

Butstrapa metodes atkarīgiem datiem

Mārcis Bratka

2010

Tapered block bootstrap

X_1, X_2, \dots ir stacionāri un vienādi sadalīti gadījuma lielumi no \mathbb{R}^m ar sadalījumu F . X_1, \dots, X_n ir dota izlase un $T_n = T_n(X_1, \dots, X_n)$ ir novērtējums parametram $T(F)$.
Piemēram,

$$T_n = f(n^{-1} \sum_{t=1}^n \phi(X_t)),$$

kur $f : \mathbb{R}^d \rightarrow \mathbb{R}$ un $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^d$.

MBB (Kunsch 1989) novērtējums $\text{Var}(\sqrt{n}T_n)$ bloka garumam b ir $\hat{\sigma}_{b, \text{MBB}}^2$. Pie dažiem momentu un jaukto procesu nosacījumiem

$$\text{Bias}(\hat{\sigma}_{b, \text{MBB}}^2) = O(1/b) \text{ un } \text{Var}(\hat{\sigma}_{b, \text{MBB}}^2) = O(b/n).$$

T_n ir lineāra, ja

$$T_n = T(F) + n^{-1} \sum_{t=1}^n \text{IF}(X_t, F) + R_n,$$

kur R_n ir mazs. Funkcija IF statistikai T_n ir definēta

$$\text{IF}(y, F) = \lim_{\epsilon \downarrow 0} \frac{T((1-\epsilon)F + \epsilon\delta_y) - T(F)}{\epsilon},$$

kur δ_y punktā y ir 1. Praktiski $\text{IF}(y, F)$ netiek lietota, jo F nav zināms, bet lieto empīrisko $\text{IF}(y, \hat{F}_n)$, kur \hat{F}_n ir empīrisks sadalījums ar $1/n$ katrā X_1, \dots, X_n .

$$\sup_y |\text{IF}(y, \hat{F}_n) - \text{IF}(y, F)| = O_P(1/\sqrt{n}).$$

$$Y_t := \text{IF}(X_t, \hat{F}_n), t = 1, \dots, n.$$

Piemērs

$T_n = \bar{X} = n^{-1} \sum_{t=1}^n X_t$, kurai $m = 1$ un f un ϕ ir identitātes funkcijas, $R_m = 0$. $\text{IF}(y, F) = y - \int x dF(x)$ un $\text{IF}(y, \hat{F}_n) = y - \hat{X}$, tātad $Y_t = \text{IF}(X_t, \hat{F}_n) = X_t - \hat{X}$.

Tapered block bootstrap (TBB) ideja ir lietot svarus $w_n(\cdot)$, $n = 1, 2, \dots$, lai samazinātu bloka sākuma un beigu elementu nozīmi. $w_n(t)$ vērtības ir $[0, 1]$ un $w_n(t) = 0, t \notin \{1, \dots, n\}$.
 $\|w_n\|_1 = \sum_{t=1}^n |w_n(t)|$ un $\|w_n\|_2 = \sqrt{\sum_{t=1}^n w_n^2(t)}$.

$$w_n(t) = w\left(\frac{t-0.5}{n}\right), \text{ kur}$$

$w(\cdot)$ apmierina nosacījumus:

- $w(t) \in [0, 1]$ visiem $t \in \mathbb{R}$, $w(t) = 0$, ja $t \notin [0, 1]$, un $w(t) > 0$ pie $t = 1/2$;
- $w(t)$ ir simetriska pret $t = 1/2$ un nedilstoša $t \in [0, 1/2]$.

Ja $w(t) = 1_{[0,1]}(t)$, tad TBB ir ekvivalents MBB. Lai TBB iegūtu novirzes uzlabojumu w ir jābūt gludai, t.i.,
 $w * w(t) = \int_{-1}^1 w(x)w(x + |t|)dx$ ir divreiz nepārtaukti
diferencējama $t = 0$.

TBB procedūra:

- izvēlamies bloka garumu b un i_0, i_1, \dots, i_{k-1} ir vienmērīgi sadalīti no $\{1, \dots, Q\}$, $Q = n - b + 1$, $k = \lfloor n/b \rfloor$;
- konstruējam Y_1^*, \dots, Y_l^* ar $l = kb$

$$Y_{mb+j}^* := w_b(j) \frac{\sqrt{b}}{\|w_b\|_2} Y_{i_m+j-1}, \text{ kur } j = 1, \dots, b, \\ m = 0, 1, \dots, k-1;$$

- $\bar{Y}_l^* = l^{-1} \sum_{i=1}^l Y_i^*$.

TBB dispersijas novērtējums $\hat{\sigma}_{b, \text{TBB}}^2 := \text{Var}^*(\sqrt{l} \bar{Y}_l^*)$ un sadalījuma $P(\sqrt{n}(T_n - T(F)) \leq x)$ novērtējums $P^*(\sqrt{l} \bar{Y}_l^* \leq x)$.

