



Une Introduction au Bootstrap

Patrice Bertail

CREST, Laboratoire de Statistique
Université Paris X, Nanterre

- 1) Idées de bases de la méthode (Origine, algorithme, Intervalles de confiance, propriétés)
- 2) Extensions dans le cas dépendant

17 juin 2003

Patrice.Bertail@ensae.fr



Les origines du Bootstrap

To pull oneself up with its own Bootstrap :

Se soulever soit même en tirant sur ses lacets
: Efron (1979)(1981).

Idée très ancienne : utilisation répétée des données
(sous-échantillons) pour estimer la variabilité d'une
statistique (Inde, 30's, Jackknife, 50's).

Méthode de calcul intensif fortement liée au
développement de l'informatique

Rééchantillonnage

Deux idées : -Estimation non-paramétrique (toute l'information est dans les données) P_n

-Algorithme de Monte-Carlo

- $(x_1, x_2, x_3, \dots, x_n)$ données i.i.d. Loi F

Ex.

(2.2, 3.2, 3.4, 2.5, 5.2, 3.3, 3.0, 3.7, 4.1, 3.1)

Statistique d'intérêt

$$T_n(x_1, \dots, x_n) = n^{-1} \sum_{i=1}^n x_i \underset{\text{Exemple}}{=} 3.37$$



Algorithme

Données :

(2.2, 3.2, 3.4, 2.5, 5.2, 3.3, 3.0, 3.7, 4.1, 3.1)

■ Répéter B fois (B grand=999)

Tirage uniforme avec remise (poids $1/n$)

dans les données $(x_1^*, x_2^*, x_3^*, \dots, x_n^*)$

Exemple = (5.2, 3.1, 3.1, 3.2, 2.5, 2.5, 3.7, 3.3, 3.7, 2.5)

Calcul de la valeur de la statistique

$$T_n^{*(1)}(x_1^*, x_2^*, x_3^*, \dots, x_n^*) \underset{\text{Exemple}}{=} 3.28$$

