

# ECHANTILLONNAGE, ESTIMATION

## I Echantillonnage, estimation ponctuelle

exemple d'introduction: on considère une population de  $N$  individus et l'on s'attache à un caractère  $X$  v.a. de moyenne  $m$ , de variance  $\sigma^2$ .

Pour étudier  $X$ , on prélève un échantillon de  $n$  individus et l'on note les caractères observés  $x_1, \dots, x_n$ : réalisations de  $n$  v.a.  $X_1, \dots, X_n$ .

• Si le tirage se fait avec remise:  $X_1, \dots, X_n$  sont i.i.d. Donc  $E(X_i) = m$ ,  $\text{var}(X_i) = \sigma^2$ .

estimation de  $m$ :  $\hat{m}_n = \frac{1}{n} \sum_{i=1}^n X_i$

estimation ponctuelle:  $\frac{1}{n} \sum_{i=1}^n x_i$   
 $E(\hat{m}_n) = m$ ,  $\text{var}(\hat{m}_n) = \frac{\sigma^2}{n}$  → erreur d'échantillonnage.

$\hat{m}_n$ : estimateur de  $m$  sans biais.

2 variables, 1 expérience globale  $(\Omega \rightarrow S^2)$   
 $\omega \mapsto (X_1(\omega), \dots, X_n(\omega))$   
 n expériences séparées  $Z: \Omega^n \rightarrow S^{2n}$   
 $(\omega_1, \dots, \omega_n) \mapsto (X_1(\omega_1), \dots, X_n(\omega_n))$   
 Si  $\Pi: \Omega^n \rightarrow S^2$ ,  $Z = (X_1 \circ \Pi_1, \dots, X_n \circ \Pi_n)$ . En posant  $X_i \circ X_0 \Pi_i$ ,  $Z \in Y$ .

• Si le tirage se fait sans remise: les  $X_1, \dots, X_n$  ne sont plus indépendants

$E(\hat{m}_n) = m$  mais  $\text{var}(\hat{m}_n) = \frac{N-n}{N-1} \frac{\sigma^2}{n} \approx (1-f) \frac{\sigma^2}{n}$   
 $N$  grand  $f = \frac{n}{N}$

Plus  $f$  est grand, plus  $\text{var}(\hat{m}_n)$  sera petit, meilleure sera l'estimation. Mais ce cas est plus difficile à étudier que le précédent.

Problème: bien souvent  $\sigma$  est inconnue et on aura besoin d'estimer  $\sigma$  pour prévoir l'erreur d'échantillonnage.

Dans la suite on ne considèrera que des tirages avec remise (non exhaustifs)

## II Estimateurs.

1) Définition: Considérons une population suivant une loi dépendant d'un paramètre  $\theta \in \Theta$  inconnu. Un estimateur est une fonction d'un échantillon i.i.d.  $(X_1, \dots, X_n)$

$\hat{\theta}_n = f(X_1, \dots, X_n)$ . (v.a.)

Une estimation ponctuelle de  $\theta$  est une valeur observée  $f(x_1, \dots, x_n)$ .

$\hat{\theta}_n$  estimateur sans biais si  $E(\hat{\theta}_n) = \theta$ . Sinon, le biais est  $B(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta$ .

Objectif: construire des estimateurs, sans biais si possible, convergeant vers  $\theta$  et le plus rapidement possible.

Exemples: rappels: si  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $V = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$

pour une population telle que  $E(X) = m$ ,  $\text{var}(X) = \sigma^2$ ,  $E(X-m)^2 = m^2$ :

$$\left. \begin{aligned} E(\bar{X}) &= m, \text{var}(\bar{X}) = \frac{\sigma^2}{n} \\ E(V) &= \frac{n-1}{n} \sigma^2, \text{var}(V) \approx \frac{m^2 - \sigma^2}{n} \\ \text{cov}(\bar{X}, V) &= \frac{n-1}{n^2} \frac{m^2}{3} \end{aligned} \right\}$$

a) estimation de la moyenne  $m$ :  $\hat{m}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $E(\hat{m}_n) = m \rightarrow$  sans biais.

b) estimation de la variance  $\sigma^2$   
 1<sup>er</sup> cas:  $m$  connue.  $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ ,  $E(\hat{\sigma}_n^2) = \sigma^2 \rightarrow$  sans biais

2<sup>e</sup> cas:  $m$  inconnue.  $\hat{A}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{m}_n)^2$ ,  $E(\hat{A}_n^2) = (1 - \frac{1}{n}) \sigma^2$   
 → biais  $B(\hat{A}_n^2) = -\frac{\sigma^2}{n} \rightarrow 0$  asymptotiquement sans biais

\*  $\hat{S}_m^2 = \frac{n}{n-1} \hat{A}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{m}_n)^2$ ,  $E(\hat{S}_m^2) = \sigma^2 \rightarrow$  sans biais.

c) estimation d'une proportion: Soit  $p = \frac{N_A}{N}$  la proportion d'individus possédant une propriété  $A$ . On prélève  $n$  individus, on dénombre  $n_A$  possédant la propriété  $A$

→ proportion observée  $\frac{m}{n}$  (estimation ponctuelle).

