

Statistika 2 - 2022/2023, Tine Markočič

Domača naloga, 8. 8. 2023

- 1 Predvidevamo, da je vzorec realnih števil realizacija neodvisnih ponovitev slučajne spremenljivke z Lebesguovo gostoto

$$f(x; a, b^2) = \frac{1}{x\sqrt{2\pi b^2}} e^{-(\ln x - a)^2 / (2b^2)} \cdot \mathbf{1}_{(0, \infty)}(x),$$

za parameter $\vartheta = (a, b^2) \in \mathbb{R} \times (0, \infty)$.

- Poiščite kompletno zadostno statistiko.
- Obravnavajte problem cenilke največjega verjetja.
- Izračunajte cenilko $\hat{\vartheta}_{MM}$ za ϑ po metodi momentov. Komentirajte jo.
- S pomočjo metode delta obravnavajte (dvorazsežno) normalno aproksimacijo $\hat{\vartheta}_{MM}$: formulirajte primeren limitni izrek.
- S pomočjo normalne aproksimacije cenilke $\hat{\vartheta}_{MM}$ konstruirajte aproksimativno območje zaupanja stopnje zaupanja 0.95: problem prevedite na večrazsežno normalno porazdelitev z znano variančno-kovariančno matriko. Dobljeno območje čim natančneje opišite.

- 2 Predvidevamo, da je naš vzorec realizacija neodvisnih ponovitev diskretne slučajne spremenljivke z verjetnostno funkcijo

$$f(x; \vartheta) = e^{-\vartheta} \cdot \frac{\vartheta^x}{x!} \cdot \mathbf{1}_{\{0, 1, 2, 3, \dots\}}(x).$$

za realnoštevilski parameter $\vartheta \in (0, \infty) \subset \mathbb{R}$.

- Poiščite kompletno zadostno statistiko.
- Kakšni preizkusi domnev $\vartheta \leq \vartheta_0$ proti alternativam $\vartheta > \vartheta_0$ so na voljo v tem modelu?
- Za domnevo $\vartheta \leq 5$ vzemite po vašem prepričanju najboljši možen preizkus stopnje značilnosti 0.05 za vzorec vaše velikosti in skicirajte graf funkcije moči.
- Na vašem vzorcu realizirajte interval zaupanja stopnje zaupanja 0.95 po pripadajočem izreku s predavanj.

- 3 a. Privzemimo model z n zaporednimi neodvisnimi ponovitvami Bernoullijeve slučajne spremenljivke s parametrom p . Za $p_0 \in (0, 1)$ naj bo $\phi_{p_0}: \mathbb{R}^n \rightarrow [0, 1]$ preizkus za $H_0: p = p_0$ proti $A: p \neq p_0$ velikosti 0.05, ki je enakomerno najmočnejši med vsemi nepristranskimi preizkusi za H_0 proti A stopnje značilnosti 0.05. Vemo, da obstajajo taki preizkusi oblike

$$\phi_{p_0}(x_1, x_2, \dots, x_n) = \begin{cases} 1, & \sum_{i=1}^n x_i < C_1(p_0) \text{ ali } C_2(p_0) < \sum_{i=1}^n x_i, \\ \gamma_j(p_0), & \sum_{i=1}^n x_i = C_j(p_0), \\ 0, & \text{sicer.} \end{cases}$$

Vaš diagram prikazuje grafa funkcij $C_1(p_0)$ in $C_2(p_0)$ za $n = 25$. S pomočjo inverzije konstruirajte (v konkretnih številkah) interval zaupanja za Bernoullijev parameter za vaš n .

- Za $p_0 = 34/100$ izračunajte pripadajoči konstanti γ_1 in γ_2 preizkusa ϕ_{p_0} iz a.
- Naj bo C območje zaupanja za parameter $\vartheta \in \Theta$. Pokritost pri $\vartheta_0 \in \Theta$ je verjetnost $P_{\vartheta_0}(\{X | \vartheta_0 \in C(X)\})$. Koeficient zaupanja območja C je seveda natančna spodnja meja pokritosti.
Za interval zaupanja iz a. napravite graf pokritosti (potrudite se!) in ocenite koeficient zaupanja. Koeficient zaupanja lahko tudi natančno izračunate.
- Izračunajte pripadajoče konstante za različico preizkusa iz a. za velikost vzorca $n+10$ (kjer je n iz točke a.) in $p_0 = 34/100$.
- (*) Za velikost vzorca $n+10$ iz d. napravite diagram kot v a.