■ Fin



Variance et distribution Bootstrap

Résultats : Collections de B valeurs $T_n^{*(i)}$, $i = 1, \dots, B$

3.17, 3.40, 3.28, 3.27, 3.37,,

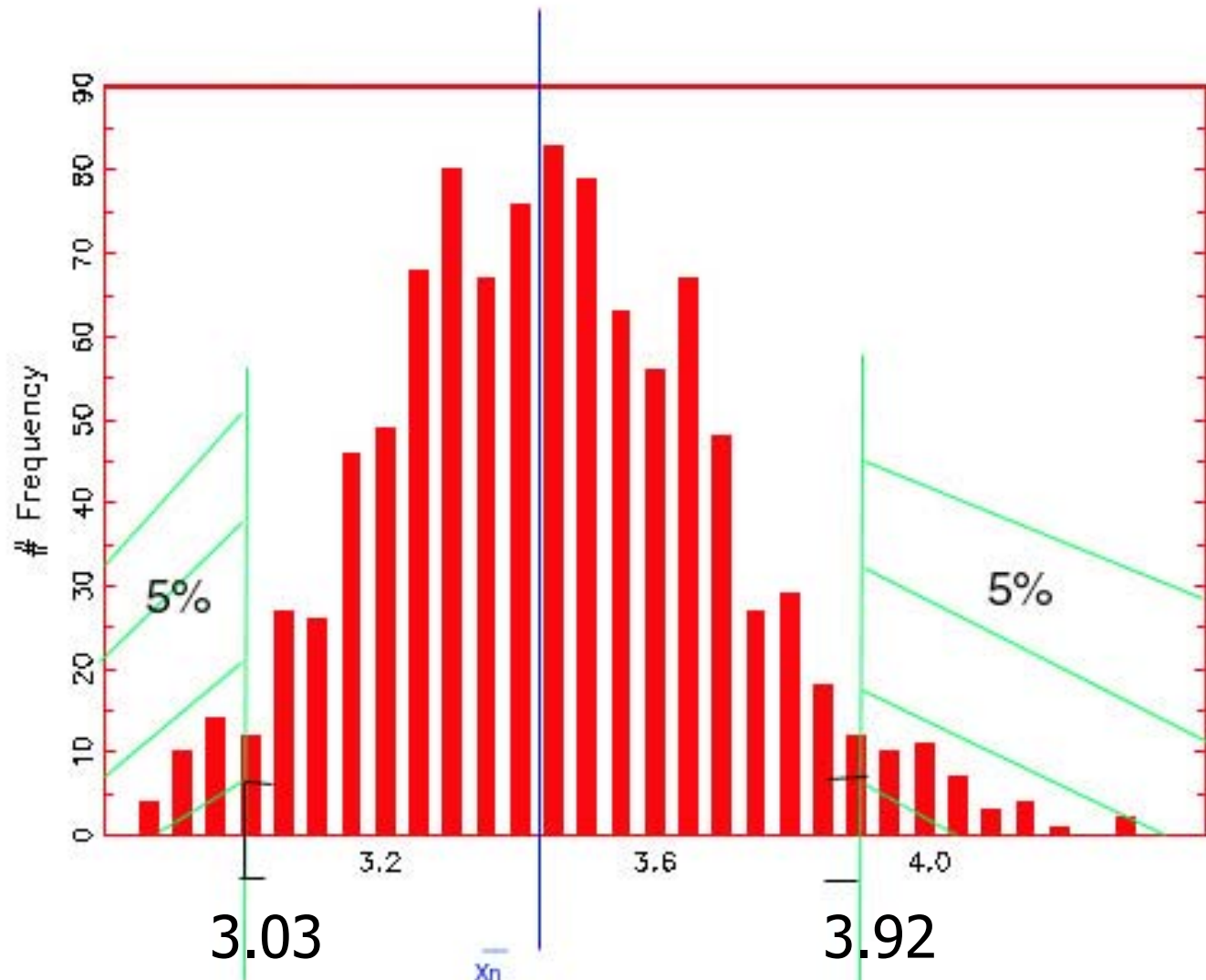
1er type d'utilisation : La variance empirique de ces valeurs
Est un estimateur de la volatilité de la statistique de départ

Exemple : 0.073 (B=999), Calcul direct 0.070

2ème type d'utilisation : La distribution empirique de ces
valeurs est un estimateur de la vraie distribution

→ Construction d'intervalles de confiance

Intervalle de confiance Bootstrap



Intervalle de confiance à 90%



Propriétés

- Facilité d'utilisation (calcul parallélisable)
- Mêmes propriétés que les intervalles de confiance asymptotique usuels.
- Possibilité d'avoir des améliorations spectaculaires en travaillant avec des statistiques standardisées (t-percentile).
- Se généralise aux données multidimensionnelles, censurées, et à toutes statistiques « lisses » (fonctionnelles Hadamard différentiables).

Variantes du Bootstrap

Bootstrap paramétrique

Si un modèle pour les données est disponible, F_θ

Ex : modèle gaussien $N(m, \sigma^2)$, $\theta = (m, \sigma^2)$

1ère étape : Estimation préalable du modèle : $\hat{\theta}_n$

2ème étape : Rééchantillonnage dans $F_{\hat{\theta}_n}$

Ex: génération de données dans $N(3.7, 0.07)$

Bootstrap lissé ou perturbé

Tirage dans les données auxquelles on rajoute une perturbation
Gaussienne $N(0, n^{-2})$ -> Meilleur pour des fractiles (VaR)