En fait  $p$  est la moyenne de la loi de Bernoulli  $B(p)$ . Soit  $X_i = 1$  { i<sup>e</sup> individu a la prop A }  
On obtient ainsi l'estimateur  $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $E(\hat{p}_n) = E(X) = P(A) = p \rightarrow$  sans biais.

$\hat{n}_A = n \hat{p}_n$  suit la loi binomiale  $B(n, p)$ .

## 2) Précision d'un estimateur

- Soit  $\hat{\theta}_n$  un estimateur sans biais de  $\theta$ . On mesurera la précision de  $\hat{\theta}_n$  à l'aide de la variance  $\text{var}(\hat{\theta}_n) = E[(\hat{\theta}_n - \theta)^2]$ .  
Si  $\hat{\theta}_{1,n}$  et  $\hat{\theta}_{2,n}$  sont deux tels estimateurs,  $\hat{\theta}_{1,n}$  est plus précis que  $\hat{\theta}_{2,n}$  si  $\text{var}(\hat{\theta}_{1,n}) < \text{var}(\hat{\theta}_{2,n})$ .
- cas général (avec biais) : on mesure la précision de  $\hat{\theta}_n$  à l'aide de la fonction de risque  $R(\theta) = B(\hat{\theta}_n)^2 + \text{var}(\hat{\theta}_n) = E[(\hat{\theta}_n - \theta)^2]$ . La précision de  $\hat{\theta}_n$  est  $\frac{1}{R(\theta)}$ .

Définition : un estimateur sans biais de variance minimale est dit efficace.

- $\hat{\theta}_n$  est un estimateur convergent si  $\hat{\theta}_n \xrightarrow{P} \theta$ . (ou "correct")
- $\hat{\theta}_n$  est un estimateur absolument correct  $n \rightarrow \infty$  s'il est correct et sans biais.

Proposition : si  $\hat{\theta}_n$  vérifie  $\left\{ \begin{array}{l} E(\hat{\theta}_n) \rightarrow \theta \\ \text{ou sans biais} \end{array} \right\}$  et  $\text{var}(\hat{\theta}_n) \rightarrow 0$  alors  $\hat{\theta}_n$  est convergent vers  $\theta$ .

(dém: avec Bienaymé-Tchebychev)

Exemples : 1)  $\hat{m}_n = \frac{1}{n} \sum_{i=1}^n X_i$  : sans biais,  $\text{var}(\hat{m}_n) = \frac{\sigma^2}{n} \rightarrow 0$  donc  $\hat{m}_n \xrightarrow{P} m$  (absolument correct)

$$\begin{aligned} 2) \hat{\sigma}_n^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - m)^2 : \text{sans biais et } \text{var}(\hat{\sigma}_n^2) = \frac{1}{n^2} \sum_{i=1}^n \text{var}[(X_i - m)^2] \\ &= \frac{1}{n} \{ E[(X_1 - m)^4] - [E[(X_1 - m)^2]]^2 \} \\ &= \frac{m^4 - \sigma^4}{n} \xrightarrow{n \rightarrow \infty} 0 \text{ donc } \hat{\sigma}_n^2 \xrightarrow{P} \sigma^2 \text{ (absolument correct)} \end{aligned}$$

$$3) \hat{s}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{m}_n)^2 \rightarrow \text{avec biais.}$$

$$\begin{aligned} \hat{S}_n^2 &= \frac{n}{n-1} \hat{s}_n^2 : \text{sans biais} & \text{var}(\hat{S}_n^2) &= \frac{m^4 - \sigma^4}{n} - 2 \frac{m^4 - 2\sigma^4}{n^2} + \frac{m^4 - 3\sigma^4}{n^3} \\ & & &= \frac{(n-1)^2}{n^3} [m^4 - \frac{n-3}{n-1} \sigma^4] \\ \text{var}(\hat{S}_n^2) &= \frac{1}{n} [m^4 - \frac{n-3}{n-1} \sigma^4] \geq \text{var}(\hat{\sigma}_n^2) = \frac{m^4 - \sigma^4}{n} \\ & & &= \left(\frac{n}{n-1}\right)^2 \text{var}(\hat{s}_n^2) \geq \text{var}(\hat{s}_n^2) \end{aligned}$$

$\hat{S}_n^2$  a une plus grande variance que  $\hat{s}_n^2$  mais  $\hat{s}_n^2$  a un biais.  $\hookrightarrow$  avec biais

## 3) Estimateur du maximum de vraisemblance

Définition : vraisemblance d'un échantillon :  $L(x_1, \dots, x_n; \theta) = P_\theta(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i)$   
ou  $\prod_{i=1}^n f(x_i; \theta)$ .

On cherche le ou les  $\theta$  qui maximisent  $\theta \mapsto L(x_1, \dots, x_n; \theta)$ .

