

Bootstrap módszerek és alkalmazásuk összefüggő adatsorokra

Varga László

Valószínűségelméleti és Statisztika Tanszék
Természettudományi Kar
Eötvös Loránd Tudományegyetem

TÁMOP Kutatószeminárium
2010. november 19.

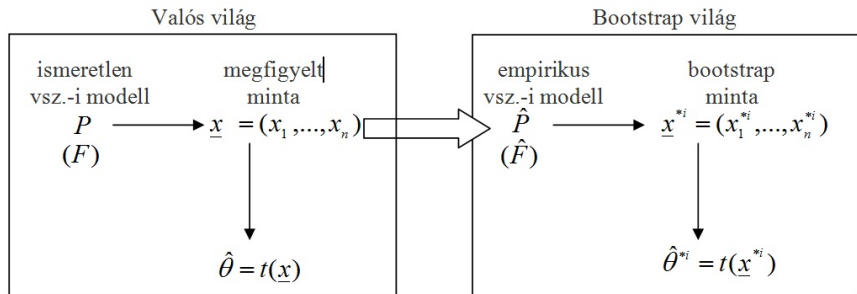
- 1 Bevezetés
- 2 Bootstrap módszerek
 - Az i.i.d. bootstrap
 - Blokk bootstrap módszerek
 - A blokkméret kiválasztásának problémája
- 3 Alkalmazás egydimenziós szélességi adatokra
- 4 További kiterjesztési lehetőségek
- 5 Goodness-of-fit tesztek
- 6 Összefoglalás

- Cél: adott x_1, \dots, x_n mintából minél több információt kisajtolni
- Jackknife (~50-es évek)
 - 1-törléses jackknife:
i-edik jackknife minta: $\mathbf{x}_i = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$ $i = 1, \dots, n$
 - d-törléses jackknife:
 $\binom{n}{d}$ mintát csinálunk d mintaelem elhagyásával
- Bootstrap (Efron 1979.)
 $\mathbf{x}_i^* = \{x_1^*, \dots, x_m^*\}$ visszatevéses mintavétellel az eredeti mintából általában $m=n$

- Cél: adott x_1, \dots, x_n mintából minél több információt kisajtolni
- Jackknife (~50-es évek)
 - 1-törléses jackknife:
i-edik jackknife minta: $\mathbf{x}_i = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$ $i = 1, \dots, n$
 - d-törléses jackknife:
 $\binom{n}{d}$ mintát csinálunk d mintaelem elhagyásával
- Bootstrap (Efron 1979.)
 $\mathbf{x}_i^* = \{x_1^*, \dots, x_m^*\}$ visszatevéses mintavétellel az eredeti mintából általában $m=n$

- Cél: adott x_1, \dots, x_n mintából minél több információt kisajtolni
- Jackknife (~50-es évek)
 - 1-törléses jackknife:
i-edik jackknife minta: $\mathbf{x}_i = \{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}$ $i = 1, \dots, n$
 - d-törléses jackknife:
 $\binom{n}{d}$ mintát csinálunk d mintaelem elhagyásával
- Bootstrap (Efron 1979.)
 $\mathbf{x}_i^* = \{x_1^*, \dots, x_m^*\}$ visszatevéses mintavétellel az eredeti mintából általában $m=n$

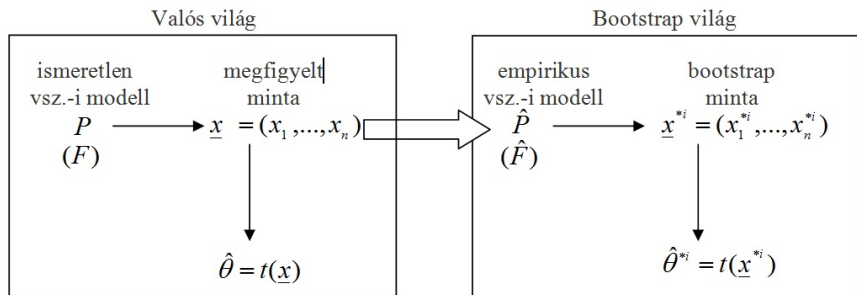
- A bootstrap módszerek lényege



- Nehézségek a gyakorlatban:

- 1 $\underline{x} \implies \hat{P}$ minden modellnél más és más
- 2 $\hat{P} \implies \underline{x}^*$ a sok ismétlés megterheli a számítógépet

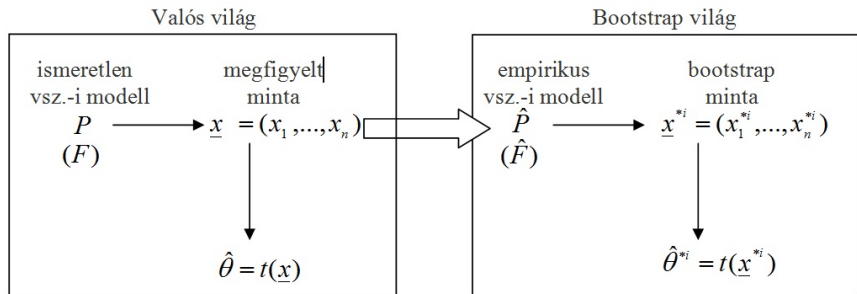
- A bootstrap módszerek lényege



- Nehézségek a gyakorlatban:

- 1 $\underline{x} \implies \hat{P}$ minden modellnél más és más
- 2 $\hat{P} \longrightarrow \underline{x}^*$ a sok ismétlés megterheli a számítógépet

- A bootstrap módszerek lényege



- Nehézségek a gyakorlatban:

- 1 $\underline{x} \implies \hat{P}$ minden modellnél más és más
- 2 $\hat{P} \implies \underline{x}^*$ a sok ismétlés megterheli a számítógépet

- 1 Bevezetés
- 2 **Bootstrap módszerek**
 - **Az i.i.d. bootstrap**
 - Blokk bootstrap módszerek
 - A blokkméret kiválasztásának problémája
- 3 Alkalmazás egydimenziós szélességi adatokra
- 4 További kiterjesztési lehetőségek
- 5 Goodness-of-fit tesztek
- 6 Összefoglalás

