

Maiņas punkta noteikšana laiktrendu analīzē

Agris Vaselāns

Latvijas Universitāte

2011. gada 10.novembrī.

Galvenie jautājumi

- Vai statistiskajā modelī notikušas izmaiņas?
- Kad modelī notikušas izmaiņas?
- Vai maiņas punkts ir tikai viens?
- Cik maiņas punkti bijuši?
- utt.

Definīcija

Laika momentu, kad statistiskajā modelī notikušas izmaiņas sauc par maiņas punktu.

Teorēma

Ja $\gamma(h)$ ir stacionāra procesa $\{x_t\}$ autokovariāciju funkcija, turklāt

$$\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty,$$

tad

$$\gamma(h) = \int_{-1/2}^{1/2} e^{2\pi i \omega h} f(\omega) d\omega, \quad h = 0, \pm 1, \pm 2, \dots$$

un procesa $\{x_t\}$ spektrālā blīvuma funkcija ir

$$f(\omega) = \sum_{h=-\infty}^{\infty} \gamma(h) e^{-2\pi i \omega h}, \quad -1/2 \leq \omega \leq 1/2.$$

Definīcija

Ja doti dati x_1, x_2, \dots, x_n , tad par diskrēto Furjē transformāciju (DFT) sauc

$$\begin{aligned}d(\omega_j) &= n^{-1/2} \sum_{t=1}^n x_t e^{-2\pi i \omega_j t} = a(j) + ib(j) \\ &= n^{-1/2} \sum_{t=1}^n V(t) (\cos(2\pi \omega_j t) + i \sin(2\pi \omega_j t)),\end{aligned}$$

kur $j = 0, 1, \dots, n - 1$ un frekvences $\omega_j = j/n$ tiek sauktas par Furjē jeb fundamentālajām frekvencēm, $a(j)$ un $b(j)$ sauc par Furjē koeficientiem.

Definīcija

Ja doti dati x_1, x_2, \dots, x_n , tad par periodogrammu sauc

$$I(\omega_j) = |d(\omega_j)|^2 = \frac{n}{2}(a(j)^2 + b(j)^2)$$

kur $j = 0, 1, \dots, n - 1$.

Īpašības

$$Ea(j) \rightarrow 0, Eb(j) \rightarrow 0, n \rightarrow \infty$$

$$Da(j) \rightarrow f(\omega_j)/2, Db(j) \rightarrow f(\omega_j)/2, n \rightarrow \infty$$

Tātad, ja stohastiskais process ir Gausa, tad $a(j)$ un $b(j)$ ir asimptotiski i.i.d. un sadalīti pēc $N(0, f(\omega_j)/2)$.

Tā kā periodogramma nav spektra būtisks novērtējums, tad nepieciešams definēt gludināto periodogrammu.

Definīcija

Par kodolu gludināto periodogrammu sauc

$$\hat{f}(\omega) = \frac{\sum_{j \in \mathcal{Z}} K\left(\frac{\omega - \omega_j}{h}\right) I(j)}{\sum_{j \in \mathcal{Z}} K\left(\frac{\omega_j}{h}\right) I(j)},$$

kur $K(\cdot)$ ir kodolsl, $I(j)$ ir periodogramma frekvencē ω_j .

Procedūras apraksts

Apskata stacionāru procesu x_t , $t = 1, \dots, T$.

- Aprēķina Furjē koeficientus, izmantojot FFT (fast fourier transform)
- Apzīmē $a(T)^* = b(T)^* = 0$ (Ja T pāra, tad $a(T/2)^* = b(T/2)^* = 0$.)
- Izmantojot kādu no butstrapa procedūrām, iegūst jaunus Furjē koeficientus $a(j)^*$ un $b(j)^*$, $j = 1, \dots, N$, $N = \lfloor (T - 1)/2 \rfloor$.
- Iegūst pārējos koeficientus, no sakarībām

$$a(j)^* = a(T - j)^*, \quad b(j)^* = b(T - j)^*.$$

- Izmanto inverso FFT, lai iegūtu centrēto laicrindu
 $z_t^* = x_t^* - \mu_t^*$.

Procedūras apraksts

Apskata stacionāru procesu x_t , $t = 1, \dots, T$.