Ceci se fait en général en résolvant l'équation  $\frac{\partial}{\partial \theta} \ln L(x_1, \dots, x_n; \theta) = 0$ .

exemples : 1) loi binomiale  $B(p, \theta)$  :  $L(x_1, \dots, x_n; \theta) = \prod_{i=1}^n C_p^{x_i} \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$

$$\frac{\partial \ln L(x_1, \dots, x_n; \theta)}{\partial \theta} = \frac{\sum_{i=1}^n x_i}{\theta} - \frac{n - \sum_{i=1}^n x_i}{1-\theta} = 0 \Rightarrow \theta = \frac{\sum_{i=1}^n x_i}{n}$$

$$\rightarrow \hat{\theta}_{MDV} = \frac{\bar{X}}{Y}, \quad E(\hat{\theta}_{MDV}) = \theta, \quad \text{var}(\hat{\theta}_{MDV}) = \frac{\theta(1-\theta)}{n}, \quad \hat{\theta}_{MDV} \text{ absolument correct.}$$

2) loi de Poisson  $\mathcal{P}(\theta)$ :  $L(x_1, \dots, x_n; \theta) = \frac{\theta^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!} e^{-n\theta}$

$\frac{\partial}{\partial \theta} \ln L(x_1, \dots, x_n; \theta) = \frac{\sum_{i=1}^n x_i}{\theta} - n = 0 \Rightarrow \theta = \frac{1}{n} \sum_{i=1}^n x_i$

$\rightarrow \hat{\theta}_{MDV} = \bar{X}$ ,  $E(\hat{\theta}_{MDV}) = \theta$ ,  $\text{var}(\hat{\theta}_{MDV}) = \frac{\theta}{n}$ ,  $\hat{\theta}_{MDV}$  absolument correct.

3) loi exponentielle  $\mathcal{E}(\theta)$ :  $L(x_1, \dots, x_n; \theta) = \theta^n e^{-\theta \sum_{i=1}^n x_i}$

$\frac{\partial}{\partial \theta} \ln L(x_1, \dots, x_n; \theta) = \frac{n}{\theta} - \sum_{i=1}^n x_i = 0 \Rightarrow \theta = \frac{n}{\sum_{i=1}^n x_i}$

$\rightarrow \hat{\theta}_{MDV} = \frac{1}{\bar{X}}$

densité de  $\hat{\theta}_{MDV}$ :  $f_{\hat{\theta}_{MDV}}(t) = \frac{(n\theta)^n e^{-n\theta/t}}{(n-1)! t^{n+1}}$

$\Rightarrow E(\hat{\theta}_{MDV}) = \frac{n}{n-1} \theta$ ,  $E(\hat{\theta}_{MDV}^2) = \frac{n^2}{(n-1)(n-2)} \theta^2$

$\text{var}(\hat{\theta}_{MDV}) = \frac{n^2}{(n-1)^2(n-2)} \theta^2$ ,  $\hat{\theta}_{MDV}$  correct mais avec biais.

4) loi normale  $\mathcal{N}(\theta, \sigma^2)$ :  $L(x_1, \dots, x_n; \theta) = \frac{\exp[-\frac{1}{2\sigma^2}(\sum_{i=1}^n x_i^2 - 2\theta \sum_{i=1}^n x_i + n\theta^2)]}{(\sqrt{2\pi}\sigma)^n}$

$\frac{\partial}{\partial \theta} \ln L(x_1, \dots, x_n; \theta) = \frac{1}{\sigma^2} [\sum_{i=1}^n x_i - n\theta] = 0 \Rightarrow \theta = \frac{1}{n} \sum_{i=1}^n x_i$

$\hat{\theta}_{MDV} = \bar{X}$

$E(\hat{\theta}_{MDV}) = \theta$ ,  $\text{var}(\hat{\theta}_{MDV}) = \frac{\sigma^2}{n}$ ,  $\hat{\theta}_{MDV}$  absolument correct

5) loi normale  $\mathcal{N}(m, \theta)$ :  $L(x_1, \dots, x_n; \theta) = \frac{\exp[-\frac{1}{2\theta} \sum_{i=1}^n (x_i - m)^2]}{(\sqrt{2\pi\theta})^n}$

$\frac{\partial}{\partial \theta} \ln L(x_1, \dots, x_n; \theta) = \frac{1}{2\theta^2} \sum_{i=1}^n (x_i - m)^2 - \frac{n}{2\theta} = 0 \Rightarrow \theta = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2$

$\rightarrow \hat{\theta}_{MDV} = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ ,  $E(\hat{\theta}_{MDV}) = \theta$ ,  $\text{var}(\hat{\theta}_{MDV}) = \frac{m^4 - \sigma^4}{n} = \frac{2\theta^2}{n}$ , abs. correct.

### III Estimation par intervalle

Définition: Un intervalle de confiance de niveau de confiance  $1-\alpha$  ( $\alpha \in ]0,1[$ ) pour un paramètre  $\theta$  est un intervalle aléatoire  $[Z_1, Z_2]$  ( $Z_1, Z_2$ : v.a.) tel que  $\mathbb{P}(\theta \in [Z_1, Z_2]) = 1-\alpha$ .  $\alpha$  est le risque que  $\theta \notin [Z_1, Z_2]$ .

#### 1) Intervalle de confiance pour le moyenne $m$ d'une population $\mathcal{N}(m, \sigma^2)$

a) cas  $\sigma^2$  connue:  $\hat{m}_n = \frac{1}{n} \sum_{i=1}^n X_i$ :  $\mathcal{N}(m, \frac{\sigma^2}{n})$  donc  $N(\frac{\hat{m}_n - m}{\sigma/\sqrt{n}}) = \mathcal{N}(0,1)$ .

