

Konkuruojančių rizikų regresinis Farlie–Gumbel–Morgenstern skirstinių šeimos išgyvenamumo modelis

Audronė JAKAITIENĖ (VDU), Danas ZUOKAS (VDU)

el. paštas: iiauja@vdu.lt

Ivadas

Konkuruojančių rizikų teorija nagrinėja objektus, kurių gedimas (techninio elemento, sistemos atveju) ar mirtis (gyvo organizmo atveju) gali priklausyti nuo kelių priežasčių. Pastarosios vadinamos konkuruojančiomis rizikomis. Statistinė konkuruojančių rizikų teorija nagrinėja cenzūruotas imtis. Norima nustatyti, kokios rizikos labiausiai įtakoja gedimą (mirčių) intensyvumą ir laiką. Išgyvenamumo uždaviniuose konkuruojančios rizikos susijusios su populiacijos individų mirtingumo tyrimais, kai visi individai yra veikiami tų pačių mirties priežasčių (konkuruojančių rizikų). Pagrindinis tokų tyrimų tikslas yra išskirti duotos priežasties ar priežasčių kompleksą poveikį populiacijai [2].

Išgyvenamumo funkcija

Tarkime, kad populiaciją vienu metu veikia m mirties priežasčių C_1, \dots, C_m . Kiekvieną mirties priežastį atitinka hipotetinis (potencialus) mirties laikas $T_1 \geq 0, \dots, T_m \geq 0$. Išgyvenamumo funkcija vektoriui $(T_1, \dots, T_m)'$ (daugiamatė išgyvenamumo funkcija) apibrėžiama taip:

$$S_{1\dots m}(t_1, \dots, t_m) = P\{T_1 > t_1, \dots, T_m > t_m\}, \quad (1)$$

kur

$$S_j(t_j) = 1 - F_j(t_j) = S_{1\dots m}(0, \dots, t_j, \dots, 0) \quad (2)$$

yra marginalinės išgyvenamumo funkcijos, o $F_j(t_j)$ – marginalinė pasiskirstymo funkcija. Čia ir toliau $j = 1, 2, \dots, m$.

Konkretus stebimas mirties laikas yra $T = \min(T_1, \dots, T_m)$. Todėl bendra išgyvenamumo funkcija apibrėžiama taip (yra vieno kintamojo funkcija):

$$S_T(t) = P\{T > t\} = S_{1\dots m}(t, \dots, t). \quad (3)$$

Taip apibrėžtos bendros išgyvenamumo funkcijos bendra (overall) rizikos funkcija yra

$$h_T(t) = -\frac{d \log S_T(t)}{dt}, \quad (4)$$

o atskiros (crude) rizikos funkcijos yra [4]

$$h_j(t_1, \dots, t_m) = -\left. \frac{\partial \log S_{1\dots m}(t_1, \dots, t_m)}{\partial t_j} \right|_{t_1=\dots=t_m=t}. \quad (5)$$

FGM skirstinių šeimos išgyvenamumo funkcija

Šiame darbe daugiamate išgyvenamumo funkcija buvo pasirinkta trimatė Farlie–Gumbel–Morgenstern (FGM) skirstinių šeimos išgyvenamumo funkcija:

$$\begin{aligned} S_{123}(t_1, t_2, t_3) = & S_1(t_1) S_2(t_2) S_3(t_3) [1 + a_{12} F_1(t_1) F_2(t_2) \\ & + a_{13} F_1(t_1) F_3(t_3) + a_{23} F_2(t_2) F_3(t_3)]. \end{aligned} \quad (6)$$

a_{12} , a_{13} ir a_{23} yra priklausomumo parametrai tarp marginaliniai skirstiniai aprašyti atsitiktinių dydžių. Jie tenkina keturių nelygybių sistemą [3]:

$$\begin{cases} 1 + a_{12} + a_{13} + a_{23} \geq 0, \\ 1 + a_{12} - a_{13} - a_{23} \geq 0, \\ 1 - a_{12} + a_{13} - a_{23} \geq 0, \\ 1 - a_{12} - a_{13} + a_{23} \geq 0. \end{cases} \quad (7)$$

Atsitiktinių dydžių T_1 , T_2 , T_3 marginaliniai skirstiniai buvo pasirinkti nupjautieji iš kairės (apibrėžti neneigiamiems atsitiktiniam dydžiams $T_1 \geq 0$, $T_2 \geq 0$, $T_3 \geq 0$) logistiniai skirstiniai, kurių išgyvenamumo funkcija yra:

$$S_j(t_j, Z) = \frac{\exp \left\{ -\frac{t_j}{\beta_j Z'} \right\} \left(1 + \exp \left\{ \frac{\alpha_j Z'}{\beta_j Z'} \right\} \right)}{1 + \exp \left\{ -\frac{t_j - \alpha_j Z'}{\beta_j Z'} \right\}}, \quad (8)$$

$\alpha_j Z' = (\alpha_{j0} \dots \alpha_{jk})(1 Z_1 \dots Z_k)' = \alpha_{j0} + \alpha_{j1}Z_1 + \dots + \alpha_{jk}Z_k \geq 0$ yra tiesinė regresorių, nusakančių atsitiktinio dydžio vietą, kombinacija, o $\beta_j Z' = (\beta_{j0} \dots \beta_{jk})(1 Z_1 \dots Z_k)' = \beta_{j0} + \beta_{j1}Z_1 + \dots + \beta_{jk}Z_k > 0$ – tiesinė regresorių, nusakančių atsitiktinio dydžio masteli, kombinacija. Todėl sakoma, kad išgyvenamumo modelis yra regresinis.