- Ģenerē standarta normālā sadalījuma gadījuma lielumus $\{s_j : 1 \leq j \leq 2N\}$
- Iegūst gadījuma lielumu butstrapa izlasi s_j^* un jaunos koeficientus

$$a(j)^* = \sqrt{\widehat{f(\omega_j)}/2s_j^*}, \quad b(j)^* = \sqrt{\widehat{f(\omega_j)}/2s_{N+j}^*},$$

kur $\widehat{f(\cdot)}$ ir gludinātais spektrs, ka

$$\sup_{\omega} |\widehat{f(\omega)} - f(\omega)| \rightarrow^P 0.$$

Hipotēžu pārbaude vidējai vērtībai (skalācijas parametram)

$$H_0 : X_i = a + e_i, i = 1, 2, \dots, M$$

$$H_1 : \exists m \in \{1, 2, \dots, M\},$$

ka

$$X_i = a + e_i, i = 1, 2, \dots, m,$$

$$X_i = a + \delta + e_i, i = m + 1, 2, \dots, M,$$

turklāt e_i ir iid.

Piezīme

Šādu m tad sauc par procesa maiņas punktu.

Vidējās vērtības izmaiņas laikrindās (AMOC - at most one change modelis)

$$Y(i) = \begin{cases} \mu_1 + V(i), & 1 \leq i \leq \tilde{k}, \\ \mu_2 + V(i), & \tilde{k} \leq i \leq T \end{cases},$$

kur $V(\cdot)$ ir stacionārs process ar $EV(\cdot) = 0$, turklāt k, μ_1, μ_2 nav zināmi.

$$H_0 : \tilde{k} < T, \mu_1 \neq \mu_2 \quad H_1 : \tilde{k} = T .$$

Tad vienkāršākā, arī saukta par CUMSUM, statistika ir

$$C_T = \max_{1 \leq k \leq T} \left| \frac{1}{\sqrt{T}} \sum_{j=1}^k (Y(j) - \bar{Y}_T) \right|.$$

Vidējās vērtības izmaiņas laukrindās (AMOC - at most one change modelis)

Kritisko vērtību iegūst izmantojot butstrapa metodes frekvenču domēnā (spektrālo analīzi). Zināms, ka

$$\frac{C_T}{\tau} \rightarrow^L \sup_{0 \leq t \leq 1} |B(t)|,$$

kur $B(\cdot)$ ir Brauna tilta process un $\tau^2 = f(0)$, kur $f(\cdot)$ ir spektrālā blīvuma funkcija procesam $\{V(\cdot)\}$. Nesen izstrādāta efektīva TFT butstrapa metode šai problemātikai (2011, Kirch C., Politis D.N.).

Lietojot TFT-bootstrapu

Novērtējam $Z(t) = V(t) - \overline{V}_T$ un

$\hat{Z}(t) = Y(t) - \hat{\mu}_1 I_{\{t \leq \hat{k}\}} - \hat{\mu}_2 I_{\{t > \hat{k}\}},$

kur $\hat{k} = \operatorname{argmax}\{|\sum_{j=1}^k (Y(j) - \overline{Y}_T)| : 1 \leq k \leq T\},$

$\hat{\mu}_1 = 1/\hat{k} \sum_{j=1}^{\hat{k}} Y(j), \hat{\mu}_2 = 1/(T - \hat{k}) \sum_{j=\hat{k}+1}^T Y(j).$

Vidējo vērtību problemātika (AMOC - at most one change model)

Tad bootstrapa statistiku definējam

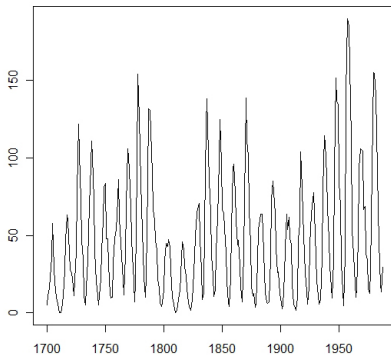
$$C_T = \left| \frac{1}{\sqrt{T}} \sum_{j=1}^k (Z^*(j)) \right|,$$

kur $\{Z^*(\cdot)\}$ ir bootstrapotie g.l., kas iegūti ar bootstrapa shēmas palīdzību.