- Legyenek X_1, X_2, \dots i.i.d. valószínűségi változók, F (ismeretlen) közös eloszlással
 - $\mathcal{X}_n = \{X_1, \dots, X_n\}$ véletlen minta
 - $T_n = t_n(\mathcal{X}_n; F)$ minket érdeklő val.változó, az eloszlása: G_n
- Cél: G_n eloszlásának becslése
- Bootstrap módszer:
 - Adott \mathcal{X} -re, visszatevéssel m elemű mintát veszünk:
 $\mathcal{X}_m^* = \{X_1^*, \dots, X_m^*\}$ (általában $m \approx n$)
 - az X_i^* -ok közös eloszlása: $F_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$
 - $T_{m,n}^* = t_m(\mathcal{X}_m^*; F_n)$
 - Ismétlések $\Rightarrow \hat{G}_{m,n}$

- Legyenek X_1, X_2, \dots i.i.d. valószínűségi változók, F (ismeretlen) közös eloszlással
 - $\mathcal{X}_n = \{X_1, \dots, X_n\}$ véletlen minta
 - $T_n = t_n(\mathcal{X}_n; F)$ minket érdeklő val.változó, az eloszlása: G_n
- Cél: G_n eloszlásának becslése
- Bootstrap módszer:
 - Adott \mathcal{X} -re, visszatevéssel m elemű mintát veszünk:
 $\mathcal{X}_m^* = \{X_1^*, \dots, X_m^*\}$ (általában $m \approx n$)
 - az X_i^* -ok közös eloszlása: $F_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$
 - $T_{m,n}^* = t_m(\mathcal{X}_m^*; F_n)$
 - Ismétlések $\Rightarrow \hat{G}_{m,n}$

- Legyenek X_1, X_2, \dots i.i.d. valószínűségi változók, F (ismeretlen) közös eloszlással
 - $\mathcal{X}_n = \{X_1, \dots, X_n\}$ véletlen minta
 - $T_n = t_n(\mathcal{X}_n; F)$ minket érdeklő val.változó, az eloszlása: G_n
- Cél: G_n eloszlásának becslése
- Bootstrap módszer:
 - Adott \mathcal{X} -re, visszatevéssel m elemű mintát veszünk:
 $\mathcal{X}_m^* = \{X_1^*, \dots, X_m^*\}$ (általában $m \approx n$)
 - az X_i^* -ok közös eloszlása: $F_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$
 - $T_{m,n}^* = t_m(\mathcal{X}_m^*; F_n)$
 - Ismétlések $\Rightarrow \hat{G}_{m,n}$

Az i.i.d. bootstrap előnyei és korlátai

- Az i.i.d. bootstrap előnyei a jackknife-hoz képest:
 - 1 Intervallumbecslésre és statisztika eloszlásának becslésére is használható
 - 2 "Bonyolultabb" statisztikák becslésére is alkalmas
- Az i.i.d. bootstrap korlátai:
 - 1 számításigényes
 - 2 bizonyos esetekben a becslés nem lesz konzisztens

Példa (Singh, 1981)

Def: $\{X_n\}_{n \geq 1}$ m -függő valamely $m \geq 0$ számra, ha $\{X_1, \dots, X_k\}$ és $\{X_{k+m+1}, \dots\}$ függetlenek minden $k \geq 0$ -ra.

Jel. $\sigma_m^2 = \text{Var}(X_1) + 2 \sum_{i=1}^{m-1} \text{Cov}(X_1, X_{1+i})$

Legyen a becslendő statisztika: $T_n = \sqrt{n}(\bar{X}_n - \mu)$

Ekkor a bootstrap alteregója: $T_{n,n}^* = \sqrt{n}(\bar{X}_n^* - \bar{X}_n)$

Tétel: Legyen $\{X_n\}_{n \geq 1}$ stacionárius m -függő v.v. sorozat, $EX_1 = \mu$, $\sigma^2 = \text{Var}(X_1) \in (0, \infty)$, $\sum_{n=1}^m \text{Cov}(X_1, X_{1+i}) \neq 0$ és $\sigma_m^2 \neq 0$ Ekkor

$$\lim_{n \rightarrow \infty} \sup_x |P_*(T_{n,n}^* \leq x) - P(T_n \leq x)| \neq 0 \text{ m.m. } \forall x \neq 0 - \text{ra}$$

Az i.i.d. bootstrap előnyei és korlátai

- Az i.i.d. bootstrap előnyei a jackknife-hoz képest:
 - 1 Intervallumbecslésre és statisztika eloszlásának becslésére is használható
 - 2 "Bonyolultabb" statisztikák becslésére is alkalmas
- Az i.i.d. bootstrap korlátai:
 - 1 számításigényes
 - 2 bizonyos esetekben a becslés nem lesz konzisztens

Példa (Singh, 1981)

Def: $\{X_n\}_{n \geq 1}$ m -függő valamely $m \geq 0$ számra, ha $\{X_1, \dots, X_k\}$ és $\{X_{k+m+1}, \dots\}$ függetlenek minden $k \geq 0$ -ra.

Jel. $\sigma_m^2 = \text{Var}(X_1) + 2 \sum_{i=1}^{m-1} \text{Cov}(X_1, X_{1+i})$

Legyen a becslendő statisztika: $T_n = \sqrt{n}(\bar{X}_n - \mu)$

Ekkor a bootstrap alteregója: $T_{n,n}^* = \sqrt{n}(\bar{X}_n^* - \bar{X}_n)$

Tétel: Legyen $\{X_n\}_{n \geq 1}$ stacionárius m -függő v.v. sorozat, $EX_1 = \mu$, $\sigma^2 = \text{Var}(X_1) \in (0, \infty)$, $\sum_{n=1}^m \text{Cov}(X_1, X_{1+i}) \neq 0$ és $\sigma_m^2 \neq 0$ Ekkor

$$\lim_{n \rightarrow \infty} \sup_x |P_*(T_{n,n}^* \leq x) - P(T_n \leq x)| \neq 0 \text{ m.m. } \forall x \neq 0 - \text{ra}$$

Az i.i.d. bootstrap előnyei és korlátai

- Az i.i.d. bootstrap előnyei a jackknife-hoz képest:
 - 1 Intervallumbecslésre és statisztika eloszlásának becslésére is használható
 - 2 "Bonyolultabb" statisztikák becslésére is alkalmas
- Az i.i.d. bootstrap korlátai:
 - 1 számításigényes
 - 2 bizonyos esetekben a becslés nem lesz konzisztens

Példa (Singh, 1981)

Def: $\{X_n\}_{n \geq 1}$ m -függő valamely $m \geq 0$ számra, ha $\{X_1, \dots, X_k\}$ és $\{X_{k+m+1}, \dots\}$ függetlenek minden $k \geq 0$ -ra.