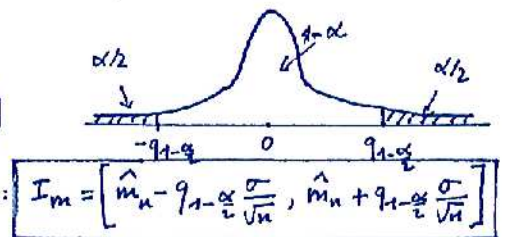
Soit  $q_{1-\frac{\alpha}{2}}$  la  $(1-\frac{\alpha}{2})$ -quantile de la loi  $\mathcal{N}(0,1)$ :  $\mathbb{P}(N \leq q_{1-\frac{\alpha}{2}}) = 1-\frac{\alpha}{2} \Rightarrow \mathbb{P}(|N| \leq q_{1-\frac{\alpha}{2}}) = 1-\alpha$

On a alors  $1-\alpha = \mathbb{P}\left(\left|\frac{\hat{m}_n - m}{\sigma/\sqrt{n}}\right| \leq q_{1-\frac{\alpha}{2}}\right)$

$= \mathbb{P}\left(m \in \left[\hat{m}_n - q_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \hat{m}_n + q_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right]\right)$   
 $= \mathbb{P}(m \in I_m)$

d'où un intervalle de confiance pour  $m$  de niveau  $1-\alpha$ :

Longueur de  $I_m = 2q_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ .



- Plus on demande un niveau de confiance élevé ( $\alpha$  plus petit), plus long( $I_m$ ) sera grande (intervalle moins précis)

- Plus on demande une précision fine (long( $I_m$ ) petite), plus  $\alpha$  sera grand (risque plus important)

- On peut alors jouer sur  $n$  (augmenter la taille de l'échantillon): si on veut une précision de  $\pm \epsilon$  (i.e. long( $I_m$ ) =  $2\epsilon$ ), choisit  $n \geq \left(\frac{q_{1-\frac{\alpha}{2}}}{\epsilon}\right)^2$ .

quelques valeurs numériques :  $\alpha = 0.01$  (niveau 99%)  $1 - \frac{\alpha}{2} = 0.995$   $q_{0.995} = 2.57$   
 $\alpha = 0.05$  (niveau 95%)  $1 - \frac{\alpha}{2} = 0.975$   $q_{0.975} = 1.96$   
 $\alpha = 0.1$  (niveau 90%)  $1 - \frac{\alpha}{2} = 0.95$   $q_{0.95} = 1.64$

b) cas  $\sigma^2$  inconnue: on utilise  $\hat{\Delta}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{m}_n)^2$  ou  $\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{m}_n)^2$ .

$$T_n \equiv \frac{\hat{m}_n - m}{\sqrt{\hat{\Delta}_n^2 / (n-1)}} = \frac{\hat{m}_n - m}{\sqrt{\hat{S}_n^2 / n}} : \text{Student } T(n-1) \text{ (même type de courbe que } N(0,1))$$

Soit  $t_{1-\frac{\alpha}{2}}(n-1)$  la  $(1-\frac{\alpha}{2})$ -quantile de  $T(n-1)$  :  $\mathbb{P}(T_n \leq t_{1-\frac{\alpha}{2}}(n-1)) = 1 - \frac{\alpha}{2}$   
 $\Rightarrow \mathbb{P}(|T_n| \leq t_{1-\frac{\alpha}{2}}(n-1)) = 1 - \alpha$ .

D'où un intervalle de confiance pour  $m$  de niveau  $1-\alpha$ :

$$I_m = \left[ \hat{m}_n - t_{1-\frac{\alpha}{2}}(n-1) \sqrt{\frac{\hat{\Delta}_n^2}{n-1}}, \hat{m}_n + t_{1-\frac{\alpha}{2}}(n-1) \sqrt{\frac{\hat{\Delta}_n^2}{n-1}} \right]$$

Pour  $n$  "grand" ( $n \geq 30$ ):  $T(n-1) \approx N(0,1)$  et  $t_{1-\frac{\alpha}{2}}(n-1) \approx q_{1-\frac{\alpha}{2}}$   $\rightarrow$  petite erreur supplémentaire

remarque: pour des échantillons sans remise il faut rajouter un coefficient d'exhaustivité  $\sqrt{1-f}$  ( $f = \frac{n}{N}$  proportion de l'échantillon) devant  $\sigma$  et  $\sqrt{\hat{\Delta}_n^2}$ .

## 2) Cas d'une population inconnue

a) si  $n$  est "grand" ( $n \geq 30$ ): approximation normale  $\frac{\hat{m}_n - m}{\sigma/\sqrt{n}} \xrightarrow{\mathcal{L}} N(0,1)$

$\rightarrow$  même type d'intervalle de confiance (avec erreur supplémentaire due à l'approximation)

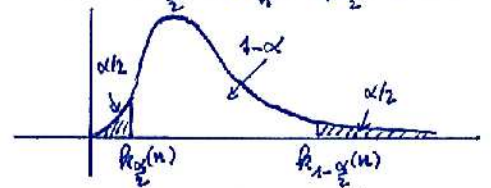
b) si  $n < 30$ : rien!

