

Convergence uniforme de la fonction de répartition dans le cas d'un plan de sondage informatif

D. BONNÉRY¹, F. J. BREIDT² and F. COQUET¹

1. CREST-ENSAI, Campus de Ker-Lann, France

2. Colorado State University, Fort Collins, USA

Ce travail a bénéficié d'un financement de la US National Science Foundation
(SES-0922142) et de l'Ensaï

09/11/2010

Plan

- Quel intérêt ?
- Résultat : convergence uniforme de la fonction de répartition.
- Exemples

Un échantillon est le résultat de deux processus aléatoires qui peuvent être dépendants :

- Un **modèle de superpopulation** génère des valeurs $(Y_{.k})_{k \in \llbracket 1, N \rrbracket}$ sur une population de taille N .
- Des individus sont sélectionnés via un plan de sondage aléatoire pour constituer un échantillon de taille aléatoire n .

Il y a deux façons naturelles de modéliser le premier processus :

- $(Y_{.k})_{k \in \llbracket 1, N \rrbracket}$ est constitué de N réalisations i.i.d. d'une même loi. (supposée appartenir à un modèle paramétrique dominé par une certaine mesure). Par la suite on supposera que $Y_{.k}$ suit une loi de densité f_θ par rapport à λ , la mesure de Lebesgue.
- $(Y_{.k})_{k \in \llbracket 1, N \rrbracket}$ est le réarrangement aléatoire d'un vecteur $(y_{.k})_{k \in \llbracket 1, N \rrbracket}$ qui joue le rôle de paramètre.

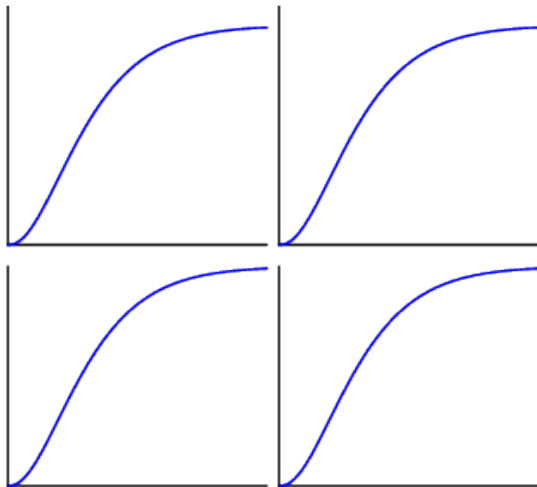
Ces deux approches se distinguent par les choix d'inférence statistique sous-jacents : sur $g(\theta)$?, $g(y_1 \dots y_N)$?.

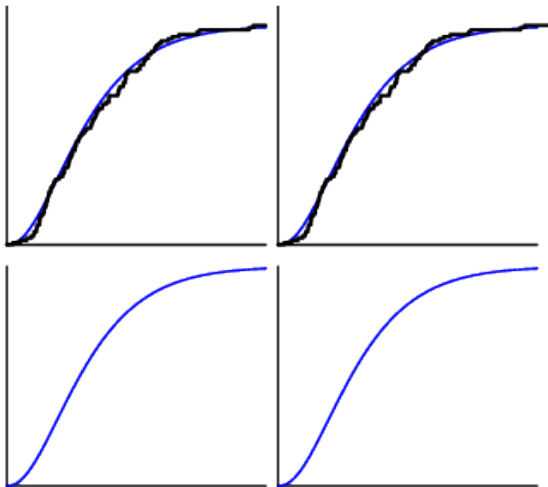
Un échantillon est un vecteur $(I_{.1}, \dots, I_{.N})$ où $I_{.k}$ représente le nombre de fois que l'individu k est sélectionné.

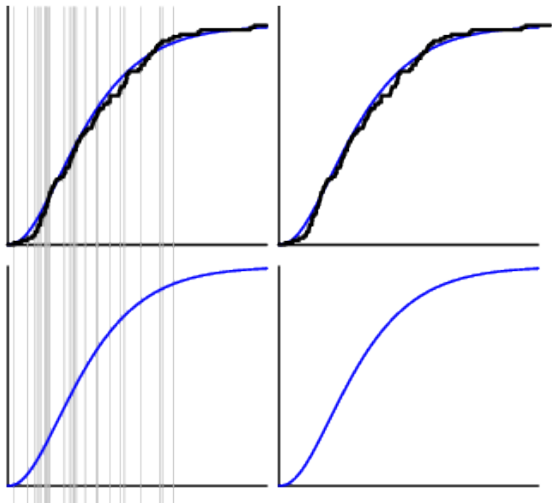
Le plan de sondage est une variable aléatoire $\Pi : \llbracket 1, N \rrbracket^{\mathbb{N}} \rightarrow [0, 1]$ qui donne la loi du vecteur $(I_{.1}, \dots, I_{.N})$ conditionnellement à la réalisation de certaines variables aléatoires auxiliaires qui peuvent être liées à $(Y_{.k})_{k \in \llbracket 1, N \rrbracket}$.