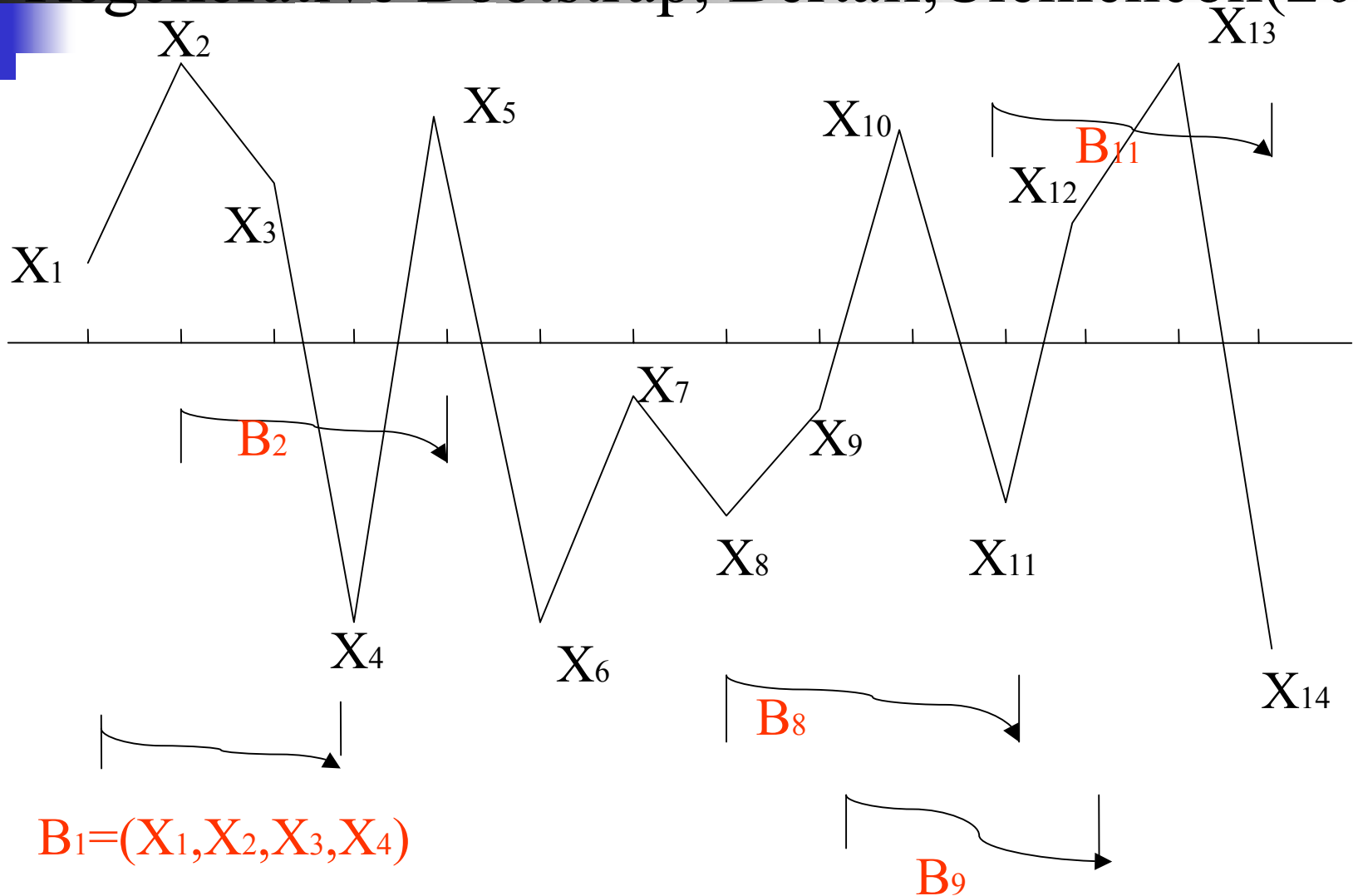


Cas non indépendant

- Essayer de se ramener au cas indépendant -
- soit en rééchantillonnant des blocs d'observations (gardant la dépendance)
 - soit en utilisant la structure du modèle (résidus i.i.d) : Ex cas de la régression linéaire, cas des modèles AR, ARMA en séries temporelles