## 3) Intervalle de confiance pour la variance $\sigma^2$ d'une population $N(m, \sigma^2)$ .

a) cas  $m$  connue:  $\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ ,  $K_n \equiv \frac{\hat{\sigma}_n^2}{\sigma^2/n} = \sum_{i=1}^n \left(\frac{X_i - m}{\sigma}\right)^2 : \chi^2(n)$

Soit  $k_{\frac{\alpha}{2}}(n)$  et  $k_{1-\frac{\alpha}{2}}(n)$  les  $\frac{\alpha}{2}$  et  $(1-\frac{\alpha}{2})$ -quantiles de  $\chi^2(n)$  :  $\mathbb{P}(K_n \leq k_{\frac{\alpha}{2}}(n)) = \frac{\alpha}{2}$   
 $\mathbb{P}(K_n \leq k_{1-\frac{\alpha}{2}}(n)) = 1 - \frac{\alpha}{2}$   
 $\Rightarrow \mathbb{P}(k_{\frac{\alpha}{2}}(n) \leq K_n \leq k_{1-\frac{\alpha}{2}}(n)) = 1 - \alpha$ .

On a alors  $1 - \alpha = \mathbb{P}\left(\frac{\hat{\sigma}_n^2}{\sigma^2/n} \in [k_{\frac{\alpha}{2}}(n), k_{1-\frac{\alpha}{2}}(n)]\right)$   
 $= \mathbb{P}\left(\sigma^2 \in \left[\frac{n \hat{\sigma}_n^2}{k_{1-\frac{\alpha}{2}}(n)}, \frac{n \hat{\sigma}_n^2}{k_{\frac{\alpha}{2}}(n)}\right]\right)$   
 $= \mathbb{P}(\sigma^2 \in I_{\sigma^2})$



d'où un intervalle de confiance pour  $\sigma^2$  de niveau  $1-\alpha$  :  $I_{\sigma^2} = \left[ \frac{n \hat{\sigma}_n^2}{k_{1-\frac{\alpha}{2}}(n)}, \frac{n \hat{\sigma}_n^2}{k_{\frac{\alpha}{2}}(n)} \right]$

b) cas  $m$  inconnue: on utilise  $\hat{\Delta}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{m}_n)^2$  ou  $\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{m}_n)^2$

Dans ce cas  $K_{n-1} \equiv \frac{\hat{\Delta}_n^2}{\sigma^2/n} = \frac{\hat{S}_n^2}{\sigma^2/(n-1)} = \sum_{i=1}^n \left(\frac{X_i - \hat{m}_n}{\sigma}\right)^2 : \chi^2(n-1)$

$$I_{\sigma^2} = \left[ \frac{n \hat{\Delta}_n^2}{k_{1-\frac{\alpha}{2}}(n-1)}, \frac{n \hat{\Delta}_n^2}{k_{\frac{\alpha}{2}}(n-1)} \right]$$

#### 4) Estimation d'une proportion

Soit  $p = \frac{N_A}{N}$  la proportion d'individus possédant un caractère particulier A.

$$\hat{P}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{où } X_i = \mathbb{1}_{\{\text{individu a la prop. A}\}} \rightarrow \text{Bernoulli } B(p), \text{ n } \hat{P}_n: B(n, p).$$

$$E(\hat{P}_n) = p \text{ et } \text{var}(\hat{P}_n) = \frac{p(1-p)}{n}.$$

lorsque n est "grand",  $\frac{\hat{P}_n - p}{\sqrt{p(1-p)/n}} \approx N(0,1).$

On trouverait donc ainsi approximativement un intervalle de confiance pour p

$$I_p = \left[ \hat{P}_n - q_{1-\frac{\alpha}{2}} \sqrt{\frac{p(1-p)}{n}}, \hat{P}_n + q_{1-\frac{\alpha}{2}} \sqrt{\frac{p(1-p)}{n}} \right], \quad q_{1-\frac{\alpha}{2}}: (1-\frac{\alpha}{2})\text{-quantile de } N(0,1).$$

Malheureusement p est inconnue ! Exutoire:

a) estimer p par  $\hat{P}_n$ :  $I_p = \left[ \hat{P}_n - q_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{P}_n(1-\hat{P}_n)}{n}}, \hat{P}_n + q_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{P}_n(1-\hat{P}_n)}{n}} \right]$   
 → très approximatif!

b) intervalle de confiance par excès:  $p(1-p) \leq \frac{1}{4} \Rightarrow \sqrt{\frac{p(1-p)}{n}} \leq \frac{1}{2\sqrt{n}}$