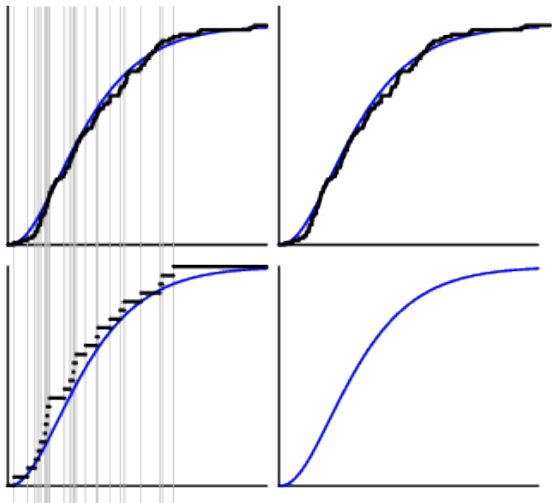
Le plan de sondage est **informatif** si

- $(I_{.1}, \dots, I_{.N})$ dépend de $(Y_{.k})_{k \in \llbracket 1, N \rrbracket}$
- On doit tenir compte de cette dépendance lors de l'inférence sur le modèle de superpopulation.



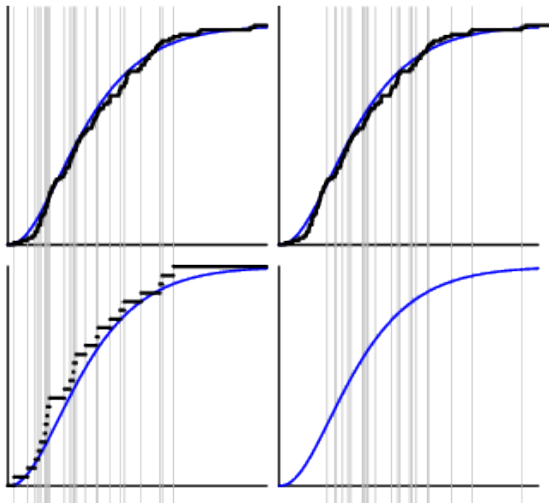






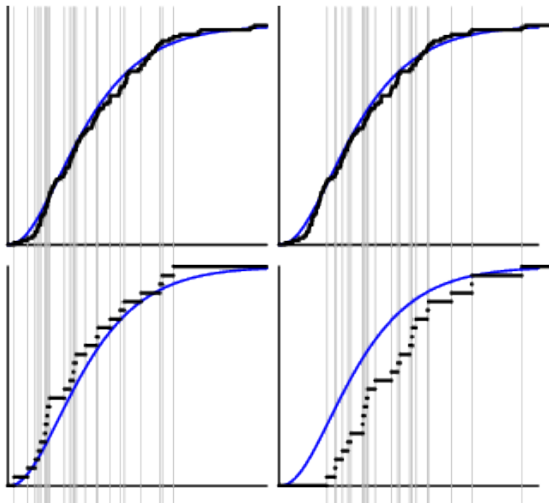
Non-Informative Selection

Informative Selection



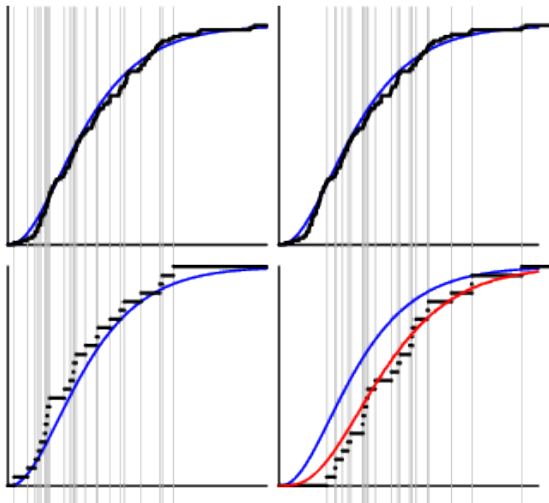
Non-Informative Selection

Informative Selection



Non-Informative Selection

Informative Selection



- Pfeiffermann et al. ont développé des méthodes d'estimation basées sur une approximation de la vraisemblance
 - comme produit de densités pondérées
 - sous de fortes hypothèses, et un cadre restrictif (n fixé, $N \rightarrow \infty$)

Nous voulons prendre du recul, pour commencer à répondre aux questions :

- Dans le cas d'un sondage informatif, la vraisemblance approchée est elle valable ?
- Quelles sont les propriétés des statistiques de bases liées à cette densité ?