Maiņas punkts saules aktivitātes datos (AMOC - at most one change model)

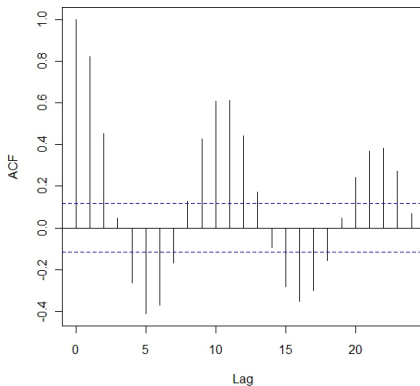
Saules plankumu (aktivitātes) dati (laikaposmā no 1700. līdz 1987.g.). Pielietojot TFT redzams, ka pastāv maiņas punkts 1935. gadā.

Sunspot Data un TFT



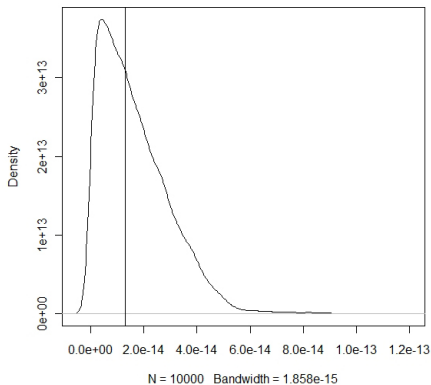
att.: Sunspot data

Sunspot Data un TFT



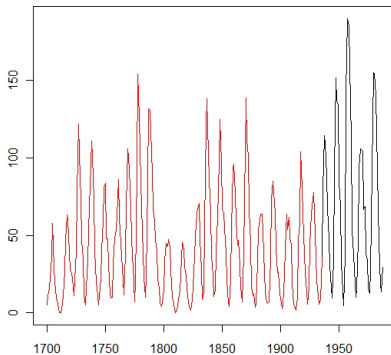
att.: Autocorrelations

Sunspot Data un TFT



att.: Test statistic

Sunspot Data un TFT



att.: Seperated data series.

Mainas punkta analīze normāla sadalījuma gadījumā
Iegūtas kritiskās vērtības un sarežģīti asimptotiskie sadalījumi.
Beijesa metode.

Mainas punkta analīze regresijas modeļos
Izmaiņas regresijas koeficientos

Mainas punkta analīze vispārīgā gadījumā ar korelāciju
atlikumos
Izmaiņas vidējā vērtībā un/vai skalācijas parametros.

Literatūra

- Antoch J. et al. Change point detection
- Kirch C., Politis D.N. TFT-bootstrap: Resampling time series in the frequency domain to obtain replicates in the time domain
- Kirch C. Resampling in the frequency domain of time series to determine critical values for change point tests.
- Taylor W.A. Change point analysis: a powerful new tool for detecting changes

Motivācija

- 1 Eksistē samērā daudz risinājumu divu izlašu lokācijas parametra starpības pārbaudei neatkarīgu datu gadījumā (Divu izlašu t-tests, Mana-Whitneya tests u.c.)
- 2 Jauni pētījumi un risinājumi neparametriskajā statistikā atkarīgu datu struktūru gadījumā (Zhang, 2011.)
- 3 Maiņas punkta analīze atkarīgiem datiem (R.Fried, 2011)
- 4 Maiņas punkta analīzes un empīriskās ticamības funkcijas (EL) metožu salīdzinājums divu izlašu lokācijas parametru starpības pārbaudei.

Nosacījumi

- Apskatīsim $X_1, X_2, \dots, X_{n1} \sim^{\text{iid}} F_1(x, \theta_0, t)$ un $Y_1, Y_2, \dots, Y_{n2} \sim^{\text{iid}} F_2(y, \theta_1, t)$
- Apskatīsim $t \rightarrow \Delta(t)$, $t \in T \subset \mathbb{R}$, $\theta(t)$ un, fiksējo t , $\theta(t) := \theta$, $\Delta(t) := \Delta$
- Visu nepieciešamo informāciju par parametriem iegūsim no nenovirzītiem funkcionāļiem:

$$E_{F_1} g_1(X, \theta_0, t) = 0, \quad E_{F_2} g_2(Y, \theta_1, t) = 0.$$

Apgalvojums

Ja $\Delta_0 = \theta_1 - \theta_0$, tad iegūstam tieši problemātiku vidējo vērtību starpībai. Apzīmējot $\theta_0 = \int x dF_1(x)$ un $\Delta_0 = \int y dF_2(y) - \int x dF_1(x)$ iegūsim funkcionāļus

$$g_1(X, \theta_0, \Delta_0, t) = X - \theta_0, \quad g_2(Y, \theta_0, \Delta_0, t) = Y - \theta_0 - \Delta_0.$$

Definīcija

Definēsim empīrisko ticamības funkciju attiecību funkciju

$$R(\Delta, \theta) = \sup_{p, q} \prod_{i=1}^{n_1} (n_1 p_i) \prod_{j=1}^{n_2} (n_2 q_j),$$

kur $p = (p_1, p_2, \dots, p_{n_1})$ un $q = (q_1, q_2, \dots, q_{n_2})$ ir sadalījuma vektori, t.i. , nenegatīvi lielumi ar summu viens, turklāt

$$\sum_{i=1}^{n_1} p_i g_1(X_i, \theta_0, t) = 0 \text{ un } \sum_{j=1}^{n_2} q_j g_2(Y_j, \theta_1, t) = 0.$$

Viens vienīgs atrisinājums eksistē, lai θ piederētu izliektajai čaulai, ko veido $g_1(X_i, \theta, \Delta, t)$ un $g_2(Y_j, \theta, \Delta, t)$. Maksimumu atrod, izmantojot Lagranža reizinātājus

Definīcija

Visbeidzot, mēs definējam empīrisko log - ticamības funkciju attiecību, kas pareizināta ar mīnuss divi

$$W(\Delta, \theta) = -2\log R(\Delta, \theta) = \\ = 2 \sum_{i=1}^{n_1} (1 + \lambda_1(\theta)g_1(X_i, \theta, \Delta, t)) + 2 \sum_{j=1}^{n_2} (1 + \lambda_2(\theta)g_2(Y_j, \theta, \Delta, t)).$$

Lai iegūtu parametru novērtējumu, maksimizē $R(\Delta, \theta)$

Apgalvojums

Tad punktveida ticamības intervāls parametram Δ katram fiksētam $t \in T$ nosakāms no $\Delta : \{R(\Delta, \hat{\theta}) > c\}$ visiem Δ_0 , kur c nosakāms no teorēmas.

Teorēma

Izpildoties gluduma nosacījumiem funkcijām $g_1, g_2, \alpha_1, \alpha_2$, eksistē tāds $\hat{\theta}$, ka $\hat{\theta}$ ir būtisks θ_0 novērtējums un $R(\Delta, \theta)$ sasniedz lokālo maksimumu pie $\hat{\theta}$, turklāt

$$\sqrt{n_1}(\hat{\theta} - \theta_0) \rightarrow_d N\left(0, \frac{\beta_1\beta_2}{\beta_2\beta_{10}^2 + k\beta_1\beta_{20}^2}\right),$$

kur $k < \infty$ ir pozitīva konstante, kura apmierina $n_2/n_1 \rightarrow k$, kad $n_1, n_2 \rightarrow \infty$, turklāt

$$-2 \log R(\Delta_0, \hat{\theta}) \rightarrow_d \chi_1^2, \quad n_1, n_2 \rightarrow \infty, \forall t \in T.$$

EL divu izlašu gadījumā - Bloku veidošana atkarīgiem datiem

Bloku izveide

- M un L veseli skaitļi, izvēlamies $Q = [(n - M)/L] + 1$, kur ar $[\cdot]$ apzīmēta veselā daļa. Tātad tiek izveidoti Q bloki no M novērojumiem katrā blokā un L atstatumu starp bloku sākumpunktiem

$$B_i = (X_{(i-1)L+1}, \dots, X_{(i-1)L+M}), \quad i = 1, \dots, Q.$$

- tad Kitamura piedāvā jau šos blokus izmantot kā novērojumus attiecīgi statistikā, proti, katram no blokiem tiek aprēķināts $T_i = \phi_M(B_i, \theta)$, kur B_i ir i -tais bloks un attēlojams $\phi_M : \mathbb{R}^M \times \theta \rightarrow \mathbb{R}$ uzdots šādā formā