- 4 Obravnavamo diskretno porazdelitev na fiksnih $m+1$ točkah ξ_0, \dots, ξ_m z verjetnostmi (p_0, p_1, \dots, p_m) . Preizkušamo domnevo čisto določene porazdelitve $H_0: (p_0, p_1, p_2, p_3, p_4) = (\frac{1}{12}, \frac{1}{12}, \frac{1}{3}, \frac{1}{4}, \frac{1}{4})$ z asimptotičnimi preizkusi pri nominalni stopnji značilnosti 0.05. Za vzorce velikosti $n = 40, 60, 80, 100$ izračunajte:

- eksaktno velikost preizkusa domneve H_0 na podlagi razmerja verjetij,
- eksaktno velikost preizkusa domneve H_0 na podlagi Pearsonove statistike $\sum_j \frac{(T_j - n\pi_j)^2}{n\pi_j}$ (za vaše podatke).

Opozorilo: Če vaši rezultati močno odstopajo od nominalne značilnosti in ste prepričani, da vaš algoritem deluje, preverite, ali je njegova implementacija numerično zanesljiva.

- (*) Obravnavajte moč na alternativah oblike $p_j = \pi_j + \delta, p_k = \pi_k - \delta$ (za vaše podatke).

- 5 Med različnimi holesteroli je holesterol LDL (*low-density lipoprotein* - lipoprotein nizke gostote) relativno težko oziroma drago meriti. Zato je še marsikje v uporabi cenilka oblike $LDL = \beta_1 \cdot TCH + \beta_2 \cdot HDL + \beta_3 \cdot TRI$ (za konkretne vrednosti parametrov $\beta_1, \beta_2, \beta_3$; ocena je znana pod imenom Friedewaldova formula), kjer so TCH (*total cholesterol* - skupni holesterol), HDL (*high-density lipoprotein* - lipoprotein visoke gostote) in TRI (trigliceridi) količine, ki jih je relativno lahko meriti. Ocenjevanje parametrov β_i je torej regresijski problem.

- Predpostavite linearni model $LDL_i = \beta_0 + \beta_1 \cdot TCH_i + \beta_2 \cdot HDL_i + \beta_3 \cdot TRI_i + \varepsilon_i$, kjer so ε_i neodvisne slučajne spremenljivke, vse porazdeljene po zakonu $N(0, \sigma^2)$.
 - Ocenite parametre $\beta_0, \beta_1, \beta_2, \beta_3$ po metodi najmanjših kvadratov.
 - Preizkusite domnevo $\beta_1 = \beta_2 = \beta_3 = 0$ na standardni način pri stopnji značilnosti 0.05.
 - Preizkusite domnevo $\beta_0 = 0$ na standardni način pri stopnji značilnosti 0.05.
- Predpostavite model brez prostega člena $LDL_i = \beta_1 \cdot TCH_i + \beta_2 \cdot HDL_i + \beta_3 \cdot TRI_i + \varepsilon_i$.
 - Ocenite parametre $\beta_1, \beta_2, \beta_3$ po metodi najmanjših kvadratov.
 - Preizkusite domnevo $(\beta_1, \beta_2, \beta_3) = (1, -1, -0.45)$.
 - Realizirajte območje zaupanja za $(\beta_1, \beta_2, \beta_3)$ stopnje zaupanja 0.95 in ga čim natančneje opišite.
 - Realizirajte hkratne intervale zaupanja za β_1, β_2 in β_3 z Bonferronijevim popravkom in primerjajte dobljeni kvader zaupanja z območjem iz prejšnje točke.

- 6 Za dani avtomobil želimo na podlagi podatka o učinkovitosti mpg („miles per gallon“) in mase weight pojasniti, ali gre za avtomobil „tujega“ izvora ($foreign=1$) ali za avtomobil „domačega“ izvora ($foreign=0$). Za dani vzorec, v katerem je 15 primerkov tujega in 85 primerkov domačega izvora predpostavimo, da je realizacija slučajnega vektorja iz modela

$$p_i = P(foreign_i = 1) = 1 - \exp(-\exp(\beta_0 + \beta_1 * weight_i + \beta_2 * mpg_i)),$$

kjer so komponente foreign_i neodvisne in porazdeljene po zakonu $B(1, p_i)$.