$(N_\gamma)_{\gamma \in \mathbb{N}}$ est une suite croissante d'entiers, $(y_k)_{k \in \mathbb{N}}$ une suite de réels. Les variables aléatoires sont définies sur l'espace probabilisé (Ω, \mathcal{A}, P) .

$$\begin{array}{c|cccccccc}
 Y & Y_{.1} & \dots & Y_{.N_1} & \dots & Y_{.N_2} & \dots & Y_{.N_\gamma} & \dots \\
 I_1 & I_{11} & \dots & I_{1N_1} & & & & & \\
 I_2 & I_{21} & \dots & \dots & \dots & I_{2N_2} & & & \\
 \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & & \\
 I_\gamma & I_{\gamma 1} & \dots & I_{\gamma k} & \dots & \dots & \dots & I_{\gamma N_\gamma} &
 \end{array}$$

Pour $\gamma \in \mathbb{N}$,

- un échantillon est un vecteur $I_\gamma \in \mathbb{N}^{N_\gamma}$, $I_{\gamma k}$ est le nombre de fois que k est sélectionné.
- $Y_\gamma = (Y_{.k})_{k \in \llbracket 1, N_\gamma \rrbracket}$

On suppose l'échangeabilité par colonne : $\forall \gamma \in \mathbb{N}$, σ une permutation de $\llbracket 1, N_\gamma \rrbracket$, (I_γ, Y_γ) et $(\sigma.I_\gamma, \sigma.Y_\gamma)$ sont identiquement distribués.

Definition

Soit $\gamma \in \mathbb{N}$, la fonction de répartition empirique de l'échantillon est le processus aléatoire

$$F_\gamma : \omega \mapsto F_\gamma(\omega) : \mathbb{R} \rightarrow [0, 1]$$

$$\alpha \mapsto \frac{\sum_{k=1}^{N_\gamma} \mathbf{1}_{]-\infty, \alpha]}(Y_k) I_{\gamma k}}{n_\gamma + \mathbf{1}_{I_\gamma=0}}$$

On définit $m_\gamma : y \mapsto E[I_{\gamma k} | Y_{.k} = y]$

La densité d'une observation sélectionnée est la fonction : $w_\gamma f$, où

$w_\gamma(y) = \frac{m_\gamma(y)}{\int m_\gamma f d\lambda}$ et la fonction de répartition des observations

sélectionnées est $y \mapsto \int_{]-\infty, \alpha]} w_\gamma f d\lambda$

Propriété

$w_\gamma f$ est une densité de probabilité.

On note $Y_{\cdot(k)}$ le k ème individu tiré. Le vecteur des observations est $(Y_{\cdot(k)})_{k \in \llbracket 1, n_\gamma \rrbracket}$.

Propriété

- Dans le cas d'un tirage à taille n_γ fixe,

$$\forall k, P^{Y_{\cdot(k)}} = w_\gamma f \cdot \lambda$$

et $E[F_\gamma(\alpha)] = \int_{]-\infty, \alpha]} w_\gamma f d\lambda$

- Dans le cas d'un tirage sans remise :

$$P^{Y_{\cdot k} | I_{\gamma k} = 1} = w_\gamma f \cdot \lambda$$

Résultats : convergence uniforme de la fonction de répartition empirique

$$m_\gamma(y) = E[I_{\gamma k} | Y_{.k} = y]$$

Hypothèses

$$A_0 \begin{cases} \exists M \in L_1(\lambda) \text{ s.t. } \forall \gamma, m_\gamma < M \\ \forall y \in \mathbb{R}, m(y) = \lim_\gamma m_\gamma(y) \text{ exists} \\ \int m f d\lambda > 0 \end{cases}$$

Definition

Pour $\gamma \in \mathbb{N}$, $\alpha \in \mathbb{R}$, on définit la densité et fonction de répartition limites lorsque A_0 est vérifié, comme

$$w(\alpha) = \lim_\gamma w_\gamma = \lim_\gamma \frac{m_\gamma(\alpha)}{\int m_\gamma f d\lambda} = \frac{m(\alpha)}{\int m f d\lambda}$$
$$F_s(\alpha) = \frac{\int \mathbf{1}_{(-\infty, \alpha]} m f d\lambda}{\int m f d\lambda} = \int_{-\infty}^{\alpha} w f d\lambda$$

On définit:

$$m_\gamma(y) = E[I_{\gamma 1} | Y_{.1} = y]$$

$$m'_\gamma(y_1, y_2) = E[I_{\gamma 2} | Y_{.1} = y_1, Y_{.2} = y_2]$$

$$v_\gamma(y) = V[I_{\gamma 1} | Y_{.1} = y]$$

$$c_\gamma(y_1, y_2) = Cov[I_{\gamma 1}, I_{\gamma 2} | Y_{.k_1} = y_1, Y_{.2} = y_2]$$

Hypothèses

$$A_1 \left\{ \begin{array}{l} \int c_\gamma(y_1, y_2) f(y_1) f(y_2) dy_1 dy_2 = o_{\gamma \rightarrow \infty}(1) \\ \int (m'_\gamma(y_1, y_2) m'_\gamma(y_2, y_1) - m_\gamma(y_1) m_\gamma(y_2)) f(y_1) f(y_2) dy_1 dy_2 \\ \quad = o_{\gamma \rightarrow \infty}(1) \\ \int v_\gamma f d\lambda = o_{\gamma \rightarrow \infty}(N_\gamma), \quad \int m_\gamma^2 f d\lambda = o_{\gamma \rightarrow \infty}(N_\gamma) \\ P(\{I_\gamma = (0 \dots 0)\}) = o_\gamma(1) \end{array} \right.$$