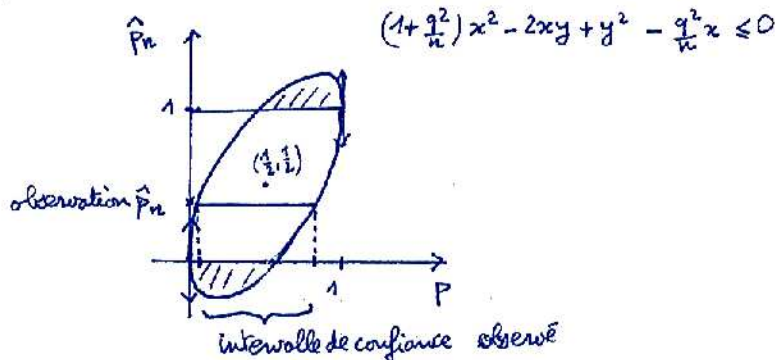
$$I_p = \left[ \hat{P}_n - \frac{q_{1-\frac{\alpha}{2}}}{2\sqrt{n}}, \hat{P}_n + \frac{q_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right] \text{ et } P(p \in I_p) \geq 1-\alpha.$$

intervalle un peu plus large (moins précis)

c) méthode de l'ellipse (valable aussi pour les petits échantillons):

$$p \in \left[ \hat{P}_n - q \sqrt{\frac{p(1-p)}{n}}, \hat{P}_n + q \sqrt{\frac{p(1-p)}{n}} \right] \Leftrightarrow |p - \hat{P}_n| \leq q \sqrt{\frac{p(1-p)}{n}} \quad \begin{matrix} (q \text{ mb donné}) \\ \text{si } n \text{ grand} \\ q = q_{1-\frac{\alpha}{2}} \end{matrix}$$

$$\Leftrightarrow (p, \hat{P}_n) \in \text{Ellipse } (x-y)^2 - \frac{q^2}{n} x(1-x) \leq 0$$



d) Utilisation de la fonction arcsinus: Soit  $Z = 2\sqrt{n} [\arcsin \sqrt{\hat{P}_n} - \arcsin \sqrt{p}]$ .

$$P(Z \leq z) = P(\hat{P}_n \leq \sin^2(\arcsin \sqrt{p} + z/2\sqrt{n})) = P(\hat{P}_n \leq [\sqrt{p} \cos(z/2\sqrt{n}) + \sqrt{1-p} \sin(z/2\sqrt{n})]^2)$$

$$= P(\hat{P}_n \leq p \cos^2(z/2\sqrt{n}) + (1-p) \sin^2(z/2\sqrt{n}) + \sqrt{p(1-p)} \sin(z/2\sqrt{n}))$$

$$= P\left(\frac{\hat{P}_n - p}{\sqrt{p(1-p)/n}} \leq \frac{\sqrt{n} (1-2p) \sin^2(z/2\sqrt{n})}{\sqrt{p(1-p)}} + \frac{\sqrt{n} \sin(z/2\sqrt{n})}{\sqrt{n}}\right) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} P(N < z) = \Phi(z).$$

d'où  $I_{\arcsin \sqrt{p}} = \left[ \arcsin \sqrt{\hat{P}_n} - \frac{q_{1-\frac{\alpha}{2}}}{2\sqrt{n}}, \arcsin \sqrt{\hat{P}_n} + \frac{q_{1-\frac{\alpha}{2}}}{2\sqrt{n}} \right] \sim z$

Bibliographie: Saporta: Probabilités, analyse des données et statistique 1990  
 Dreesenbeke: Eléments de statistique 1991  
 Fougereand et Fuchs: Statistique 1967  
 Aide-mémoire statistique 1995  
 Polycope INSA d'Angoumnd-Fauchon et al.

Complément : estimation par intervalle directe à partir d'une seule observation.

Théorème :

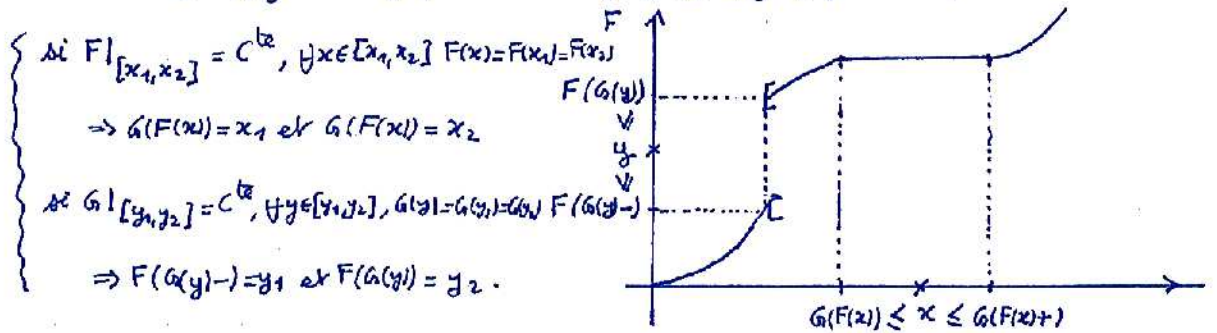
Soit  $F_\theta$  la f.r. de  $X$  :  $F_\theta(x) = P_\theta(X \leq x)$ ,  $F_\theta^-(x) = F_\theta(x-) = P_\theta(X < x)$ .  
 On suppose que  $\forall x, \theta \in (\theta_0, \theta_1) \mapsto F_\theta(x)$  est continue strictement décroissante (resp. croissante) ainsi que  $\theta \mapsto F_\theta^-(x)$ , et que  $\lim_{\theta \rightarrow \theta_0} F_\theta(x) = 1$  (resp.  $0$ )  
 Soit  $\alpha \in ]0, 1[$ . On a :

$$\forall x, \exists ! \theta_\alpha^+(x) \in (\theta_0, \theta_1) / F_{\theta_\alpha^+(x)}(x) = \alpha \text{ et } \exists ! \theta_{1-\alpha}^-(x) \in (\theta_0, \theta_1) / F_{\theta_{1-\alpha}^-(x)}^-(x) = 1 - \alpha.$$