Jel. $\sigma_m^2 = \text{Var}(X_1) + 2 \sum_{i=1}^{m-1} \text{Cov}(X_1, X_{1+i})$

Legyen a becslendő statisztika: $T_n = \sqrt{n}(\bar{X}_n - \mu)$

Ekkor a bootstrap alteregója: $T_{n,n}^* = \sqrt{n}(\bar{X}_n^* - \bar{X}_n)$

Tétel: Legyen $\{X_n\}_{n \geq 1}$ stacionárius m -függő v.v. sorozat, $EX_1 = \mu$, $\sigma^2 = \text{Var}(X_1) \in (0, \infty)$, $\sum_{n=1}^m \text{Cov}(X_1, X_{1+i}) \neq 0$ és $\sigma_m^2 \neq 0$ Ekkor

$$\lim_{n \rightarrow \infty} \sup_x |P_*(T_{n,n}^* \leq x) - P(T_n \leq x)| \neq 0 \text{ m.m. } \forall x \neq 0 - \text{ra}$$

Az i.i.d. bootstrap előnyei és korlátai

- Az i.i.d. bootstrap előnyei a jackknife-hoz képest:
 - 1 Intervallumbecslésre és statisztika eloszlásának becslésére is használható
 - 2 "Bonyolultabb" statisztikák becslésére is alkalmas
- Az i.i.d. bootstrap korlátai:
 - 1 számításigényes
 - 2 bizonyos esetekben a becslés nem lesz konzisztens

Példa (Singh, 1981)

Def: $\{X_n\}_{n \geq 1}$ m -függő valamely $m \geq 0$ számra, ha $\{X_1, \dots, X_k\}$ és $\{X_{k+m+1}, \dots\}$ függetlenek minden $k \geq 0$ -ra.

Jel. $\sigma_m^2 = \text{Var}(X_1) + 2 \sum_{i=1}^{m-1} \text{Cov}(X_1, X_{1+i})$

Legyen a becslendő statisztika: $T_n = \sqrt{n}(\bar{X}_n - \mu)$

Ekkor a bootstrap alteregója: $T_{n,n}^* = \sqrt{n}(\bar{X}_n^* - \bar{X}_n)$

Tétel: Legyen $\{X_n\}_{n \geq 1}$ stacionárius m -függő v.v. sorozat, $EX_1 = \mu$, $\sigma^2 = \text{Var}(X_1) \in (0, \infty)$, $\sum_{n=1}^m \text{Cov}(X_1, X_{1+i}) \neq 0$ és $\sigma_m^2 \neq 0$ Ekkor

$$\lim_{n \rightarrow \infty} \sup_x |P_*(T_{n,n}^* \leq x) - P(T_n \leq x)| \neq 0 \text{ m.m. } \forall x \neq 0 - \text{ra}$$

- 1 Bevezetés
- 2 **Bootstrap módszerek**
 - Az i.i.d. bootstrap
 - **Blokk bootstrap módszerek**
 - A blokkméret kiválasztásának problémája
- 3 Alkalmazás egydimenziós szélességi adatokra
- 4 További kiterjesztési lehetőségek
- 5 Goodness-of-fit tesztek
- 6 Összefoglalás

- 1 $Y_t = X_{t \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- 2 Legyenek i_1, i_2, \dots minta az $\{1, \dots, n\}$ halmazon egyenletes eloszlásból
- 3 Adott b blokkméretre készítsünk $n' = mb$ ($n' \approx n$) pszeudo-megfigyelést:

$$Y_{(k-1)b+j}^* = Y_{i_m+j-1} \quad \text{ahol } j=1, \dots, b; \quad k=1, \dots, m$$

- 4 A minket érdeklő statisztika kiszámítása a pszeudo-megfigyelésekből:

$$\bar{Y}_{n'}^* = (n')^{-1} (Y_1^* + \dots + Y_{n'}^*)$$

Circular blokk bootstrap (CBB)

- 1 $Y_t = X_{t \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- 2 Legyenek i_1, i_2, \dots minta az $\{1, \dots, n\}$ halmazon egyenletes eloszlásból
- 3 Adott b blokkméretre készítsünk $n' = mb$ ($n' \approx n$) pszeudo-megfigyelést:

$$Y_{(k-1)b+j}^* = Y_{i_m+j-1} \quad \text{ahol } j=1, \dots, b; \quad k=1, \dots, m$$

- 4 A minket érdeklő statisztika kiszámítása a pszeudo-megfigyelésekből:

$$\bar{Y}_{n'}^* = (n')^{-1} (Y_1^* + \dots + Y_{n'}^*)$$

Circular blokk bootstrap (CBB)

- 1 $Y_t = X_{t_{\text{mod}(n)}}$ azaz periodikusan kiterjesztjük a mintát
- 2 Legyenek i_1, i_2, \dots minta az $\{1, \dots, n\}$ halmazon egyenletes eloszlásból
- 3 Adott b blokkméretre készítsünk $n' = mb$ ($n' \approx n$) pszeudo-megfigyelést:

$$Y_{(k-1)b+j}^* = Y_{i_m+j-1} \quad \text{ahol } j=1, \dots, b; \quad k=1, \dots, m$$

- 4 A minket érdeklő statisztika kiszámítása a pszeudo-megfigyelésekből:

$$\bar{Y}_{n'}^* = (n')^{-1} (Y_1^* + \dots + Y_{n'}^*)$$

Circular blokk bootstrap (CBB)

- 1 $Y_t = X_{t_{\text{mod}(n)}}$ azaz periodikusan kiterjesztjük a mintát
- 2 Legyenek i_1, i_2, \dots minta az $\{1, \dots, n\}$ halmazon egyenletes eloszlásból
- 3 Adott b blokkméretre készítsünk $n' = mb$ ($n' \approx n$) pszeudo-megfigyelést:

$$Y_{(k-1)b+j}^* = Y_{i_m+j-1} \quad \text{ahol } j=1, \dots, b; \quad k=1, \dots, m$$

- 4 A minket érdeklő statisztika kiszámítása a pszeudo-megfigyelésekből:

$$\bar{Y}_{n'}^* = (n')^{-1} (Y_1^* + \dots + Y_{n'}^*)$$

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_j és J_j nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_j és J_j függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_j és J_j függetlenek