Modelio parametru įvertinimas

Nežinomų modelio parametrų (marginalinių funkcijų vietas ir mastelio, ir ryšio tarp atsитiktinių dydžių parametrus) įvertinimui taikomas didžiausio tikėtinumo metodas. Įveskime cenzūravimo indikatorių:

$$\delta_j = \begin{cases} 1, & \text{jei } T \leq T^+, \text{ ir } T = T_j, \text{ t.y., mirties priežastis yra } j, X = T_j, \\ 0, & \text{jei stebėjimas yra cenzūruotas, } X = T^+, \end{cases} \quad (9)$$

čia T^+ yra cenzūravimo laikas. Taigi stebima N tūrio imtis, o kiekvienas imties objektas turi stebėjimo duomenis $(X_i, \delta_{j_i}, j_i, Z_i)$, $X_i = \min(T_i, T_i^+)$ – mirties arba cenzūravimo laikas, δ_{j_i} cenzūravimo indikatorius, j_i – mirties priežasties indeksas ir Z_i – regresorių vektorius ($j_i = 1, 2, \dots, m$, $i = 1, 2, \dots, N$). Naudojant išgyvenamumo ir rizikos funkcijas sukonstruojama tikėtinumo funkcija:

$$\begin{aligned} L(X_1, \dots, X_N, \delta_{j_1}, \dots, \delta_{j_N}, Z_i, \theta) \\ = \prod_{i=1}^N (h_{j_i}(X_i, Z_i, \theta))^{\delta_{j_i}} S_T(X_i, Z_i, \theta). \end{aligned} \quad (10)$$

Darbe pasiūlytas ir kitas tikėtinumo funkcijos konstravimo būdas:

$$\begin{aligned} L(X_1, \dots, X_N, \delta_{j_1}, \dots, \delta_{j_N}, Z_i, \theta) \\ = \prod_{i=1}^N (f_{j_i}(X_i, Z_i, \theta) \bar{S}_{j_i}(X_i, Z_i, \theta))^{\delta_{j_i}} S_T(X_i, Z_i, \theta)^{1-\delta_{j_i}}, \end{aligned} \quad (11)$$

kur $f_j(X_i, Z_i, \theta)$ – (2) formule išreikštą išgyvenamumo funkciją atitinkantis vienmatis marginalinis j -osios mirties priežasties potencialaus mirties laiko pasiskirstymo tankis, apskaičiuotas i -ojo stebėjimo laiku X_i . θ – vertinamų parametru vektorius (priklausomumo (7), padėties ir mastelio (8)). $\bar{S}_{j_i}(X_i, Z_i, \theta)$, pavadinimė ją “salygine” išgyvenamumo funkcija, apibrėžiama taip:

$$\begin{aligned} \bar{S}_j(t_j) &= P\{T_k > t_k, k \in \{1, \dots, m\} \setminus \{j\} \mid T_j \leq t_j\} \\ &= \frac{\int_{t_j}^{\infty} \int_{t_j}^{\infty} \int_0^{t_j} f_{1\dots m}(u_1, \dots, u_m) du_k du_j}{\int_0^{t_j} f_j(u_j) du_j}, \quad (k \in \{1, \dots, m\} \setminus \{j\}). \end{aligned} \quad (12)$$

Čia $f_{1\dots m}(t_1, \dots, t_m)$ yra daugiamatis, (1) formule išreikštą daugiamatę pasiskirstymo funkciją atitinkantis, potencialių mirties laikų vektoriaus pasiskirstymo tankis. Tokio tikėtinumo funkcijos konstravimo būdo prasmė yra ta, kad atsižvelgiama ne tik į tai, kad objektas mirė nuo j -osios priežasties, bet ir į tai, kad jis nemirė nuo likusių priežasčių, t.y., yra cenzūruotas stebėjimas kitų mirties priežasčių atžvilgiu.