On définit la version presque sûre des hypothèses précédentes :

Hypothèses

$$A_2 : \forall y \in \mathbb{R}^N \text{ t.q. } \sup_{\alpha} \left\| \frac{\sum_{k=1}^{N_{\gamma}} b_{\alpha}(y_k)}{N_{\gamma}} - \int b_{\alpha} f d\lambda \right\| = o_{\gamma}(1), \forall \alpha \in \mathbb{R},$$

$$\begin{cases} \sum_{k,l=1}^{N_{\gamma}} b_{\alpha}(y_k) b_{\alpha}(y_l) \text{Cov} [I_{\gamma k}, I_{\gamma l} | Y_{\gamma} = (y_1 \dots y_{N_{\gamma}})] = o_{\gamma}(N_{\gamma}^2) \\ \sum_{k=1}^{N_{\gamma}} b_{\alpha}(y_k) (E [I_{\gamma k} | Y_{\gamma} = (y_1 \dots y_{N_{\gamma}})] - m_{\gamma}(y_k)) = o_{\gamma}(N_{\gamma}) \\ g_{\gamma}(y, (0 \dots 0)) = o_{\gamma}(1) \end{cases}$$

Théorème

- Si A_0 et A_1 sont vérifiées, alors

$$\|F_\gamma - F_s\|_\infty \xrightarrow{\gamma \rightarrow \infty} 0$$

avec $\|F_\gamma - F_s\|_\infty = \sup_{\alpha \in \mathbb{R}} \{|F_\gamma(\alpha) - F_s(\alpha)|\}$

- Si A_0 et A_2 sont vérifiées, alors il existe une suite de variables aléatoires $(I'_{\gamma k})_{\gamma \in \mathbb{N}, k \in U_\gamma}$, $(Y'_k)_{k \in \mathbb{N}}$ définies sur $(\Omega \times [0, 1], \mathcal{A} \otimes \mathcal{B}_{[0,1]}, P' = P \otimes \lambda_{[0,1]})$ telles que
 - $\|F'_\gamma - F_s\|_\infty$ converge P' -p.s. vers 0
 - $\forall \gamma \in \mathbb{N}$, (I'_γ, Y'_γ) et (I_γ, Y_γ) ont la même loi
 - $\forall \gamma \in \mathbb{N}$, $\omega \in \Omega$, $x \in [0, 1]$, $Y'_\gamma(\omega, x) = Y_\gamma(\omega)$

avec

$$F'_\gamma : \mathbb{R} \rightarrow [0, 1], \alpha \mapsto F'_\gamma(\alpha) = \frac{\sum_{k \in U_\gamma} \mathbf{1}_{(-\infty, \alpha]}(Y'_k) I'_{\gamma k}}{\sum_{k \in U_\gamma} I'_{\gamma k} + \mathbf{1}_{I'_\gamma=0}}.$$

Exemples et simulations

Exemple de non convergence

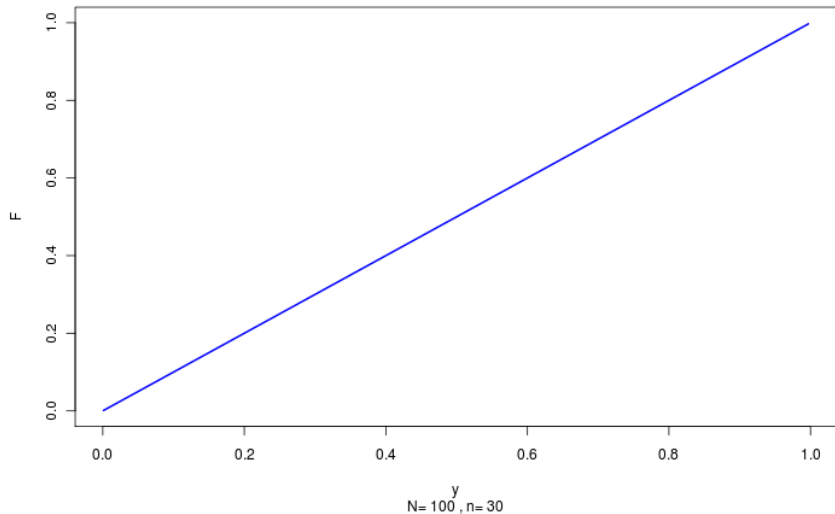
La limite F_s peut exister (A_0 est vérifiée), sans que F_γ ne converge (A_1, A_2 ne sont pas vérifiées):

On suppose que $Y \sim U_{[0,1]}$. On considère le plan de sondage suivant :