- Zapišite funkcijo verjetja.
- Ocenite parametre β_0, β_1 in β_2 po metodi največjega verjetja.
- Izračunajte Fisherjevo informacijsko matriko.
- Izračunajte standardne napake za parametre β_0, β_1 in β_2 .
- Na standardni način preizkusite domnevo $H_0 : \beta_1 = \beta_2 = 0$.

7 Predpostavimo linearni regresijski model (za neskončni vzorec) oblike

$$\mathbb{X} = Z \cdot \beta + \varepsilon,$$

kjer je $Z \in \mathbb{R}^{\infty \times d}$ fiksna matrika z neskončno mnogo vrsticami in so komponente ε_i vektorja ε neodvisne in enako porazdeljene slučajne spremenljivke s končno (neznano) disperzijo σ^2 in pričakovano vrednostjo 0 in je $\beta = [\beta_1 \ \dots \ \beta_d]^T$ stolpec (preostalih) parametrov. Naj Z_i označuje i -to vrstico matrike Z in naj bo $Z^{(n)} \in \mathbb{R}^{n \times d}$ matrika, sestavljena iz prvih n vrstic matrike Z . Analogno naj bo $\mathbb{X}^{(n)}$ stolpec, sestavljen iz prvih n komponent vektorja \mathbb{X} . Privzemite, da imajo matrike $Z^{(n)T} Z^{(n)}$ poln rang za vsa dovolj velika števila n in naj bo

$$\hat{\beta}^{(n)} = (Z^{(n)T} Z^{(n)})^{-1} Z^{(n)T} \mathbb{X}^{(n)}$$

cenilka za β (po metodi najmanjših kvadratov) za restrikcijo na vzorec velikosti n . Naj bo

$$\zeta^{(n)} = \frac{1}{\sigma} (Z^{(n)T} Z^{(n)})^{-1/2} Z^{(n)T} \varepsilon^{(n)}$$

standardizacija (prepričajte se, da to drži) cenilke $\hat{\beta}^{(n)}$.

a. Dokažite, da pri pogoju

$$\lim_{n \rightarrow \infty} \max_{1 \leq j \leq n} \left\langle (Z^{(n)T} Z^{(n)})^{-1} Z_j^T, Z_j^T \right\rangle = 0$$

vektorji $\zeta^{(n)}$ konvergirajo k d -razsežni standardni normalni porazdelitvi.

- Dokažite, da je pogoj iz prejšnje točke izpolnjen, če zaporedje $d \times d$ -matrik $Z^{(n)T} Z^{(n)} / n$ konvergira k neki pozitivno definitni matriki P in za neko število $M > 0$ velja $\|Z_j\| \leq M$ za vsa naravna števila j .
- Kaj pomeni pogoj iz točke b. v primeru enostavne linearne regresije ($X_i = \beta_0 + \beta_1 \cdot z_i + \varepsilon_i$)?
- Ali je pogoj iz točke b. v praksi razumen? (Pri razmišljanju o tem bodite inventivni; lahko se odločite za bolj praktično ali bolj teoretično utemeljevanje.)
- (*) Pokažite, da je pogoj omejenosti $\|Z_j\| \leq M$ iz točke b. v resnici odveč. Splošni primer dovolj preprosto sledi iz posebnega primera enostavne linearne regresije.

Za točko a. si lahko pomagata z naslednjim izrekom: naj bodo za vsako naravno število n slučajne spremenljivke $Y_1^{(n)}, \dots, Y_n^{(n)}$ neodvisne in enako porazdeljene s pričakovano vrednostjo 0 in disperzijo 1. Dalje naj bodo $\mathbf{c}^{(n)} \in \mathbb{R}^n$ enotski vektorji z lastnostjo

$$\lim_{n \rightarrow \infty} \max_{1 \leq j \leq n} (\mathbf{c}_j^{(n)})^2 = 0.$$

Tedaj zaporedje slučajnih spremenljivk $\sum_{j=1}^n \mathbf{c}_j^{(n)} Y_j^{(n)}$ konvergira k standardni normalni porazdelitvi.