Praktinis modelio pritaikymas

Tikėtinumo funkcijos sprendinys ieškotas skaitmeniškai, naudojant BFGS optimizavimo algoritmą su apribojimais [1]. Yra žinoma, kad didžiausio tikėtinumo metodu rastas įvertis, esant atsitiktiniams cenzūravimui, asymptotiskai pasiskirstęs pagal daugiamati nor maluji dėsnį su vidurkiu θ (vertinamų parametru vektorius) ir kovariacine matrica $I^{-1}(\theta)$, t.y.,

$$\hat{\theta} \sim N(\theta, I^{-1}(\theta)) \quad \theta \in \mathbb{R}^p. \quad (13)$$

I^{-1} , Fišerio informacinės matricos atvirkštinė, maksimumo taške nusako kovariaciją tarp maksimalaus tikėtinumo įverčio $\hat{\theta}$ komponenčių. Tikrinant hipotezę apie parametrus $H_0: \hat{\theta} = \theta_0$, buvo naudojamas Neyman–Pearson–Wilks kriterijus

$$k = -2 \log \frac{L(\theta_0)}{L(\hat{\theta})}, \quad (14)$$

ir yra žinoma, kad jis turi χ_p^2 asymptotinį pasiskirstymą [4].

Rezultatai

Šiame darbe buvo sudarytas trimatis regresiniis konkuruojančių rizikų išgyvenamumo modelis FGM skirstinių šeimai su logistiniais iš kairės nupjautaisiais marginaliniai skirstiniai. Pasiūlytas tikėtinumo funkcijos konstravimo būdas išgyvenamumo modelio parametrams įvertinti. Sudarytas išgyvenamumo modelis pritaikytas 1972 metais pradėtos

1 lentelė

Stebėtų mirčių bei pagal modelį dviem variantais ((10) ir (11))
paskaičiuotų mirčių skaičius

Mirties priežastis	Koronarinės ligos	Vėžys	Kitos	Viso
Pirmas variantas, (10) formulė				
Stebėtų mirčių skaičius	412	247	201	860
Modeliavimo rezultatai	533	353	239	1125
Skirtumas	121	106	38	265
% nuo stebėtų mirčių	29	43	19	31
Antras variantas, (11) formulė				
Stebėtų mirčių skaičius	412	247	201	860
Modeliavimo rezultatai	469	320	245	1034
Skirtumas	57	73	44	174
% nuo stebėtų mirčių	14	30	22	20

ir 22 metus trukusios Pasaulio sveikatos organizacijos koordinuotos, Kauno–Roterdamo profilaktinės metodinės programos, pavadintos “KRIS” (The Kaunas Rotterdam Intervention Study), metu surinktiems duomenims.

1 lentelėje yra pateikti stebėtų mirčių nuo kiekvienos priežasties bei pagal modelį paskaičiuotų mirčių nuo kiekvienos priežasties skaičiai abiem variantais ((10) ir (11)). Modeliavimo mirčių skaičius rastas pagal tokią formulę:

$$\text{MS}_j = \sum_{i=1}^N F_{j_i}(X_i) \mathbf{1}\{j_i = j\}, \quad j = 1, 2, 3, \quad (15)$$

t.y., buvo susumuotos visų stebėjimų tikimybės mirti nuo konkrečios priežasties per visą stebėjimo laikotarpį. Tokiu būdu gautas tam tikras skaičius, nusakantis kiekvienos mirties dažnį.

Iš 1 lentelėje pateiktų rezultatų matome, kad pagal pasiūlytą variantą (ivedant “salyginę” išgyvenamumo funkciją (11)) sukonstruotos tikėtinumo funkcijos maksimizavimas duoda geresnius rezultatus nei pagal literatūroje rastą variantą (naudojant išgyvenamumo ir rizikos funkcijas (10)) sukonstruotos tikėtinumo funkcijos maksimizavimas, taip pat prasme, kad pagal modelį paskaičiuotų mirčių skaičius pirmu atveju buvo artimesnis stebėtų mirčių skaičiui nuo koronarinių ligų ir vėžio, ir šiek tiek didesnis mirtims nuo kitų priežasčių.

Literatūra

- [1] T. Coleman, M.A. Branch, A. Grace, *Optimization Toolbox User's Guide*, MathWorks (1990–1999).
- [2] A. Jakaitienė, *Konkuruojančių riziku regresinių modelių algoritmai*, Daktaro disertacija, Kaunas (2001).
- [3] S. Kotz, N. Balakrishnan, N.L. Johnson, *Continuous Multivariate Distributions: Models and Applications*, Vol. 1, John Wiley and Sons, New York (2000).
- [4] R.G. Miller, Jr., *Survival Analysis*, John Wiley and Sons, New York (1981).

The survival regression model of competing risks for the family of Farlie–Gumbel–Morgenstern distributions

A. Jakaitienė, D. Zuokas

In this paper the trivariate survival regression model for FGM family of distributions is constructed with marginal left-truncated logistic distributions. Two methods (using survival and hazard functions in the first case, and distributional density and “conditional” survival function in the second case) are used when constructing likelihood function for model parameter estimation. Constructed survival model was run with the data of the “KRIS” (The Kaunas Rotterdam Intervention Study), which lasted for 22 years from 1972. The results show, that using second case for likelihood function construction gives better approximation for the data, because some additional conditions are considered.