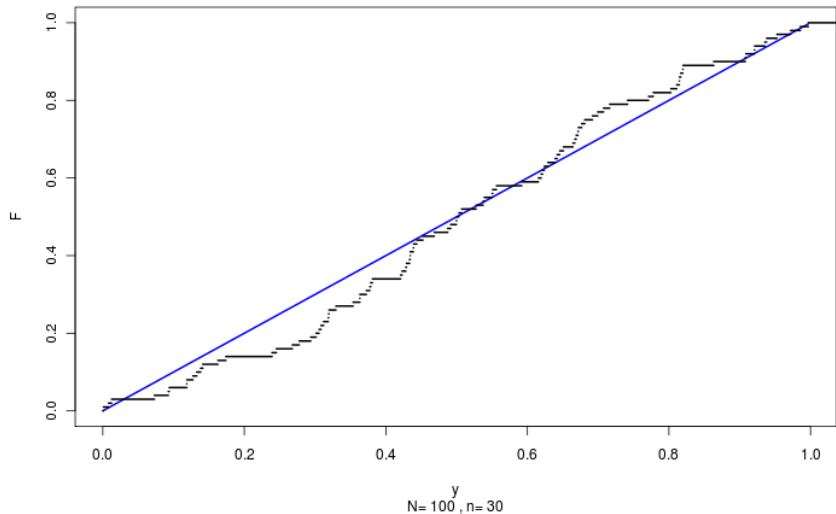
$$P(I_\gamma) \left| \begin{array}{cccccccccc} Y_{(1)} & \dots & Y_{(n_\gamma)} & Y_{(n_\gamma+1)} & \dots & Y_{(N_\gamma-n_\gamma)} & Y_{(N_\gamma-n_\gamma+1)} & \dots & Y_{(N_\gamma)} \\ \frac{1}{2} & & & & & & & & & & \\ \frac{1}{2} & & & & & & & & & & \\ \frac{1}{2} & & & & & & & & & & \end{array} \right. \begin{array}{cccccccccc} 1 & \dots & 1 & 0 & \dots & & & \dots & & & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 & & & 1 & \dots & 1 \end{array}$$

Dans ce cas, A_1, A_2 ne sont pas vérifiées.