$$\phi_M(B_i, \theta) = \sum_{n=1}^M g(X_{(i-1)L+n}, \theta) / M.$$

Hipotēžu pārbaude divu izlašu lokācija parametru starpībai

$$H_0 : \Delta_0 = 0, H_1 : \Delta_0 \neq 0$$

un statistiku

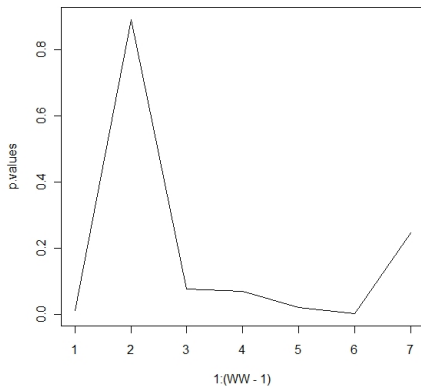
$$-2 \log R(\Delta_0, \hat{\theta}) \rightarrow_d \chi_1^2, Q_1, Q_2 \rightarrow \infty, \forall t \in T,$$

jāņem vērā, ka mēs apskatīsim statistiku bloku izlasēm.

Ideja

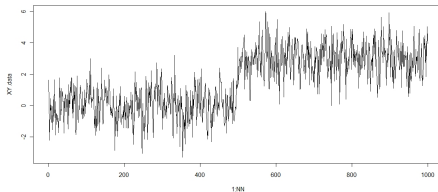
Divu izlašu vietā aplūkosim laikrindu un divus logus. Slīdošos logus uzskatīsim par izlasēm un apskatīsim p-vērtību grafikus. Punktā, kur grafiks sasniedz pietiekami mazas vērtības, laikrindā var būt maiņas punkts.

Sunspot Data un EL



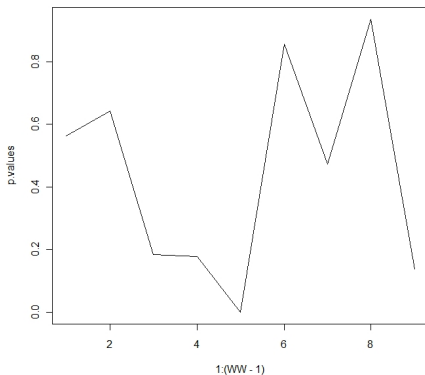
att.: Test statistic

Simulācijas AR1 ($\theta = 0.3$)



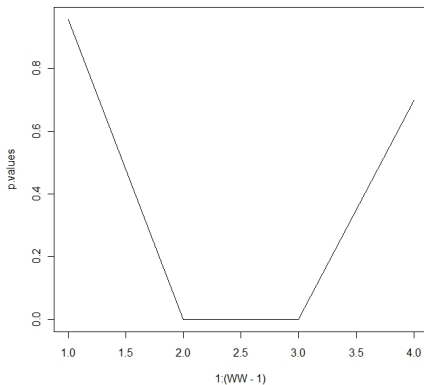
att.: Laikrinda ar shiftu

Simulācijas AR1 ($\theta = 0.3$)



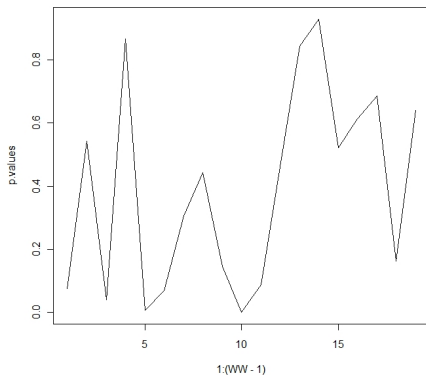
att.: Samērā laba izvēle

Simulācijas AR1 ($\theta = 0.3$)



att.: Pārāk liela loga platuma izvēle

Simulācijas AR1 ($\theta = 0.3$)



att.: Pārāk maza loga platuma izvēle

Paldies par uzmanību!

Jautājumi?