Teorēma

$$E\hat{\sigma}_{b,\text{TBB}}^2 = \sigma_\infty^2 + \Gamma/b^2 + o(1/b^2) \text{ un } \text{Var}(\hat{\sigma}_{b,\text{TBB}}^2) = \Delta \frac{b}{n} + o(b/n),$$

kur Γ un Δ ir konstantes. Pie tam

$$\sup_{\mathbf{x}} |P^*(\sqrt{l}(\bar{Y}_1^* - E^*\bar{Y}_1^*) \leq \mathbf{x}) - P(\sqrt{n}(T_n - T(F)) \leq \mathbf{x})| \rightarrow_p 0.$$

$\text{MSE}(\hat{\sigma}_{b, \text{TBB}}^2) = \frac{\Gamma^2}{b^4} + \Delta \frac{b}{N} + \dots$ MSE tiek minimizēts pie
 $b_{\text{opt}} = \left(\frac{4\Gamma^2}{\Delta}\right)^{1/5} n^{1/5}$. Tad minimālā asimptotiskā MSE

$$\text{MSE}_{\text{opt}} = \left(\Gamma^{2/5} \Delta^{4/5} \frac{5}{4^{4/5}}\right) n^{-4/5}.$$

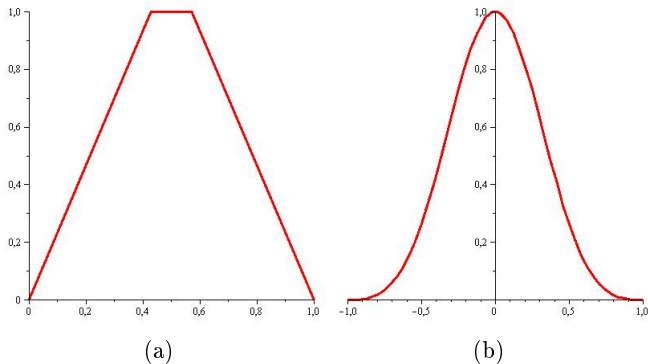
Δ un Γ ir atkarīgas no w un kovariācijām. w tiek izvēlēts tā, lai minimizētu

$$|\tilde{w}''(0)| * \|\tilde{w}\|_2^4, \text{ kur}$$

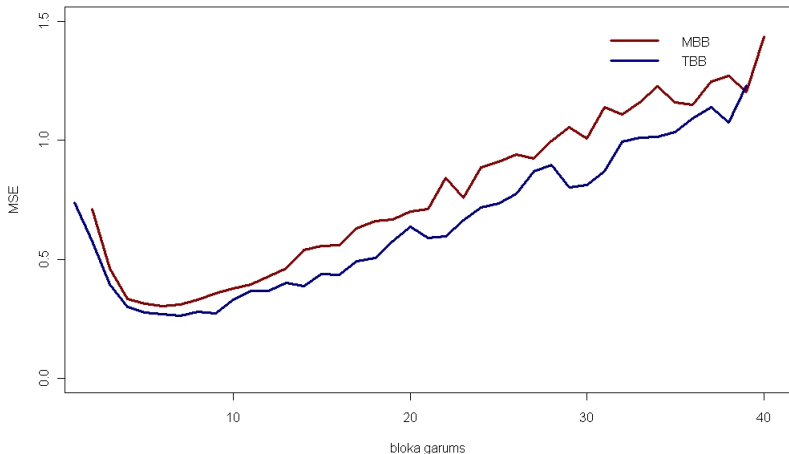
$$\tilde{w}(t) = (w * w)(t) / (w * w)(0).$$

$$w_c^{\text{TRAP}}(t) = \begin{cases} t/c & t \in [0, c]; \\ 1 & t \in (c, 1 - c); \\ (1 - c)/t & t \in (1 - c, 1]. \end{cases}$$

kur $c \simeq 0.43$.



$X_t = 0.6\sin(X_{t-1}) + Z_t$, kur Z_t ir $N(0, 1)$ un $n = 200$. $T_n = \bar{X}$.



Dependent wild bootstrap

Doti X_1, \dots, X_n ir stacionāri ar $\mu = E(X_t)$ un $\gamma_k = \text{cov}(X_0, X_k)$.
Dependent wild bootstrap (DWB) izlases elementi ir

$$X_i^* = \bar{X}_n + (X_i - \bar{X}_n)W_i, \quad i = 1, \dots, n, \quad \text{kur}$$

$\bar{X}_n = n^{-1} \sum_{t=1}^n X_t$ un $\{W_i\}_{i=1}^n$ ir gadījuma lielumi.

Pieņēmums

Gadījuma lielumi $\{W_t\}_{t=1}^n$ ir neatkarīgi no X_1, \dots, X_n ,
 $E(W_t) = 0$ un $\text{var}(W_t) = 1$, $t = 1, \dots, n$. W_t ir stacionāri ar
 $\text{cov}(W_t, W_{t'}) = a((t - t')/l)$, kur $a(\cdot)$ ir kodola funkcija un $l = l_n$.

