sigma}{\sqrt{100}} = 1,96 \times \frac{2}{10} = 0,392.$$

D'où l'intervalle: $[2,5 - 0,392; 2,5 + 0,392]$
 $[2,108; 2,892]$

Quand $X \sim \mathcal{N}(\mu, \sigma)$ avec σ inconnu

X_1, \dots, X_n n v.d.c. avec $X_i \sim \mathcal{N}(\mu, \sigma)$

$\bar{X} = \frac{X_1 + \dots + X_n}{n}$: moyenne empirique

$\bar{S}^2 = \frac{1}{n} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right)$: variance empirique

$\hat{S}^2 = \frac{1}{n-1} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right) = \frac{n}{n-1} \bar{S}^2$

On rappelle que :

$$\frac{\bar{X} - \mu}{\sqrt{\frac{\bar{S}^2}{n-1}}} = \frac{\bar{X} - \mu}{\sqrt{\frac{S^2}{n}}} \sim \text{St}(n-1)$$

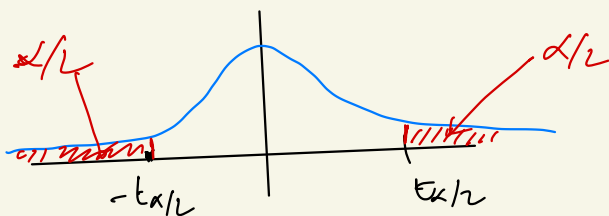
↑
loi de Student à $n-1$ degrés de liberté

La loi de Student a pour densité une fonction paire.
→ probabilités bilatérales.

On note $t_{\alpha/2}$ le nombre positif tel que :

$$P(|Y| \geq t_{\alpha/2}) = \alpha$$

Ici $Y \sim \text{St}(n)$ où n est donné.



exemple: $n=4$
 $r=3$
 $\alpha=0,20 (=20\%)$
 $t_{\alpha/2} = 1,64$

Ainsi

$$P\left(\left|\frac{\bar{X} - \mu}{\sqrt{\frac{\bar{S}^2}{n}}}\right| \geq t_{\alpha/2}\right) = \alpha$$

On arrive à l'intervalle (au risque α):

(12)

$$\bar{x} - t_{\alpha/2} \frac{\hat{\Delta}}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{\alpha/2} \frac{\hat{\Delta}}{\sqrt{n}}$$

où $\bar{x} = \frac{x_1 + \dots + x_n}{n}$: moyenne de l'échantillon

$$\hat{\Delta}^2 = \frac{1}{n-1} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)$$

$\hat{\Delta} = \sqrt{\hat{\Delta}^2}$: estimation de l'écart-type

On remarque que $\sqrt{\frac{n-1}{n}} \hat{\Delta} =$ écart-type de l'échantillon.

$t_{\alpha/2}$ est tel que

$$P(|Y| \geq t_{\alpha/2}) = \alpha \quad \text{quand } Y \sim St(n-1)$$

le cas d'une loi quelconque ~~X~~

(13)

On utilise le **théorème limite central**.

~~X~~ de moyenne μ et d'écart type σ .

X_1, \dots, X_n n v.a.i avec $X_i \sim X$.

Alors $\bar{X} = \frac{X_1 + \dots + X_n}{n}$ proche de $\mathcal{N}(\mu, \sigma/\sqrt{n})$.

quand n est grand.

Comme on prend "n grand", la variance σ^2 sera donnée par l'estimateur :

$$\hat{S}^2 = \frac{1}{n-1} \left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \right).$$

On arrive ainsi à l'approximation suivante :

$$P \left(\left| \frac{\bar{X} - \mu}{\hat{S}/\sqrt{n}} \right| \leq t_{\alpha/2} \right) = 1 - \alpha$$

donné par
la table de la loi
 $\mathcal{N}(0,1)$

d'in l'intervalle;

(15)

$$\bar{x} - t_{\alpha/2} \hat{\sigma} / \sqrt{n} \leq \mu \leq \bar{x} + t_{\alpha/2} \hat{\sigma} / \sqrt{n}$$

où $\bullet \bar{x} = \frac{x_1 + \dots + x_n}{n}$ moyenne de l'échantillon

$$\bullet \hat{\sigma}^2 = \frac{1}{n-1} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)$$

estimation de la variance

$\bullet \hat{\sigma} = \sqrt{\hat{\sigma}^2}$: estimation de l'écart-type

$\bullet t_{\alpha/2}$ est tel que

$$P(|Y| \leq t_{\alpha/2}) = 1 - \alpha$$

lorsque $Y \sim \mathcal{N}(0, 1)$

4) Intervalle de confiance par une fréquence (15)

Soit $X \sim \mathcal{B}(p)$ une loi de Bernoulli.