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_j és J_j nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_j és J_j függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_j és J_j függetlenek

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_j és J_j nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_j és J_j függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_j és J_j függetlenek

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_j és J_j nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_j és J_j függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_j és J_j függetlenek

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_i és J_i nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_i és J_i függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_i és J_i függetlenek

Általánosított blokk bootstrap (GBB)

- $Y_i = X_{i \bmod(n)}$ azaz periodikusan kiterjesztjük a mintát
- I_1, I_2, \dots az adott $A \subseteq \{1, \dots, n\}$ halmazon független, egyenletes eloszlású v.v. sorozat \rightarrow blokkok kezdőindexei
- J_1, J_2, \dots i.i.d. v.v. sorozat \rightarrow blokkok hosszai
- I_i és J_i nem feltétlenül függetlenek!
- Blokkok kiválasztása $B(I_1, J_1) = \{Y_{I_1}, Y_{I_1+1}, \dots, Y_{I_1+J_1-1}\}$

Gyakran használt speciális esetek:

1 CBB (circular block bootstrap)

- $A = \{1, \dots, n\}$
- $P(J_1=b)=1 \rightarrow$ fix blokkméret
- I_i és J_i függetlenek

2 SBB(stationary block bootstrap)

- $A = \{1, \dots, n\}$
- $J_1 \sim \text{Geo}(p)$ (p -t egy "elvárt" blokkmérethez határozzuk meg)
- I_i és J_i függetlenek

- 1 Bevezetés
- 2 **Bootstrap módszerek**
 - Az i.i.d. bootstrap
 - Blokk bootstrap módszerek
 - **A blokkméret kiválasztásának problémája**
- 3 Alkalmazás egydimenziós szélességi adatokra
- 4 További kiterjesztési lehetőségek
- 5 Goodness-of-fit tesztek
- 6 Összefoglalás

Több tényezőtől is függ

- az adatgeneráló folyamat - van-e összefüggőség:
 - térbeli
 - időbeli
- a bennünket érdeklő statisztika, például
 - átlag
 - medián
 - lag(k) korreláció
- a statisztikából mit szeretnénk számolni, ha például
 - szórást $\implies b_{opt} = C \cdot n^{1/3}$
 - egyoldali eloszlást $\implies b_{opt} = C \cdot n^{1/4}$

De ha tudjuk is a nagyságrendet, $C=?$

Több tényezőtől is függ

- az adatgeneráló folyamat - van-e összefüggőség:
 - térbeli
 - időbeli
- a bennünket érdeklő statisztika, például
 - átlag
 - medián
 - lag(k) korreláció
- a statisztikából mit szeretnénk számolni, ha például
 - szórást $\implies b_{opt} = C \cdot n^{1/3}$
 - egyoldali eloszlást $\implies b_{opt} = C \cdot n^{1/4}$

De ha tudjuk is a nagyságrendet, $C=?$

Több tényezőtől is függ

- az adatgeneráló folyamat - van-e összefüggőség:
 - térbeli
 - időbeli
- a bennünket érdeklő statisztika, például
 - átlag
 - medián
 - lag(k) korreláció
- a statisztikából mit szeretnénk számolni, ha például
 - szórást $\implies b_{opt} = C \cdot n^{1/3}$
 - egyoldali eloszlást $\implies b_{opt} = C \cdot n^{1/4}$

De ha tudjuk is a nagyságrendet, $C=?$

Több tényezőtől is függ

- az adatgeneráló folyamat - van-e összefüggőség:
 - térbeli
 - időbeli
- a bennünket érdeklő statisztika, például
 - átlag
 - medián
 - lag(k) korreláció
- a statisztikából mit szeretnénk számolni, ha például
 - szórást $\implies b_{opt} = C \cdot n^{1/3}$
 - egyoldali eloszlást $\implies b_{opt} = C \cdot n^{1/4}$

De ha tudjuk is a nagyságrendet, **C=?**

Blokkméret kiválasztása (Politis & White)

Jel. $\mathcal{F}_{-\infty}^0 = \{X_n : n \leq 0\}$, $\mathcal{F}_k^\infty = \{X_n : n \geq k\}$

Def.: $\{X_t : t \in \mathbb{Z}\}$ erősen keverő, ha $\alpha_X(k) \rightarrow 0$ ($k \rightarrow \infty$), ahol

$$\alpha_X(k) = \sup\{|P(A \cap B) - P(A)P(B)| : A \in \mathcal{F}_{-\infty}^0, B \in \mathcal{F}_k^\infty\}$$

Tétel :

Tegyük fel, hogy $E|X_t|^{6+\delta} < \infty$, $\sum_{k=1}^{\infty} k^2(\alpha_X(k))^{\frac{\delta}{6+\delta}} < \infty$

valamely $\delta > 0$ -ra

Legyen $b = o(n^{1/2})$, $n \rightarrow \infty$ esetén $b \rightarrow \infty$.

$$\text{Ekkor } MSE(\sigma_{b,\bar{X}}^2) = \frac{G^2}{b^2} + D\frac{b}{n} + o(b^{-2}) + o\left(\frac{b}{n}\right)$$

$$\text{ahol } D = \frac{4}{3}g^2(0) \text{ és } G = \sum_{k=-\infty}^{\infty} |k|R(k)$$

$g(\cdot)$: spektrális sűrűségfüggvény

$R(\cdot)$: autokovariancia függvény

Blokkméret kiválasztása (Politis & White)

Jel. $\mathcal{F}_{-\infty}^0 = \{X_n : n \leq 0\}$, $\mathcal{F}_k^\infty = \{X_n : n \geq k\}$

Def.: $\{X_t : t \in \mathbb{Z}\}$ erősen keverő, ha $\alpha_X(k) \rightarrow 0$ ($k \rightarrow \infty$), ahol

$$\alpha_X(k) = \sup\{|P(A \cap B) - P(A)P(B)| : A \in \mathcal{F}_{-\infty}^0, B \in \mathcal{F}_k^\infty\}$$

Tétel:

Tegyük fel, hogy $E|X_t|^{6+\delta} < \infty$, $\sum_{k=1}^{\infty} k^2(\alpha_X(k))^{\frac{\delta}{6+\delta}} < \infty$

valamely $\delta > 0$ -ra

Legyen $b = o(n^{1/2})$, $n \rightarrow \infty$ esetén $b \rightarrow \infty$.