Pie momentu un kovariāciju nosacījumiem

$T_n := \sqrt{n}(\bar{X}_n - \mu) \rightarrow_D N(0, \sigma_\infty^2)$, kur $\sigma_\infty^2 = \sum_{j=-\infty}^{\infty} \gamma_j$.

$$\begin{aligned} \hat{\sigma}_{l,DWB}^2 &= n^{-1} \sum_{t,t'=1}^n (X_t - \bar{X}_n)(X_{t'} - \bar{X}_n) \text{cov}^*(W_t, W_{t'}) \\ &= n^{-1} \sum_{h=1-n}^{n-1} \sum_{t=\max(1,1-h)}^{\min(n,n-h)} (X_t - \bar{X}_n)(X_{t+h} - \bar{X}_n) a(h/l) \\ &= 2\pi \hat{f}_n(0), \text{ kur} \end{aligned}$$

$\hat{f}_n(\lambda) = (2\pi)^{-1} \sum_{k=1-n}^{n-1} a(k/l) \hat{\gamma}_k \cos(k\lambda)$ ir spektra novērtējums un $a(\cdot)$ ir lag window funkcija.

Teorēma

$$E(\hat{\sigma}_{l,DWB}^2) = \sigma_{\infty}^2 + \Gamma/l^2 + o(1/l^2) \text{ un } \text{var}(\hat{\sigma}_{l,DWB}^2) = \Delta \frac{1}{n} + o(1/n),$$

kur Γ un Δ ir konstantes.

- $\hat{\sigma}_{l,DWB}^2$ un $\hat{\sigma}_{l,TBB}^2$ ir asimptotiski vienādi, ja $a(\mathbf{x}) = w * w(\mathbf{x})/w * w(0)$ un l ir vienādi abām metodēm.
- Ja $a(\mathbf{x}) = (1 - |\mathbf{x}|)1(|\mathbf{x}| \leq 1)$, tad DWB ir vienāds ar MBB.

Vispārīgā gadījumā DBW ir lietojams lineārām statistikām.

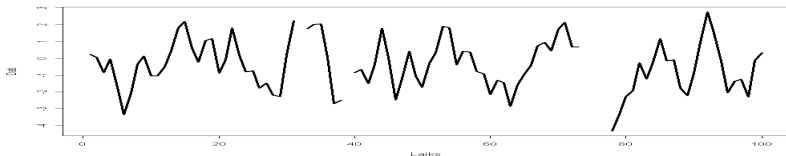
- Dotai izlasei X_1, \dots, X_n ar sadalījumu F parametra $\theta = T(F)$ novērtējums ir $\hat{\theta} = T(\rho_n)$, kur $\rho_n = n^{-1} \sum_{t=1}^n \delta_{X_t}$.
- $T_n(\rho_n) = T(F) + n^{-1} \sum_{t=1}^n IF(X_t, F) + R_n$.
- Pie nosacījumiem $\text{nvar}(\hat{\theta}) = n^{-1} \text{var}(\sum_{t=1}^n IF(X_t, \rho_n)) + o(1)$.
- Butstrapotais mērs $\rho_n^* = n^{-1} \sum_{i=1}^n (W_i + 1 - \bar{W}_n) \delta_{X_i}$, kur $\bar{W}_n = n^{-1} \sum_{t=1}^n W_t$ un $\{W_t\}_{t=1}^n$ ir gadījuma lielumi.

Ja $T(F) = \int x dF$, tad $T(\rho_n) = \bar{X}_n$ un $T(\rho_n^*) = n^{-1} \sum_{t=1}^n (W_t + 1 - \bar{W}_t) X_t = \bar{X}_n + n^{-1} \sum_{t=1}^n W_t (X_t - \bar{X}_n)$.

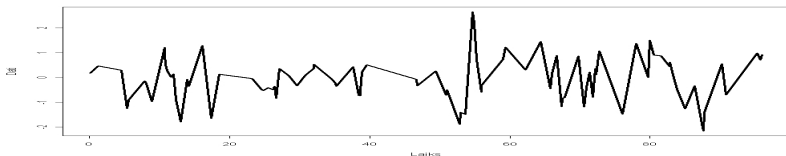
Vispārīgākām statistikām var būt problēmas ar butstrapa izlases iegūšanu, jo ρ_n^* nav varbūtību mērs.

- Šobrīd nav iegūti rezultāti par second-order correctness, bet ir zināms, ja W_t tiek izvēlēts ar normālo sadalījumu, tad vispārīgā gadījumā second-order correctness netiks sasniegts.
- Ja X_t ir normālais sadalījums, tad, iespējams, DWB būs second-order correctness ar W_t normālo sadalījumu, gludas funkcijas modeļa gadījumā.
- Šobrīd nav rekomendācijas, kā izvēlēties W_t pie dotiem X_1, \dots, X_n , lai sasniegtu second-order correctness.
- No precizitātes viedokļa svarīgāk ir izvēlēties $a(\cdot)$ un l (first-order), kā W_t .

- Trūkstoši elementi (missing values).



- Laika intervāli ir gadījuma lielumi.



$X_t = 0.6\sin(X_{t-1}) + Z_t$, kur Z_t ir $N(0, 1)$ un $n = 200$. $T_n = \bar{X}$ ar
 $a(x) = (1 - |x|)1(|x| \leq 1)$.