X_1, \dots, X_n des v.a.i. $X_i \sim \mathcal{B}(p)$.

$\bar{X} = \frac{X_1 + \dots + X_n}{n}$: "nombre moyen de succès".

loi binomiale

Pour n grand, on peut utiliser le théorème central limite.

Pour n grand, \bar{X} est proche de

$$\mathcal{N}\left(p, \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right)$$

c'est à dire

$$\frac{\bar{X} - p}{\sqrt{p(1-p)}} \sqrt{n} \text{ proche de } \mathcal{N}(0, 1)$$

Pour n grand $B(n, p)$ proche de $N(np, \sqrt{npq})$
bi-linéaire

np et $nq \geq 15$

ou grand

$n \geq 30$ et np et $nq \geq 5$.

Pour n grand on obtient l'intervalle de confiance

$$\bar{x} - t_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \bar{x} + t_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

donné par $N(0, 1)$

problème. L'intervalle dépend de p .

une méthode radicale

Lemme. Soit $0 \leq x \leq 1$. Alors $x(1-x) \leq 1/4$.

Ce qui permet d'obtenir l'intervalle

$$\bar{x} - t_{\alpha/2} \frac{1}{2\sqrt{n}} \leq p \leq \bar{x} + t_{\alpha/2} \frac{1}{2\sqrt{n}}$$

en estimant la variance (ici n est grand) (H)

$p(1-p)$ est proche de $\frac{1}{n-1} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)$

proche de $\frac{1}{n} \left((x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right)$

↑ variance de l'échantillon
estimation biaisée.

$$\bar{x}(1-\bar{x})$$

On obtient l'intervalle:

$$\bar{x} - t_{\alpha/2} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \leq p \leq \bar{x} + t_{\alpha/2} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}}$$

- \bar{x} = moyenne de l'échantillon
- $\bar{x}(1-\bar{x})$ = variance de l'échantillon
- $t_{\alpha/2}$ est donné par le tableau de $\mathcal{N}(0, 1)$

Exercice lors d'une élection entre deux (18)

Candidats, on fait un sondage sur 800 personnes pour "estimer" la proportion p d'électeurs qui voteront pour

A et q pour B. Ici $q = 1 - p$.

le sondage donne A à 48% et B à 52%.

① Donner un intervalle de confiance pour p
(avec $\alpha = 5\% = 0,05$)

② Quelle devrait être la taille de l'échantillon pour avoir un intervalle de longueur 1%?
(avec $\alpha = 5\%$).

① $\Omega =$ les électeurs. $(\mathcal{P}(\Omega))$: loi de Bernoulli.
 $n = 800$
 $\begin{matrix} p \rightarrow A \\ q \rightarrow B \end{matrix}$

Ici $\alpha = 0,05$ et $t_{\alpha/2} = 1,96$.

$$\bar{x} = 0,48$$

d'où l'intervalle

$$\bar{x} - t_{\alpha/2} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \leq p \leq \bar{x} + t_{\alpha/2} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}}$$

Ici $t_{\alpha/2} \sqrt{\frac{\bar{n}(n-\bar{n})}{n}} = 1,96 \times \sqrt{\frac{0,48 \times 0,52}{800}} \approx 0,0348$ (15)

$$0,4453 \leq p \leq 0,5134$$

longueur de l'intervalle
 $\approx 7\%$

(2) En utilisant le "lemme", on obtient un intervalle de longueur:

$$t_{\alpha/2} \times \frac{1}{\sqrt{n}}$$

"méthode radicale"

(ne dépend pas de \bar{n})

On veut $\frac{t_{\alpha/2}}{\sqrt{n}} \leq 0,01,$

c'est à dire :

$$\sqrt{n} \geq \frac{t_{\alpha/2}}{0,01} = 100 \times 1,96 = 196$$

d'où $n \geq 38416$

(20)

Remarque. Partant de $X \sim \mathcal{N}(\mu, \sigma)$, il est possible de donner un intervalle de confiance pour σ (que l'on connaisse μ ou non). Cela passe par

$$S^2 = (X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2 \text{ ou}$$

$$S_0^2 = (X_1 - \mu)^2 + \dots + (X_n - \mu)^2$$

puis la loi du **chi. deux** χ^2 .

Mais avec une attention particulière: la fonction de densité du χ^2 n'est pas paire (pas de symétrie).