$$\text{Ekkor } MSE(\sigma_{b,\bar{X}}^2) = \frac{G^2}{b^2} + D\frac{b}{n} + o(b^{-2}) + o\left(\frac{b}{n}\right)$$

$$\text{ahol } D = \frac{4}{3}g^2(0) \text{ és } G = \sum_{k=-\infty}^{\infty} |k|R(k)$$

$g(\cdot)$: spektrális sűrűségfüggvény

$R(\cdot)$: autokovariancia függvény

Jel. $\mathcal{F}_{-\infty}^0 = \{X_n : n \leq 0\}$, $\mathcal{F}_k^\infty = \{X_n : n \geq k\}$

Def.: $\{X_t : t \in \mathbb{Z}\}$ erősen keverő, ha $\alpha_X(k) \rightarrow 0$ ($k \rightarrow \infty$), ahol

$$\alpha_X(k) = \sup\{|P(A \cap B) - P(A)P(B)| : A \in \mathcal{F}_{-\infty}^0, B \in \mathcal{F}_k^\infty\}$$

Tétel :

Tegyük fel, hogy $E|X_t|^{6+\delta} < \infty$, $\sum_{k=1}^{\infty} k^2(\alpha_X(k))^{\frac{\delta}{6+\delta}} < \infty$

valamely $\delta > 0$ -ra

Legyen $b = o(n^{1/2})$, $n \rightarrow \infty$ esetén $b \rightarrow \infty$.

Ekkor $MSE(\sigma_{b,\bar{X}}^2) = \frac{G^2}{b^2} + D\frac{b}{n} + o(b^{-2}) + o(\frac{b}{n})$

ahol $D = \frac{4}{3}g^2(0)$ és $G = \sum_{k=-\infty}^{\infty} |k|R(k)$

$g(\cdot)$: spektrális sűrűségfüggvény

$R(\cdot)$: autokovariancia függvény

Optimális blokkméret: $b_{opt} = \lceil (\frac{2G^2}{D})n^{1/3} \rceil$

Kérdés: hogyan becsüljük G -t és D -t

$$\hat{D} = \frac{4}{3} \hat{g}^2(0)$$

$$\hat{G} = \sum_{k=-M}^M \lambda\left(\frac{k}{M}\right) |k| \hat{R}(k)$$

$$\text{ahol } \hat{R}(k) = n^{-1} \sum_{k=1}^{n-|k|} (X_i - \bar{X}_n)(X_{i+|k|} - \bar{X}_n)$$

$$\lambda(t) = \begin{cases} 1 & \text{ha } |t| \in [0, 1/2] \\ 2(1 - |t|) & \text{ha } |t| \in [1/2, 1] \\ 0 & \text{különben} \end{cases}$$

$M = 2\hat{m}$, ahol \hat{m} : ahonnan a korrelogram "lényegében" 0

Blokkméret kiválasztása (Politis & White)

Optimális blokkméret: $b_{opt} = [(\frac{2G^2}{D})n^{1/3}]$

Kérdés: hogyan becsüljük G -t és D -t

$$\hat{D} = \frac{4}{3}\hat{g}^2(0)$$

$$\hat{G} = \sum_{k=-M}^M \lambda\left(\frac{k}{M}\right)|k|\hat{R}(k)$$

$$\text{ahol } \hat{R}(k) = n^{-1} \sum_{k=1}^{n-|k|} (X_i - \bar{X}_n)(X_{i+|k|} - \bar{X}_n)$$

$$\lambda(t) = \begin{cases} 1 & \text{ha } |t| \in [0, 1/2] \\ 2(1 - |t|) & \text{ha } |t| \in [1/2, 1] \\ 0 & \text{különben} \end{cases}$$

$M = 2\hat{m}$, ahol \hat{m} : ahonnan a korrelogram "lényegében" 0

Optimális blokkméret: $b_{opt} = [(\frac{2G^2}{D})n^{1/3}]$

Kérdés: hogyan becsüljük G -t és D -t

$$\hat{D} = \frac{4}{3}\hat{g}^2(0)$$

$$\hat{G} = \sum_{k=-M}^M \lambda\left(\frac{k}{M}\right)|k|\hat{R}(k)$$

$$\text{ahol } \hat{R}(k) = n^{-1} \sum_{i=1}^{n-|k|} (X_i - \bar{X}_n)(X_{i+|k|} - \bar{X}_n)$$

$$\lambda(t) = \begin{cases} 1 & \text{ha } |t| \in [0, 1/2] \\ 2(1 - |t|) & \text{ha } |t| \in [1/2, 1] \\ 0 & \text{különben} \end{cases}$$

$M = 2\hat{m}$, ahol \hat{m} : ahonnan a korrelogram "lényegében" 0

- Legyenek X_1, X_2, \dots, X_n egyváltozós stacionárius megfigyelések;
 $EX_i = \mu, \text{Var}(X_i) = \sigma^2$
- Ha X_1, X_2, \dots, X_n i.i.d., akkor $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$
- Sorozat összefüggőség \rightarrow nagyobb variancia
- Effektív mintaméret (n_e):

$$n_e = \frac{\sigma^2}{\text{Var}^*(\bar{X})}$$

ahol $\text{Var}^*(\bar{X})$: becsült variancia \leftarrow bootstrap

- Legyenek X_1, X_2, \dots, X_n egyváltozós stacionárius megfigyelések;
 $EX_i = \mu, \text{Var}(X_i) = \sigma^2$
- Ha X_1, X_2, \dots, X_n i.i.d., akkor $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$
- Sorozat összefüggőség \rightarrow nagyobb variancia
- Effektív mintaméret (n_e):

$$n_e = \frac{\sigma^2}{\text{Var}^*(\bar{X})}$$

ahol $\text{Var}^*(\bar{X})$: becsült variancia \leftarrow bootstrap

- Legyenek X_1, X_2, \dots, X_n egyváltozós stacionárius megfigyelések;
 $EX_j = \mu, \text{Var}(X_j) = \sigma^2$
- Ha X_1, X_2, \dots, X_n i.i.d., akkor $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$
- Sorozat összefüggőség \rightarrow nagyobb variancia
- Effektív mintaméret (n_e):

$$n_e = \frac{\sigma^2}{\text{Var}^*(\bar{X})}$$

ahol $\text{Var}^*(\bar{X})$: becsült variancia \leftarrow bootstrap

- Legyenek X_1, X_2, \dots, X_n egyváltozós stacionárius megfigyelések;
 $EX_j = \mu, \text{Var}(X_j) = \sigma^2$
- Ha X_1, X_2, \dots, X_n i.i.d., akkor $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$
- Sorozat összefüggőség \rightarrow nagyobb variancia
- Effektív mintaméret (n_e):

$$n_e = \frac{\sigma^2}{\text{Var}^*(\bar{X})}$$

ahol $\text{Var}^*(\bar{X})$: becsült variancia \leftarrow **bootstrap**

Politis & White algoritmusával

- $n=2591$ megfigyelés (szélességi adatok heti maximumai) 5 német településre
- Automatikus blokkméret-kiválasztó algoritmus eredménye:

