

DAUGIAMATIS RETŲ ĮVYKIŲ TIKIMYBIŲ VERTINIMO ALGORITMAS

Leonidas Sakalauskas, Ingrida Vaičiulytė

Vilniaus universitetas, Matematikos ir informatikos institutas,
Akademijos g. 4, LT-08663, Vilnius
sakal@ktl.mii.lt, ingrida_vaiciulyte@yahoo.com

Anotacija. Šiame darbe sudarytas Monte–Karlo Markovo grandinių algoritmas daugiamačių retų įvykių dažniams vertinti. Retų įvykių tikimybės logitai modeliuojami Puasono skirstiniu, kurio parametrai pasiskirstę pagal daugiamatį normalųjį dėsnį su nežinomais vidurkių vektoriumi ir kovariacijų matrica. Nežinomų parametrų įverčiai gaunami didžiausio tikėtino metodo. Išvestos lygtys, kurias turi tenkinti modelio didžiausio tikėtino parametrų įverčiai. Vertinamos kovariacijų matricos teigiamas apibrėžtumas kontroliuojamas apskaičiuojant matricos didžiausios ir mažiausios tikrinių reikšmių santykį.

Pagrindiniai žodžiai: Monte–Karlo Markovo grandinė, Bajeso metodas, tikėtino metodo, Puasono–Gauso modelis.

Įvadas

Bajeso metodas plačiai taikomas įvairiems sprendimo priėmimo uždaviniams spręsti bei verslo ar finansinių rodiklių vertinimui (Bradley ir Thomas, 2000, Clayton ir Kaldor, 1987, Chen, 2009; Liseo ir Loperfido, 2003; Tsutakawa ir kt. 1985). Monte–Karlo Markovo grandinės (Markov Chain Monte Carlo, MCMC) yra kompiuterinio imitavimo būdas dažnai taikomas nežinomiems parametrams įvertinti. Šiame darbe pasinaudojant Puasono–Gauso modeliu konstruojamas Monte–Karlo Markovo sekų algoritmas, skirtas kelių retų įvykių tikimybėms vertinti. Pavyzdžiui, gali būti nagrinėjami draudiminiai įvykiai (skirtingų automobilių rūšių avarių skaičius), susirgimų ar mirties atvejai populiacijoje. Apsteriorinis įvykių dažnių skirstinys yra sudaromas Bajeso metodu. Darbe nagrinėjamas logit modelis, kuriame nepriklausomas kintamasis α išreiškiamas per nagrinėjamo įvykio tikimybę P :

$$\alpha = \ln \frac{P}{1-P}, \quad (1)$$

čia P yra tikimybė, kad priklausomas kintamasis įgys reikšmę 1, o $(1-P)$ – tikimybė, kad priklausomas kintamasis įgys reikšmę 0 (Altaleb ir Chauveau, 2002; Pearce, 2006). Vieno įvykio vertinimo algoritmas panaudojant Bajeso metodą buvo sudarytas Bradley ir Thomas (2000), Clayton ir Kaldor (1987), Tsutakawa ir kt. (1985).

1. Puasono–Gauso modelis

Tegul turime generalinę imtį $\Theta = (\theta_1, \theta_2, \dots, \theta_K)$, sudarytą iš K populiacijų ir kiekviena populiacija θ_j turi N_j individų, $j = \overline{1, K}$. Tarkime, kad stebėjimo metu populiacijoje gali įvykti tam tikri nepalankūs įvykiai (draudiminio įvykio, susirgimo, mirties,...). Tiek

statistinio modelio parametrai, tiek stebėjimai iš esmės nesiskiria – abu yra atsitiktiniai kintamieji. Tikslas yra įvertinti nepalankių įvykių tikimybes P_j^m , kai Y_j^m yra stebėtas m -tojo įvykio pasirodymų skaičius, $j = \overline{1, K}$ ir $m = \overline{1, M}$. Dėl didelių populiacijos apimties N_j skirtumų paprastas santykinis įvertis $\frac{Y_j^m}{N_j}$ ne visada tinka. Todėl nagrinėjamas empirinis Bajeso metodas, taikant retų įvykių tikimybėms modeliuoti empirinį Bajeso Puasono–Gauso modelį.