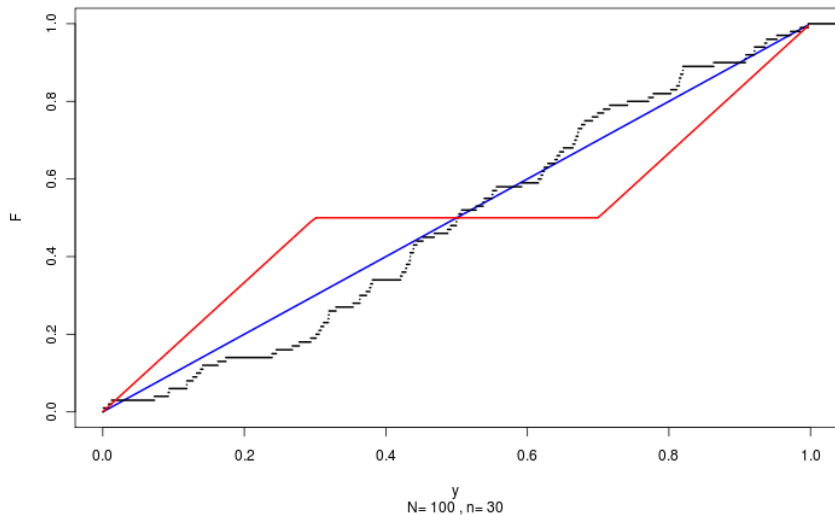
Non convergence



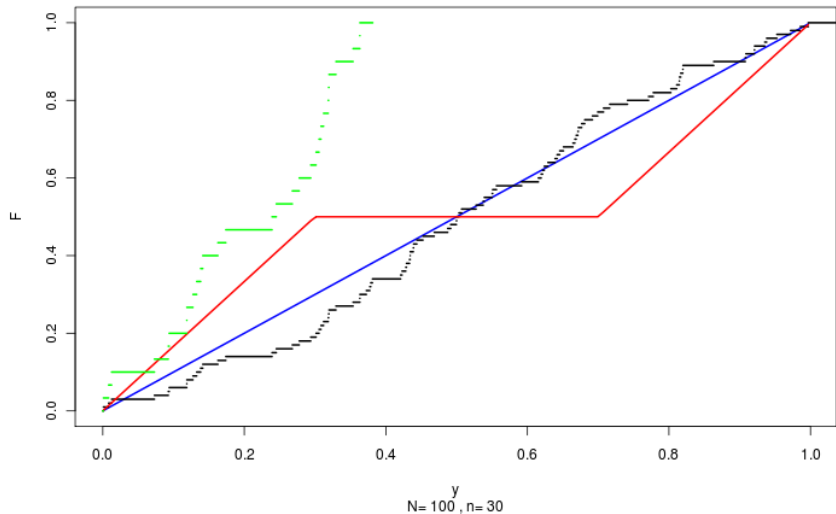
Non convergence



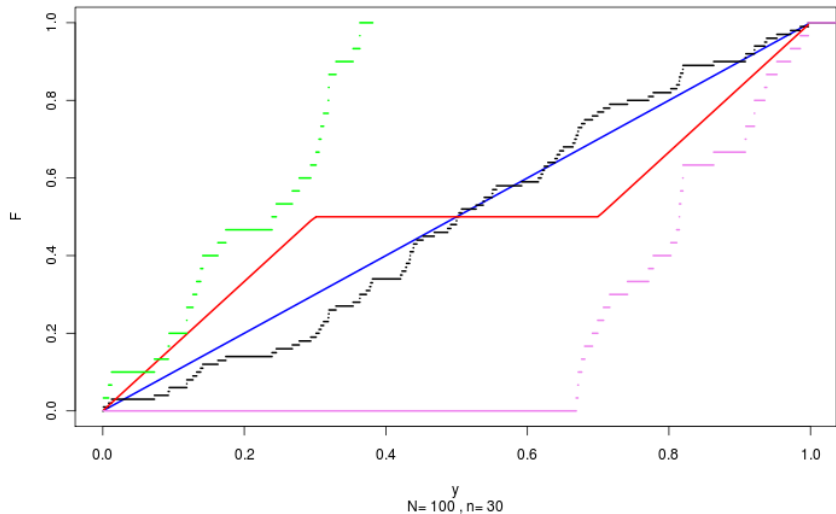
Non convergence



Non convergence



Non convergence

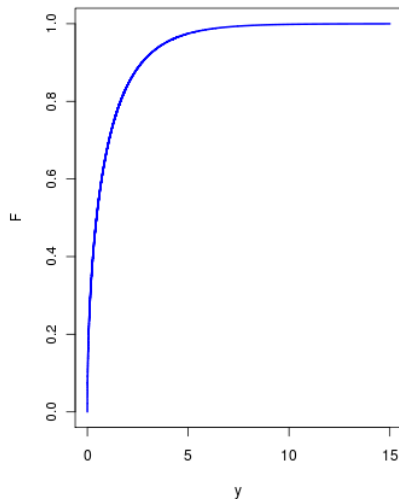


Sondage avec remise et probabilités proportionnelles à la taille

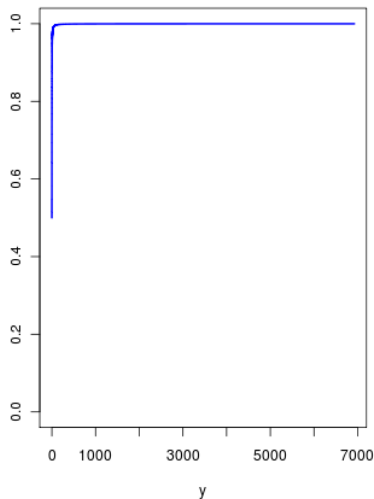
- On considère un tirage avec remise de taille n_γ avec tirages indépendants et probabilités proportionnelles à la taille :

$$p_\gamma = \frac{Y_\gamma}{\sum_{k=1}^{N_\gamma} Y_{\gamma k}}$$
 Alors A_0, A_1, A_2 sont vérifiées si $\lim \frac{n_\gamma}{N_\gamma} > 0$
 et $Y \in L_6$
- Sans ces conditions sur les moments, A_0, A_1, A_2 peuvent ne pas être vérifiées.