Település	Optimális blokkméret
Hamburg	31
Hannover	11
Bremerhaven	28
Fehmarn	31
Schleswig	15

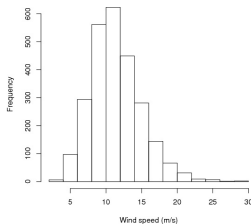


- 30 hét \approx 0,58 év \rightarrow meteorológiailag értelmetlen

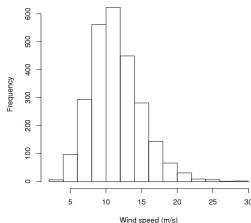
Módszer:

- 1 AR(1) modell illesztése az adatokra:
 $X_t = \mu + \phi X_{t-1} + Z_t$, $Z_t \sim$ extrém érték eloszlás
- 2 Az elméleti $Var(\bar{X}_n)$ kiszámolása az AR(1) paramétereiből:
$$Var(\bar{X}_n) = \frac{\sigma^2}{n(1-\hat{\phi}^2)} \frac{2\hat{\phi}^{n+1} - n\hat{\phi}^2 - 2\hat{\phi} + n}{n(\hat{\phi}-1)^2}$$
- 3 b^* optimális blokkméret: ahol az átlag szimulált varianciája először elmetszi az elméleti értéket

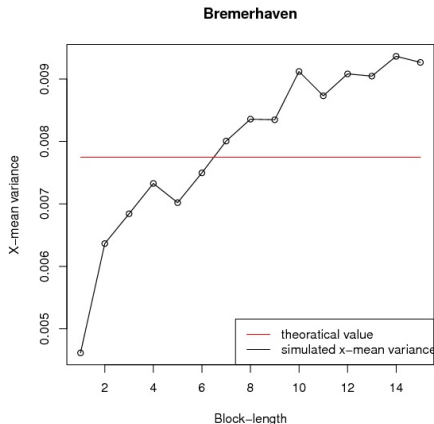
Histogram of wind speed maxima (Bremerhaven)



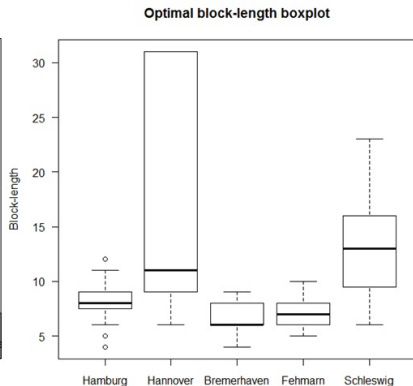
Histogram of wind speed maxima (Bremerhaven)



A bootstrap szimuláció eredményei



$$b^* = 6$$



Bootstrap szimulációs eredmények:

Település	Optimális blokkméret	Bootstrap $Var(\bar{X})$	Elméleti $Var(\bar{X})$	i.i.d. $Var(\bar{X})$	Mintaméret csökkenése
Hamburg	8-9	0,0038	0,0034	0,0020	1,85
Hannover	7	0,0067	0,0071	0,0042	1,59
Bremerhaven	6	0,0073	0,0077	0,0043	1,71
Fehmarn	7	0,0035	0,0034	0,0020	1,74
Schleswig	13	0,0037	0,0030	0,0018	2,09

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig

$$\hat{\varepsilon}_t = R_t - \bar{R}_t, \quad \text{ahol } t = \hat{p} + 1, \dots, n$$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig

$$\hat{\varepsilon}_t = R_t - \bar{R}_t, \quad \text{ahol } t = \hat{p} + 1, \dots, n$$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_{\hat{p}} = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig

$$\hat{\varepsilon}_t = R_t - \bar{R}_t, \quad \text{ahol } t = \hat{p} + 1, \dots, n$$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig

$$\hat{\varepsilon}_t = R_t - \bar{R}_t, \quad \text{ahol } t = \hat{p} + 1, \dots, n$$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig
 $\hat{\varepsilon}_t = R_t - \bar{R}_t$, ahol $t = \hat{p} + 1, \dots, n$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig
 $\hat{\varepsilon}_t = R_t - \bar{R}_t$, ahol $t = \hat{p} + 1, \dots, n$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig
 $\hat{\varepsilon}_t = R_t - \bar{R}_t$, ahol $t = \hat{p} + 1, \dots, n$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig
 $\hat{\varepsilon}_t = R_t - \bar{R}_t$, ahol $t = \hat{p} + 1, \dots, n$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

AR-sieve bootstrap

- Feltétel: a folyamat stacionárius és jól becsülhető $AR(p)$ modellel:

$$X_t - \mu_X = \sum_{j=1}^p \phi_j (X_{t-j} - \mu_X) + \varepsilon_t, \quad t \in \mathbb{Z}$$

$$\text{ahol } \mu_X = EX_t$$

$(\varepsilon_t)_{t \in \mathbb{Z}}$ i.i.d., $E(\varepsilon_t) = 0$ és ε_t független $\{X_s; s < t\}$ -től

- Paraméterek és hibák becslése:

① $\hat{p} = ? \rightarrow \text{AIC}$

② $\hat{\mu}_X = n^{-1} \sum_{t=1}^n X_t$

③ $\hat{\phi}_1, \dots, \hat{\phi}_p = ? \rightarrow \text{Yule-Walker módszer}$

④ $R_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\phi}_j X_{t-j}$, ahol $t = \hat{p} + 1, \dots, n$ ebből pedig
 $\hat{\varepsilon}_t = R_t - \bar{R}_t$, ahol $t = \hat{p} + 1, \dots, n$