Empiriniame Bajeso metode priimama, kad m -tojo įvykio j -tojoje populiacijoje skaičius Y_j^m pasiskirstęs pagal Puasono dėsnį su parametru $\lambda_j^m = N_j \cdot P_j^m$, $1 \leq m \leq M$, t.y., su tankiu (Bradley ir Thomas, 2000; Clayton ir Kaldor, 1987; Tsutakawa ir kt., 1985):

$$f(Y_j^m, \lambda_j^m) = e^{-\lambda_j^m} \frac{(\lambda_j^m)^{Y_j^m}}{(Y_j^m)!} \quad j = 1, \dots, K. \quad (2)$$

Puasono–Gauso modelyje priimama, kad įvykių tikimybių logitai populiacijose yra pasiskirstę pagal daugiamačią Gauso dėsnį su parametrais μ, Ω (Bradley ir Thomas, 2000), t.y., logit skirstinio (1) tankis

$$g(\alpha, \mu, \Omega) = \frac{\exp\left(-(\alpha - \mu)^T \Omega^{-1} (\alpha - \mu)\right)}{\sqrt{|\Omega|} \cdot (2\pi)^{\frac{M}{2}}}. \quad (3)$$

Tikimybės P_j^m įvertis yra apskaičiuojamas aposterioriniu vidurkiu

$$P_j^m = \frac{\int_{-\infty}^{+\infty} \frac{1}{1 + e^{-\alpha_m}} \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha}{D_j(\mu, \Omega)}, \quad (4)$$

čia

$$D_j(\mu, \Omega) = \int_{-\infty}^{+\infty} \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha, \quad j = \overline{1, K}, \quad m = \overline{1, M} \quad (5)$$

Taikant Bajeso metodą statistikoje dažnai tenka minimizuoti tam tikras funkcijas, išreikštas per aposteriorinio tankio integralą. Taigi, tegul nežinomi parametrai μ, Ω yra vertinami didžiausio tikėtimumo metodu (Bradley ir Thomas, 2000; Tsutakawa ir kt., 1985). Gaunama, kad logaritminė tikėtimumo funkcija yra:

$$L(\mu, \Omega) = -\sum_{j=1}^K \ln \left(\int_{-\infty}^{+\infty} \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha \right) = -\sum_{j=1}^K \ln(D_j(\mu, \Omega)), \quad (6)$$

kurią minimizavus gaunami parametrai μ, Ω įverčiai.

2. Didžiausio tikėtimumo funkcijos išvestinės

Tikėtimumo funkcija (6) gali būti diferencijuojama daugelį kartų pagal parametrus μ, Ω . Nesunku įsitikinti, kad šios funkcijos pirmos eilės išvestinės yra:

$$\frac{\partial L(\mu, \Omega)}{\partial \mu} = \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} \Sigma^{-1}(\alpha - \mu) \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha}{D_j(\mu, \Omega)}, \quad (7)$$

$$\frac{\partial L(\mu, \Omega)}{\partial \Omega} = - \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} (\Omega^{-1} - \Omega^{-1}(\alpha - \mu)(\alpha - \mu)^T \Omega^{-1}) \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha}{D_j(\mu, \Omega)} \quad (8)$$

Išvestines (7), (8) prilyginus nuliui, gaunamos lygtys, kurias tenkina Puasono–Gauso modelio parametrų įverčiai:

$$\bar{\mu} = \frac{1}{K} \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} \alpha \cdot \prod_{m=1}^M f\left(Y_j^m, \frac{N_j}{1 + e^{-\alpha_m}}\right) g(\alpha, \mu, \Omega) d\alpha}{D_{j,k}^m(\mu, \Omega)}, \quad (9)$$

$$\bar{\Omega} = \frac{1}{K} \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} (\alpha - \mu)(\alpha - \mu)^T \prod_{m=1}^M f\left(Y_{j,k}^m, \frac{N_j}{1 + e^{-\alpha}}\right) g(\alpha, \mu, \Omega) d\alpha}{D_{j,k}^m(\mu, \Omega)} \quad (10)$$

3. Puasono–Gauso modelio parametrų įverčiai

Didžiausio tikėtinumo įverčius $\bar{\mu}, \bar{\Omega}$ galima apskaičiuoti kvazi-Niutono metodu (Dennis ir Schnabel, 1996), pasinaudojus (6), (7), (8) išraiškomis. Be to, didžiausio tikėtinumo funkcijos minimizavimą ir integravimą galima atlikti pasinaudojus matematine programine įranga MATHCAD, MAPLE ir pan. Panagrinėsime „fiksuoto taško iteracijų“ metodą parametrų μ, Ω didžiausio tikėtinumo įverčiams rasti, pasinaudojus lygybėmis (9), (10) (Kantorovich ir Akilov, 1982):

$$\mu_{t+1} = \frac{1}{K} \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} \alpha \cdot f\left(Y_j, \frac{N_j}{1 + e^{-\alpha}}\right) g(\alpha, \mu_t, \Omega_t) d\alpha}{D_j(\Omega_t, \Sigma_t)}, \quad (11)$$

$$\Omega_{t+1} = \frac{1}{K} \frac{\sum_{j=1}^K \int_{-\infty}^{+\infty} (\alpha - \mu_t)(\alpha - \mu_t)^T f\left(Y_j, \frac{N_j}{1 + e^{-\alpha}}\right) g(\alpha, \mu_t, \Omega_t) d\alpha}{D_j(\mu_t, \Omega_t)}. \quad (12)$$