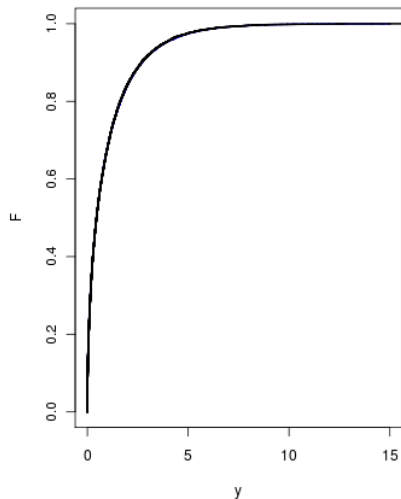
Chi2 N= 10000 , n= 300



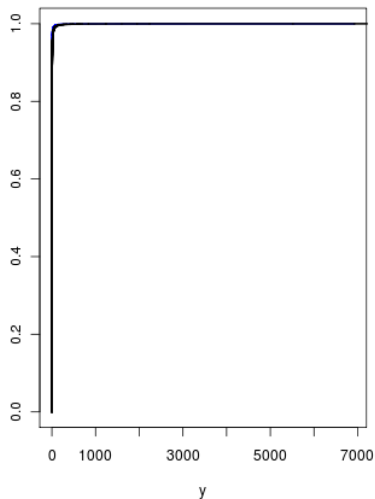
Cauchy N= 10000 , n= 300



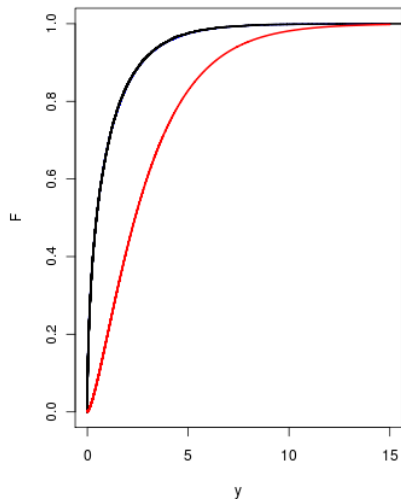
Chi2 N= 10000 , n= 300



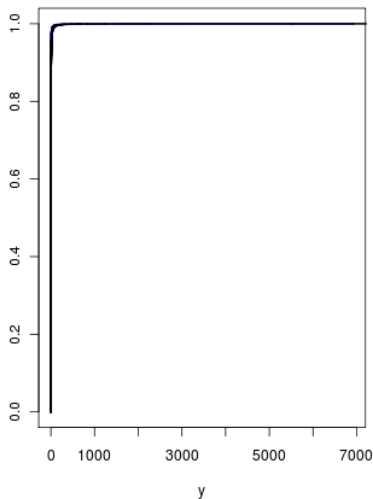
Cauchy N= 10000 , n= 300



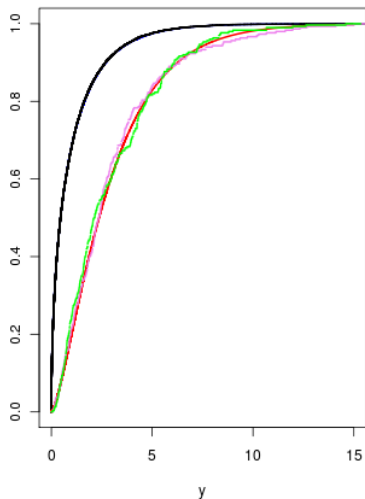
Chi2 N= 10000 , n= 300



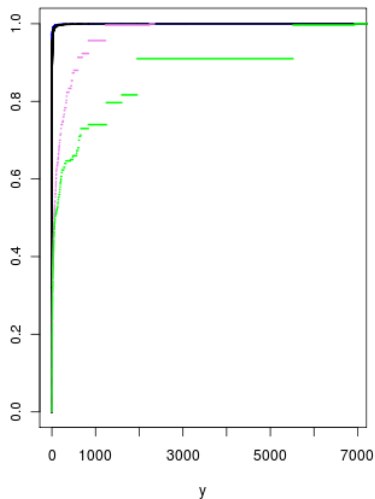
Cauchy N= 10000 , n= 300



Chi2 N= 10000 , n= 300



Cauchy N= 10000 , n= 300



Stratification basée sur Y

Pour $\gamma \in \mathbb{N}$, la population est stratifiée par tranches de valeurs de Y en H_γ strates. La strate h de taille $N_{\gamma h}$ est

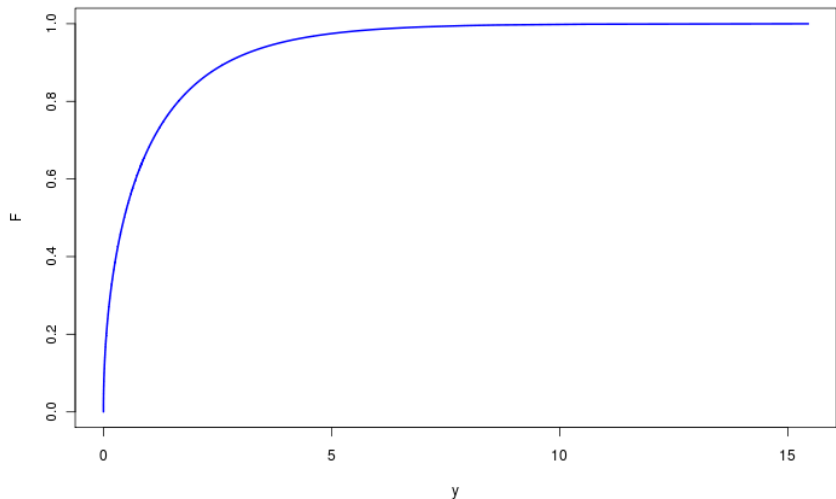
$$\{k \in \llbracket 1, N \rrbracket \mid Y_{(\sum_{h' < h} N_{\gamma h'})} < Y_k \leq Y_{(\sum_{h' \leq h} N_{\gamma h'})}\}$$

Si $\lim_{\gamma} \sum_{h=1}^{H_\gamma} \frac{n_{\gamma h}}{N_{\gamma h}} \mathbb{1}_{\left[\frac{\sum_{h' < h} N_{\gamma h'}}{N_\gamma}, \frac{\sum_{h' \leq h} N_{\gamma h'}}{N_\gamma} \right]}$ exists, alors A_0, A_1, A_2