- Bootstrap minta konstruálásának lépései:

- ε_t^* : véletlen elem $\{\hat{\varepsilon}_{\hat{p}+1}, \dots, \hat{\varepsilon}_n\}$ halmazból

- Nagy u -ra $(X_{-u}^*, \dots, X_{-u+\hat{p}-1}^*) = (\hat{\mu}_X, \dots, \hat{\mu}_X)$ (a folyamat indítása)

- $X_t^* = \mu_X + \sum_{j=1}^p \phi_j (X_{t-j}^* - \mu_X) + \varepsilon_t^* \quad t \in \mathbb{Z}$

- Ebből a bootstrap minta: $\{X_1^*, \dots, X_n^*\}$

Lag(1) korreláció becslése

3 módszert próbáltam ki:

- 1 "Naiv" (sima) blokk bootstrap ($b=10$)
- 2 Dupla blokk bootstrap ($b_1 = 2$ és $b_2=10$)
- 3 AR sieve (szita) bootstrap

Futtatási eredmények lag(1) korrelációra:

Módszer	Hamburg	Hannover	Bremerhaven	Fehmarn	Schleswig
Mintából	0,2755	0,2528	0,2901	0,2548	0,2588
Naiv blokk bootstrap	0,2485	0,2277	0,2598	0,2281	0,2306
Dupla blokk bootstrap	0,2755	0,252	0,2907	0,255	0,2565
AR-sieve bootstrap	0,274	0,2514	0,2887	0,2531	0,2588

3 módszert próbáltam ki:

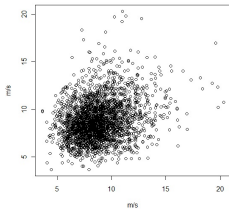
- 1 "Naiv" (sima) blokk bootstrap ($b=10$)
- 2 Dupla blokk bootstrap ($b_1 = 2$ és $b_2=10$)
- 3 AR sieve (szita) bootstrap

Futtatási eredmények lag(1) korrelációra:

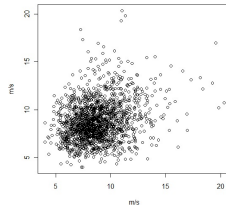
Módszer	Hamburg	Hannover	Bremerhaven	Fehmarn	Schleswig
Mintából	0,2755	0,2528	0,2901	0,2548	0,2588
Naiv blokk bootstrap	0,2485	0,2277	0,2598	0,2281	0,2306
Dupla blokk bootstrap	0,2755	0,252	0,2907	0,255	0,2565
AR-sieve bootstrap	0,274	0,2514	0,2887	0,2531	0,2588

Lag(1) pontdiagramok Schleswig szélességi adataira

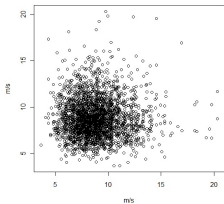
Lag(1) pontdiagram az eredeti mintára



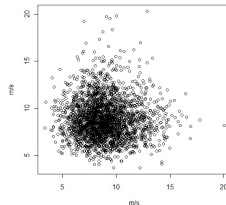
Lag(1) pontdiagram dupla blokk bootstrap esetén



Lag(1) pontdiagram 'naliv' blokk bootstrap esetén



Lag(1) pontdiagram sieve bootstrap esetén



Goodness-of-fit (GoF) tesztek

- A modell paramétereinek becslése
- Goodness-of-fit teszt: $H_0: C \in \mathbb{C}_0 = \{ C_\theta, \theta \in \Theta \}$

- Cramér-von Mises tesztstatisztika:

$$T = n \int_{-\infty}^{\infty} (F_n(x) - F(x))^2 \phi(x) dF(x)$$

- F : eloszlásfv.
- F_n : empirikus eloszlásfv.
- Φ : súlyfv. – Anderson-Darling: $\Phi(x) = \frac{1}{F(x)(1-F(x))}$

- Kritikus érték - szimuláció

- 1 Minta szimulálása a \mathbb{C}_0 copula modellből H_0 esetén
- 2 $\hat{\theta}$ újbóli becslése ML-módszerrel
- 3 A tesztstatisztika kiszámítása
Ismétlések, p-érték becslése

Több dimenzióban:

- Probability integral transformation (PIT) - összehúzás a d -dimenziós egységkockára
- Kendall-transzformáció (K -függvény):

$$K(\theta, t) = P(C_\theta(F_1(X_1), \dots, F_d(X_d)) \leq t)$$

Goodness-of-fit (GoF) tesztek

- A modell paramétereinek becslése
- Goodness-of-fit teszt: $H_0: C \in \mathbb{C}_0 = \{ C_\theta, \theta \in \Theta \}$

- Cramér-von Mises tesztstatisztika:

$$T = n \int_{-\infty}^{\infty} (F_n(x) - F(x))^2 \phi(x) dF(x)$$

- F : eloszlásfv.
- F_n : empirikus eloszlásfv.
- Φ : súlyfv. – Anderson-Darling: $\Phi(x) = \frac{1}{F(x)(1-F(x))}$

- Kritikus érték - szimuláció

- 1 Minta szimulálása a \mathbb{C}_0 copula modellből H_0 esetén
- 2 $\hat{\theta}$ újbóli becslése ML-módszerrel
- 3 A tesztstatisztika kiszámítása
Ismétlések, p-érték becslése

Több dimenzióban:

- Probability integral transformation (PIT) - összehúzás a d -dimenziós egységkockára
- Kendall-transzformáció (K -függvény):

$$K(\theta, t) = P(C_\theta(F_1(X_1), \dots, F_d(X_d)) \leq t)$$

Goodness-of-fit (GoF) tesztek

- A modell paramétereinek becslése
- Goodness-of-fit teszt: $H_0: C \in \mathbb{C}_0 = \{ C_\theta, \theta \in \Theta \}$

- Cramér-von Mises tesztstatisztika:

$$T = n \int_{-\infty}^{\infty} (F_n(x) - F(x))^2 \phi(x) dF(x)$$

- F : eloszlásfv.
- F_n : empirikus eloszlásfv.
- Φ : súlyfv. – Anderson-Darling: $\Phi(x) = \frac{1}{F(x)(1-F(x))}$