Integralus (4)–(6), (11)–(12) galima apskaičiuoti pasinaudojus Ermito–Gauso kvadratūrinėmis formulėmis (Abramowitz ir Stegun, 1964). Lygybėse (11), (12) galima paimti tokį pradinį tašką (μ_0, Ω_0) :

$$\mu_0 = P, \quad (13)$$

$$\Omega_0 = \frac{\sum_{j=1}^K (Y_j - N_j P)^T \cdot (Y_j - N_j P)}{\sum_{j=1}^K N_j}, \quad (14)$$

$$\text{čia } P = \left(\frac{\sum_{j=1}^K Y_j^1}{\sum_{j=1}^K N_j}, \frac{\sum_{j=1}^K Y_j^2}{\sum_{j=1}^K N_j}, \dots, \frac{\sum_{j=1}^K Y_j^M}{\sum_{j=1}^K N_j} \right), Y_j = (Y_j^1, Y_j^2, \dots, Y_j^M).$$

4. MCMC algoritmas

„Fiksuoto taško“ algoritmą (11), (12) galima realizuoti Monte–Karlo Markovo grandinių metodu. Tegul, sugeneruota t grandžių ir kiekvienoje apskaičiuoti įverčiai μ_t, Ω_t , paėmus pradines reikšmes (13), (14). Kiekvienoje grandyje generuojami daugiamačiai Gauso vektoriai $\alpha_{j,k} \sim N(\mu^t, \Omega^t)$, $k=1, \dots, N^t$, čia N^t yra imties tūris t -toje grandyje. Siekiant išvengti skaičiuojamųjų problemų, galinčių atsirasti dėl labai mažų tarpinių rezultatų reikšmių, skaičiuojant integralus tiesiog pagal formules (4)–(12) įvedama pagalbinė funkcija

$$r_j(\alpha) = \ln \left(\prod_{m=1}^M f_j(Y_j^m, \frac{N_j}{1+e^{-\alpha_m}}) / \prod_{m=1}^M f_j(Y_j^m, \frac{N_j}{1+e^{-\mu_m}}) \right) = \sum_{m=1}^M \frac{N_j (e^{-\alpha_m} - e^{-\mu_m})}{(1+e^{-\mu_m})(1+e^{-\alpha_m})} + Y_j^m \cdot \ln \left(\frac{1+e^{-\mu_m}}{1+e^{-\alpha_m}} \right). \quad (15)$$

Toliau apskaičiuojamos sumos:

$$\tilde{D}_j^t = \sum_{k=1}^{N^t} r_j(\alpha_{j,k}), \quad (16)$$

$$\tilde{D}2_j^t = \sum_{k=1}^{N^t} \left(r_j(\alpha_{j,k}) - \frac{\tilde{D}_j^t}{N^t} \right)^2, \quad (17)$$

$$\tilde{m}_j^t = \sum_{k=1}^{N^t} \alpha_{j,k} \cdot r(\alpha_{j,k}), \quad (18)$$

$$\tilde{S}_j^t = \sum_{k=1}^{N^t} (\alpha_{j,k} - \tilde{m}_j^t) \cdot (\alpha_{j,k} - \tilde{m}_j^t)^T \cdot r(\alpha_{j,k}), \quad (19)$$

$$P_{j,m}^t = \sum_{k=1}^{N^t} \frac{r(\alpha_{j,k})}{1+e^{-\alpha_{j,k,m}}}, \quad (20)$$

pagal kurias gaunami sekančios iteracijos įverčiai:

$$\mu^{t+1} = \frac{1}{K} \sum_{j=1}^K \frac{\tilde{m}_j^t}{\tilde{D}_j^t}, \quad (21)$$

$$\Sigma^{t+1} = \frac{1}{K} \sum_{j=1}^K \frac{\tilde{S}_j^t}{\tilde{D}_j^t}, \quad (22)$$

tikėtimumo funkcijos įvertis:

$$L^t = -\sum_{j=1}^K \ln \tilde{D}_j^t, \quad (23)$$

jos imties dispersijos įvertis:

$$d^t = \sum_{j=1}^K \left(\frac{\tilde{D}2_j^t \cdot N^t}{(\tilde{D}_j^t)^2} - 1 \right), \quad (24)$$