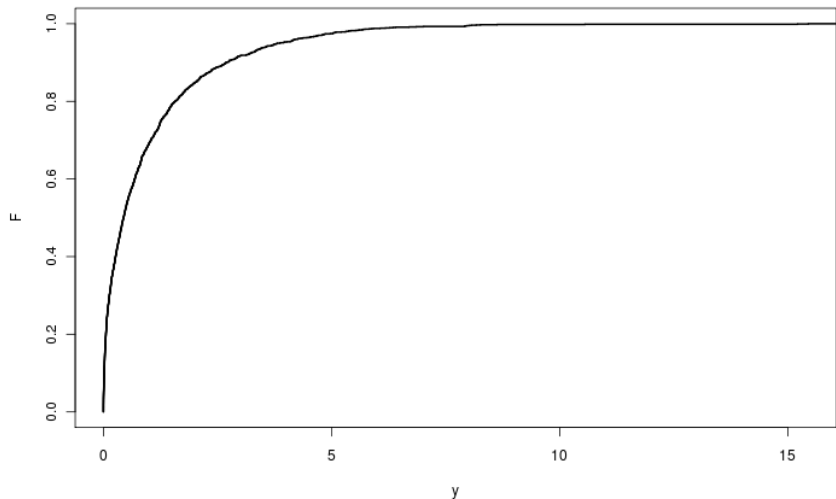
sont vérifiées.

On considère désormais un sondage stratifié avec une allocation $n_{\gamma h}$ croissante avec h .

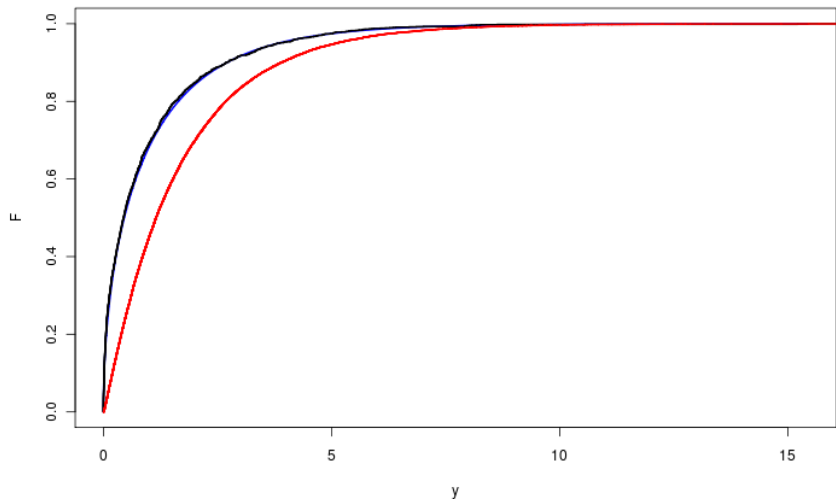
Stratified sampling - Chi2 N= 787 , n= 188



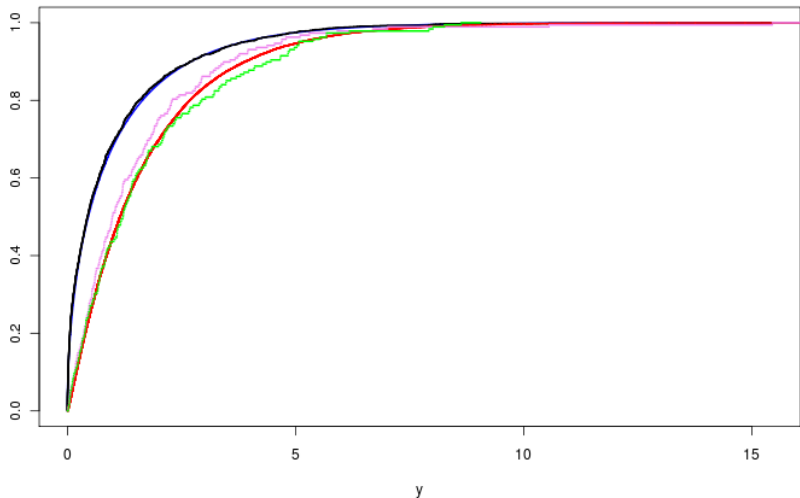
Stratified sampling - Chi2 N= 787 , n= 188



Stratified sampling - Chi2 N= 787 , n= 188



Stratified sampling - Chi2 N= 787 , n= 188






Résumé et discussion

Résumé:

- Conditions simples et vérifiables sur la suite de plans de sondage.
- Résultats de convergence uniforme de la fonction de répartition.
- Exemples

Travaux en cours :

- Autres propriétés asymptotiques de statistiques du même type dans le même cadre.
- Normalité asymptotique.

-  Danny Pfeffermann, Abba M. Krieger, and Yosef Rinott.
Parametric distributions of complex survey data under
informative probability sampling.
Statist. Sinica, 8(4):1087–1114, 1998.
-  Donald B. Rubin.
Inference and missing data.
Biometrika, 63(3):581–592, 1976.
-  R. A. Sugden and T. M. F. Smith.
Ignorable and informative designs in survey sampling inference.
Biometrika, 71(3):495–506, 1984.