Alors  $(\theta_0, \theta_\alpha^+(x)]$ ,  $[\theta_{1-\alpha}^-(x), \theta_1)$  et  $[\theta_{1-\frac{\alpha}{2}}^-(x), \theta_{\frac{\alpha}{2}}^+(x)]$  sont des intervalles de confiance de  $\theta$  de niveau  $\geq 1 - \alpha$ .

Démonstration : 1)  $P\{\theta \in (\theta_0, \theta_\alpha^+(x)]\} = P(F_\theta(x) \geq \alpha)$ .

Soit  $G_\theta$  l'inverse continue à gauche de  $F_\theta$  :  $G_\theta(y) = \inf\{x : F_\theta(x) \geq y\}$ .  
 On a  $F_\theta(x) \geq y \Leftrightarrow G_\theta(y) \leq x$  et  $F_\theta^-(G_\theta(y)) \leq y$  (cf dessin)



D'où  $P\{\theta \in (\theta_0, \theta_\alpha^+(x)]\} = P(X \geq G_\theta(\alpha)) = 1 - P(X < G_\theta(\alpha)) = 1 - F_\theta^-(G_\theta(\alpha)) \geq 1 - \alpha$ .

2) le 2<sup>e</sup> cas se traite avec  $Y = -X$ ,  $\varphi = -\theta$ ,  $G_\varphi(x) = P_\varphi(Y \leq x) = P_{-\varphi}(X \geq -x) = 1 - F_\theta^-(x)$ .

3)  $P(\theta \in [\theta_{1-\frac{\alpha}{2}}^-(x), \theta_{\frac{\alpha}{2}}^+(x)]) = P(A \cap B)$  avec  $A = \{\theta \in [\theta_{1-\frac{\alpha}{2}}^-(x), \theta_1)\}$ ,  $B = \{\theta \in (\theta_0, \theta_{\frac{\alpha}{2}}^+(x)]\}$   
 $= P(A) + P(B) - P(A \cup B)$   $P(A), P(B) \geq 1 - \frac{\alpha}{2}$   
 $\geq 1 - \alpha$  .  $\square$

Exemples : 1) loi binomiale  $B(n, \theta)$ .  $F_\theta(k) = \sum_{i=0}^k C_n^i \theta^i (1-\theta)^{n-i}$

$\theta \in [0, 1] \mapsto F_\theta(k)$  est continue strictement décroissante :  
 $\frac{d}{d\theta} F_\theta(k) = \sum_{i=0}^k C_n^i [i \theta^{i-1} (1-\theta)^{n-i} - (n-i) \theta^i (1-\theta)^{n-i-1}] = n! \left[ \sum_{i=1}^k \frac{\theta^{i-1} (1-\theta)^{n-i}}{(i-1)!(n-i)!} - \sum_{i=0}^k \frac{\theta^i (1-\theta)^{n-i-1}}{i!(n-i-1)!} \right]$   
 $= -\frac{n!}{k!(n-k-1)!} \theta^k (1-\theta)^{n-k-1} < 0$ .

Equation  $F_\theta(k) = \alpha$  difficile à résoudre!

2) Loi de Poisson  $P(\theta)$ .  $F_\theta(m) = \sum_{i=0}^m e^{-\theta} \frac{\theta^i}{i!} = 1 - \int_0^\theta \frac{x^n e^{-x}}{n!} dx$  (lien Poisson - Gamma)

$\theta \in [0, +\infty[ \mapsto F_\theta(m)$  continue strict. décroissante.  $F_{X^2(2n+2)}(2\theta) \rightarrow$  f.r.  $X^2$ .

$F_\theta(m) = \alpha \Leftrightarrow \int_0^\theta \frac{x^n}{n!} e^{-x} dx = 1 - \alpha \Leftrightarrow \theta = \theta_\alpha^+(m) = \frac{1}{2} \chi_{1-\alpha}^2(2n+2)$ .

$F_\theta^-(m) = 1 - \alpha \Leftrightarrow F_\theta(m-1) = 1 - \alpha \Leftrightarrow \theta = \theta_\alpha^-(m) = \frac{1}{2} \chi_\alpha^2(2n)$ .

D'où des intervalles de confiance pour  $\theta$  :  $[0, \frac{1}{2} \chi_{1-\alpha}^2(2X+2)]$ ,  $[\frac{1}{2} \chi_\alpha^2(2X), +\infty[$ ,  
 $[\frac{1}{2} \chi_{\frac{\alpha}{2}}^2(2X), \frac{1}{2} \chi_{1-\frac{\alpha}{2}}^2(2X+2)]$  de niveau  $\geq 1 - \alpha$ .