Moving-Block Bootstrap, Kunsch(1989)

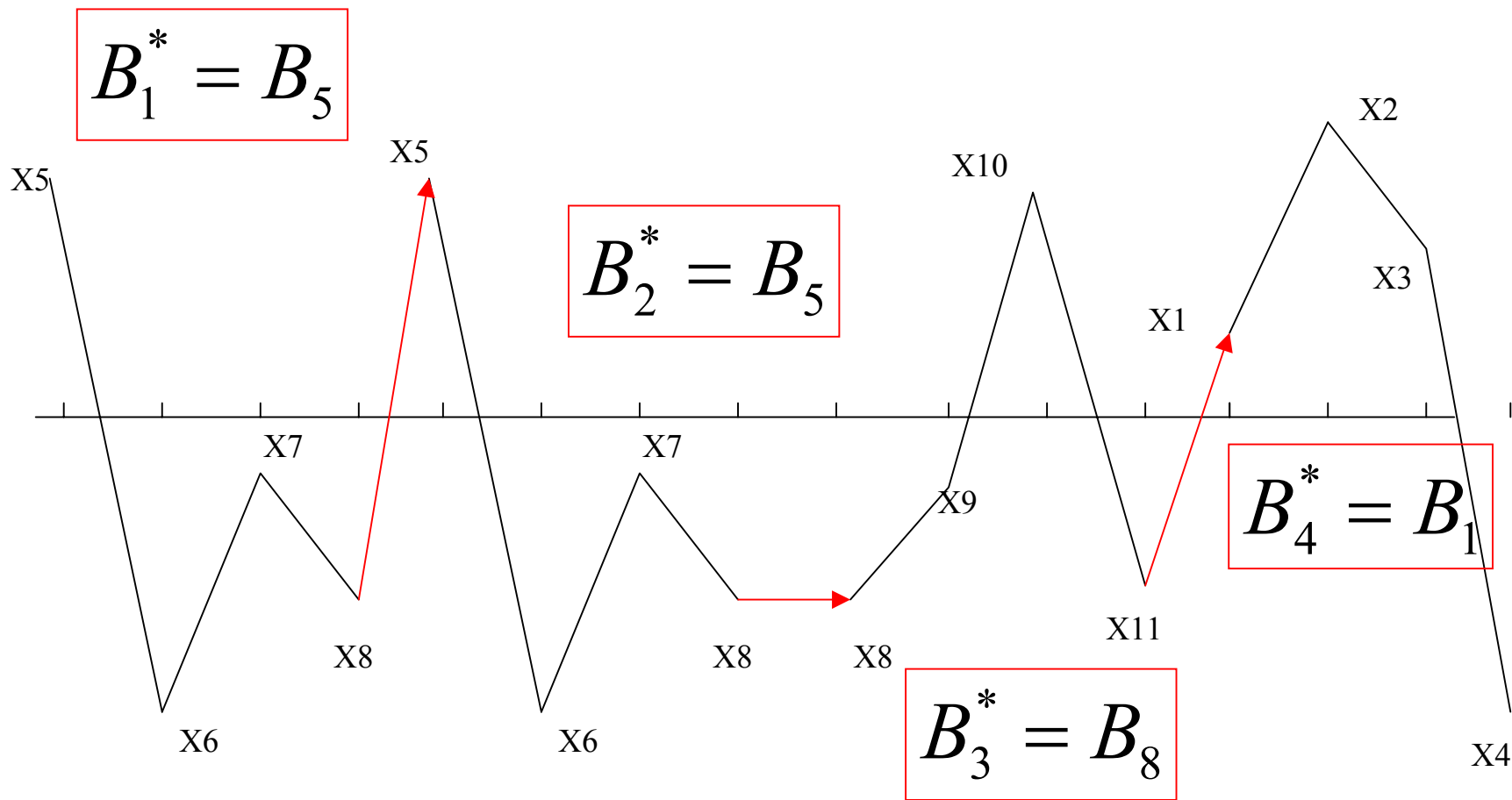
Regenerative Bootstrap, Bertail, Cléménçon(2003)