- Kritikus érték - szimuláció

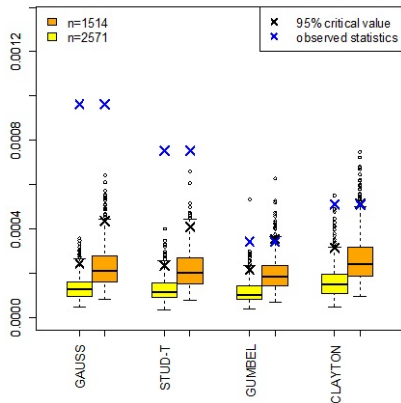
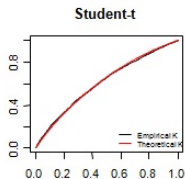
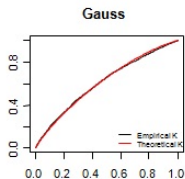
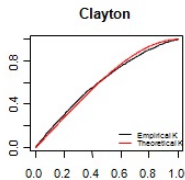
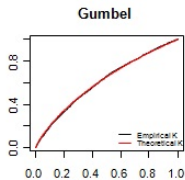
- 1 Minta szimulálása a \mathbb{C}_0 copula modellből H_0 esetén
- 2 $\hat{\theta}$ újbóli becslése ML-módszerrel
- 3 A tesztstatisztika kiszámítása
Ismétlések, p-érték becslése

Több dimenzióban:

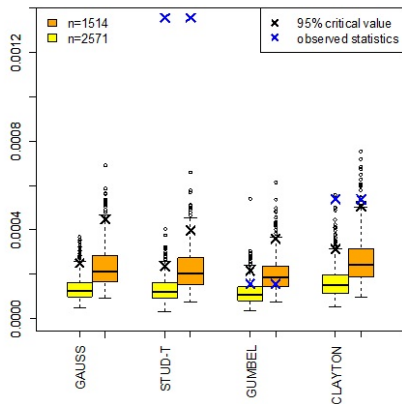
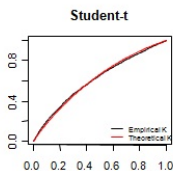
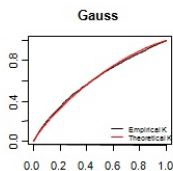
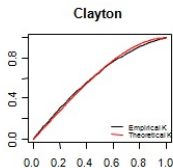
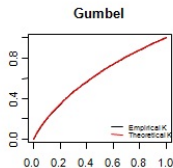
- Probability integral transformation (PIT) - összehúzás a d -dimenziós egységkockára
- Kendall-transzformáció (K -függvény):

$$K(\theta, t) = P(C_\theta(F_1(X_1), \dots, F_d(X_d)) \leq t)$$

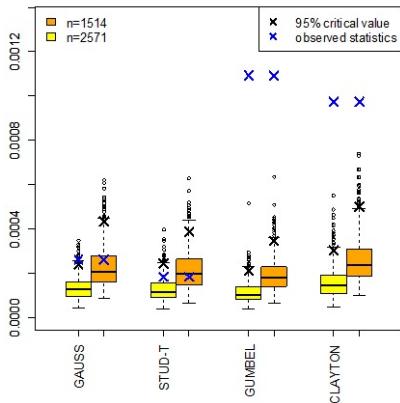
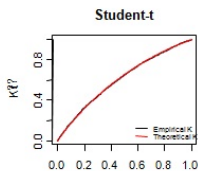
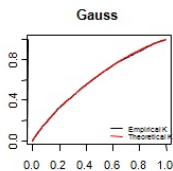
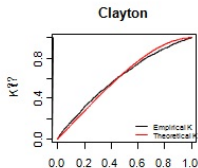
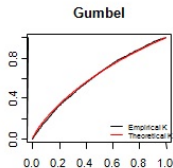
Bremerhaven & Fehmarn



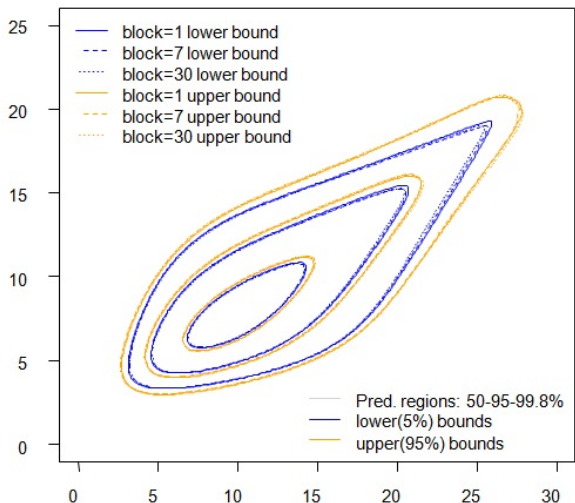
Bremerhaven & Schleswig



Fehrnann & Schleswig




Előrejelzési régiók (Bremerhaven & Fehmarn)





- A bootstrap módszer kiválasztásánál figyelembe kell venni:
 - milyen az adatok struktúrája
 - mi a bennünket érdeklő statisztika
 - mire akarjuk felhasználni
- nincs feltétlen legjobb módszer - egy (fél)paraméteres vagy egy paraméter nélküli is jó lehet


Köszönetnyilvánítás

- Zempléni András
- Rakonczai Pál

 S.N. Lahiri:
Resampling Methods for Dependent Data.
Springer, 2003.

 L. Kish:
Survey Sampling.
J. Wiley, 1965.

 P. Rakonczai, L. Varga, A. Zempléni:
Copula Fitting to Autocorrelated Data, with Applications to Wind
Speed Modelling.
Working paper, november 11., 2010.

 P. Rakonczai, A. Zempléni:
Copulas and goodness of fit tests. Recent advances in stochastic
modeling and data analysis.
World Scientific, pp. 198-206, 2007.



P. Bühlmann:

Bootstraps for Time Series.

Statistical Science, Vol. 17, pp. 52-72, 2002.



P. Bühlmann, H.R. Künsch:

Block length selection in the bootstrap for time series.

Computational Statistics Data Analysis, pp. 295-310, 1999.



D.N. Politis, H. White:

Automatic Block-Length Selection for the Dependent Bootstrap.

Econometric Reviews, Vol. 23, pp. 53-70, 2004. 2000.



C. Genest, B. Rémillard:

Validity of the parametric bootstrap for goodness-of-fit testing in semiparametric models.

Annales de l'Institut Henri Poincaré - Probabilités et Statistiques, Vol. 44, No. 6, 10961127, 2008. 2000.