3) Loi géométrique  $G(\theta)$ .  $F_\theta(m) = \sum_{i=1}^m \theta(1-\theta)^{i-1} = 1 - (1-\theta)^m$

$\theta \in ]0, 1[ \mapsto F_\theta(m)$  continue strictement croissante.

$$F_\theta(m) = \alpha \Leftrightarrow \theta = \theta_\alpha^+(m) = 1 - (1-\alpha)^{\frac{1}{m}}$$

$$F_\theta(m) = 1-\alpha \Leftrightarrow F_\theta(m-1) = 1-\alpha \Leftrightarrow \theta = \theta_\alpha^-(m) = 1 - \alpha^{\frac{1}{m-1}}$$

D'où les intervalles de confiance :  $[0, 1 - \alpha^{\frac{1}{m-1}}]$ ,  $[1 - (1-\alpha)^{\frac{1}{m}}, 1]$ ,  $[1 - (1-\frac{\alpha}{2})^{\frac{1}{m}}, 1 - (\frac{\alpha}{2})^{\frac{1}{m-1}}]$ .

4) Loi exponentielle  $E(\theta)$ .  $F_\theta(x) = 1 - e^{-\theta x}$

$\theta \in [0, +\infty[ \mapsto F_\theta(x)$  continue strict. croissante.

$$F_\theta(x) = \alpha \Leftrightarrow \theta = \theta_\alpha^+(x) = -\frac{1}{x} \ln(1-\alpha)$$

D'où les intervalles de confiance :  $[0, -\frac{1}{x} \ln \alpha]$ ,  $[-\frac{1}{x} \ln(1-\alpha), +\infty[$ ,  $[-\frac{1}{x} \ln(1-\frac{\alpha}{2}), -\frac{1}{x} \ln \frac{\alpha}{2}]$ .

5) Loi normale  $N(\theta, \sigma^2)$ .  $F_\theta(x) = \int_{-\infty}^{x-\theta} e^{-\frac{y^2}{2\sigma^2}} \frac{dy}{\sqrt{2\pi}\sigma} = \int_{-\infty}^{\frac{x-\theta}{\sigma}} e^{-\frac{y^2}{2}} \frac{dy}{\sqrt{2\pi}}$

$\theta \in \mathbb{R} \mapsto F_\theta(x)$  continue strict. décroissante. Equation

$$F_\theta(x) = \alpha \Leftrightarrow \Phi\left(\frac{x-\theta}{\sigma}\right) = \alpha \Leftrightarrow \theta = \theta_\alpha^+(x) = x - \sigma q_\alpha = x + \sigma q_{1-\alpha}$$

D'où les intervalles de confiance :  $]-\infty, x + \sigma q_{1-\alpha}]$ ,  $[x - \sigma q_{1-\alpha}, +\infty[$ ,  $[x - \sigma q_{\frac{\alpha}{2}}, x + \sigma q_{\frac{\alpha}{2}}]$ .

Condition suffisante de monotonie de  $\theta \mapsto F_\theta(x)$  (rapport de vraisemblance monotone)

Soit  $f_{\theta\theta}$  une densité  $> 0$  de  $X$  par rapport à une mesure  $\mu$  ne dépendant pas de  $\theta$ , et  $F_\theta$  : f.r. de  $X$ . On suppose que  $\forall \theta, \theta', \theta < \theta' \Rightarrow (y \mapsto \frac{f_{\theta'}(y)}{f_\theta(y)})$  est strict. croissante (resp. décroissante).  
Alors  $\forall x$  tel que  $\forall \theta, 0 < F_\theta(x) < 1$ ,  $\theta \mapsto F_\theta(x)$  est strictement décroissante (resp. croissante).

démonstration: soit  $\theta < \theta'$ ,  $a = F_\theta(x)$ ,  $b = \frac{f_{\theta'}(x)}{f_\theta(x)}$ .

$$f_{\theta'}(y) - b f_\theta(y) = f_\theta(y) \left( \frac{f_{\theta'}(y)}{f_\theta(y)} - \frac{f_{\theta'}(x)}{f_\theta(x)} \right) \begin{cases} > 0 & \text{si } y > x \text{ (resp. } y < x) \\ < 0 & \text{si } y < x \text{ (resp. } y > x) \end{cases}$$

$$\begin{aligned} F_{\theta'}(x) - F_\theta(x) &= (F_{\theta'}(x) - b F_\theta(x)) + (b-1) F_\theta(x) \\ &= \int_{-\infty}^x [f_{\theta'}(y) - b f_\theta(y)] d\mu(y) - a \int_{-\infty}^{+\infty} [f_{\theta'}(y) - b f_\theta(y)] d\mu(y) \\ &= \underbrace{(1-\alpha)}_{>0} \underbrace{\int_{-\infty}^x [f_{\theta'}(y) - b f_\theta(y)] d\mu(y)}_{<0 \text{ (resp. } >0)} - \underbrace{\alpha}_{>0} \underbrace{\int_x^{+\infty} [f_{\theta'}(y) - b f_\theta(y)] d\mu(y)}_{>0 \text{ (resp. } <0)} \\ &< 0 \text{ (resp. } > 0). \quad \square \end{aligned}$$