Longueur des blocs : $o(\sqrt{n})$

$$(B_1^*, B_2^*, B_3^*, \dots, B_{[n/b_n]}^*) \text{ i.i.d } F_B = \frac{1}{(n - b_n + 1)} \sum_{i=1}^{n - b_n + 1} \delta_{\cdot}$$

Rééchantillonnage de bloc : Reconstruction de séries artificielles



Bootstrap semi-paramétrique

Modèle autorégressif

$$X_0 = x_0, \quad V(\varepsilon_t) = \sigma^2, \quad c = x_0/\sigma$$

$$X_t = X_{t-1}\rho + \varepsilon_t \quad t = 1, \dots, T$$

Estimation du modèle

$\widehat{\rho}_T$: estimateur des m.c.o

$$\widehat{\varepsilon}_t = X_t - \widehat{\rho}_T X_{t-1} \quad \widetilde{\varepsilon}_t = \widehat{\varepsilon}_t - \frac{1}{T} \sum_{t=1}^T \widehat{\varepsilon}_t$$

Rééchantillonnage des résidus

Répéter B fois

$(\varepsilon_1^*, \varepsilon_2^*, \dots, \varepsilon_T^*)$ tirés avec remise dans
 $(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_T)$

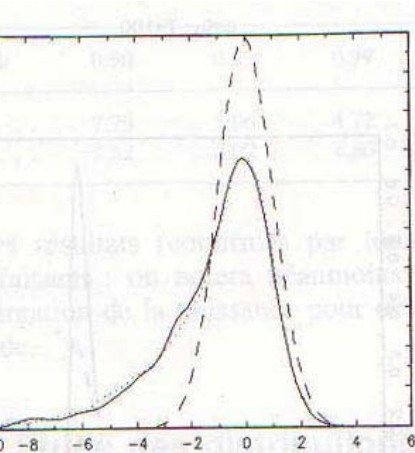
Reconstruction des X_t^* :

$$X_0^* = x_0,$$

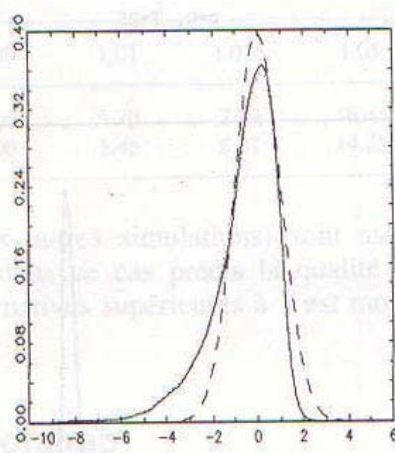
$$X_t^* = X_{t-1}^* \widehat{\rho}_T + \varepsilon_t^* \quad t = 1, \dots, T$$

Recalculer $\widehat{\rho}_T^*$: estimateur des m.c.o

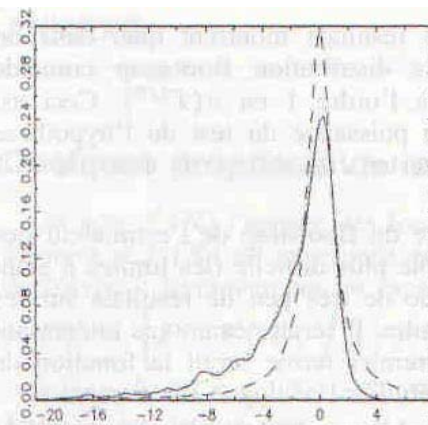
Fin



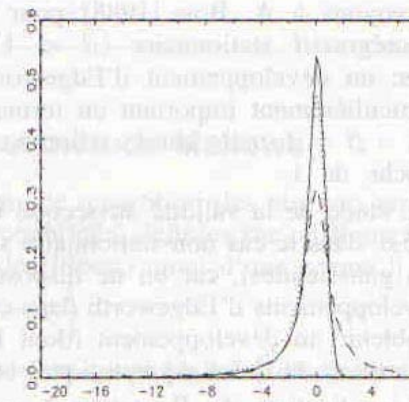
$c=0, T=25$



$c=0, T=100$



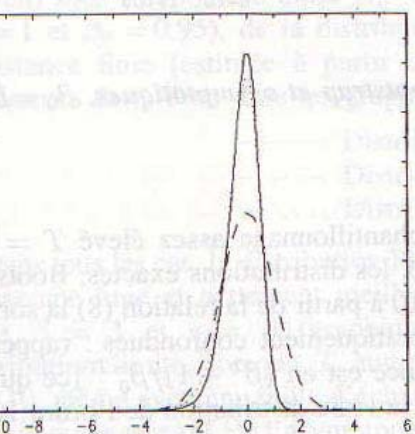
$c=10, T=100$



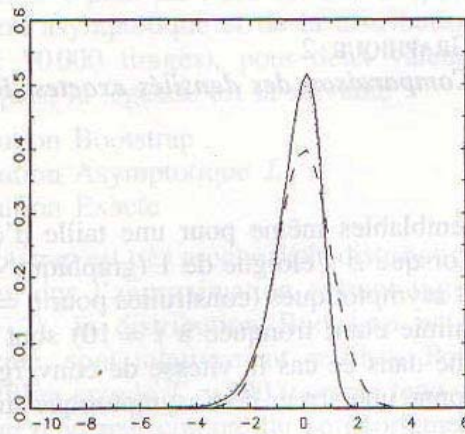
$c=10, T=100$

GRAPHIQUE 4

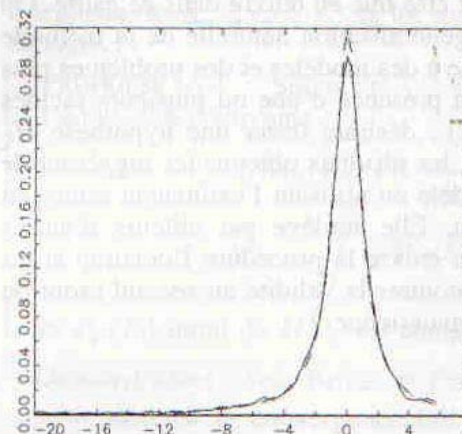
Comparaison des densités exactes, Bootstrap et asymptotiques, $\beta_0 = 1.$



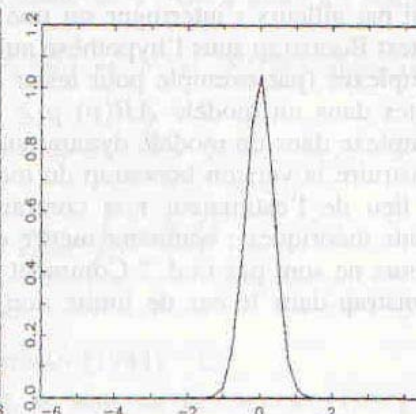
$c=10, T=25$



$c=10, T=100$



$c=0, T=100$



$c=10, T=100$

GRAPHIQUE

Comparaison des densités exactes, Bootstrap et asymptotiques, $\beta_0 = 1.$



Bibliographie succincte

■ Livres introductifs

- . Efron (1981), SIAM
- . Efron, Tibshirani (1997), Chapman, Hall
- . Davison, Hinkley (2000),

■ Livres théoriques

- . Hall (1992), Springer
- . Shao, Tu (1995), Springer (biblio)
- . Barbe, Bertail (1995), Springer
- . Giné (1996), in Ecole Proba St Flour, Springer

Sous-échantillonnage

Politis, Romano, Wolf (2000), Springer