įvykių tikimybių populiacijose įverčiai:

$$\tilde{P}_{j,m}^t = \frac{P_{j,m}^t}{\tilde{D}_j^t}. \quad (25)$$

Algoritmo stabdymui įvestas kriterijus, skirtas patikrinti hipotezei apie vidurkių ir kovariacijų matricų sutapimą dviejose gretimose iteracijose:

$$H^t = \frac{K}{\frac{1}{K} \sum_{j=1}^K \frac{\tilde{D}2_j^t}{(\tilde{D}_j^t)^2}} \cdot \left(\ln \frac{|\Omega^{k+1}|}{|\Omega^k|} + SP(\Omega^{k+1} \cdot (\Omega^k)^{-1}) + (\mu^{k+1} - \mu^k)^T \cdot (\Omega^k)^{-1} \cdot (\mu^{k+1} - \mu^k) - M \right) \quad (26)$$

Eksperimentiniai tyrimai ir teorinė analizė parodė, kad esant pakankamai dideliam imties tūriui tikėtinumo funkcijos maksimumo taške šis kriterijus yra pasiskirstęs pagal Fišerio dėsnį su $\nu = \frac{M(M+3)}{2}$ laisvės laipsnių.

Algoritmas stabdomas, jei stabdymo kriterijus neprieštaruja hipotezei apie vidurkių ir kovariacijų matricų sutapimą dviejose gretimose iteracijose:

$$H^t \leq F_{\alpha,\nu} \quad (27)$$

tikėtinumo funkcijos pasiklivimo intervalas yra mažesnis už pasirinktą mažą reikšmę ε :

$$2 \cdot \eta_\gamma \cdot \sqrt{\frac{d^t}{N^t}} \leq \varepsilon, \quad (28)$$

α, γ yra atitinkami Fišerio ir normaliojo skirstinio kvantiliai, ir kovariacijų matricų įverčių didžiausios ir mažiausios tikrinių reikšmių santykis ne daugiau 10.

Jei nors vienas stabdymo sąlyga nėra tenkinama, generuojama nauja Monte–Karlo imtis. Imties tūrį galima reguliuoti pagal taisyklę (Sakalauskas, 2000):

$$N^{t+1} \geq \frac{N^t \cdot \nu}{H^t} \cdot F_{\beta,\nu}, \quad (29)$$

čia $F_{\beta,\nu}$ – Fišerio skirstinio kvantilis, β – pasiklivimo lygmuo. Šitokios taisyklės taikymas leidžia racionaliai parinkti imčių tūrį Monte–Karlo Markovo grandinėje ir kartu užtikrina konvergavimą į tikėtinumo funkcijos maksimumą.

5. Kompiuterinis modeliavimas

Sudarytam algoritmui patikrinti buvo atliktas kompiuterinis eksperimentas. Sugeneruota imtis $\Theta = (\theta_1, \theta_2, \dots, \theta_K)$, sudaryta iš $K = 10$ populiacijų, kuriose galėjo įvykti $M = 3$ įvykių, kai tikimybių logitai pasiskirstę pagal daugiamatį Gauso dėsnį su parametrais:

$$\mu = \begin{pmatrix} -3 \\ -4 \\ -5 \end{pmatrix}; \quad \Sigma = \begin{pmatrix} 0,25 & 0 & 0 \\ 0 & 0,25 & 0 \\ 0 & 0 & 0,25 \end{pmatrix}. \quad (30)$$

Toliau buvo sugeneruota $t = 100$ Monte–Karlo Markovo grandžių pagal (15)–(20), (29) išraiškas. Siekiant išvengti labai mažų arba labai didelių imties tūrio reikšmių, buvo taikomos ribos: $500 \leq N^k \leq 17000$.

Sprendinys pradėjo konverguoti jau po $t = 6$ iteracijų. Gautos parametų reikšmės pateiktos 1 lentelėje:

1 lentelė. Parametų įverčiai.

Iteracija	μ_1	μ_2	μ_3	Tikslo funkcija	Pasiklovimo intervalas	Imties tūris	Statistika
1	-2,96	-4,29	-5,52	-62,90	5,57	500	9,55
2	-2,89	-4,04	-5,27	-396,58	4,81	500	6,18
3	-2,91	-4,03	-5,19	-420,42	2,97	500	3,86
4	-2,90	-4,04	-5,16	-424,87	3,2	500	0,35
5	-2,91	-4,04	-5,13	-428,05	1,57	2 963	1,41
6	-2,90	-4,04	-5,14	-427,57	1,32	4 383	0,32
7	-2,91	-4,04	-5,13	-425,54	0,75	13 986	0,40
8	-2,91	-4,04	-5,14	-425,33	0,75	14 345	0,40
9	-2,91	-4,04	-5,13	-425,71	0,75	13 525	0,84
10	-2,91	-4,04	-5,13	-426,47	0,75	15 135	0,22

Išvados

Darbe sudarytas MCMC algoritmas, pasinaudojus statistiniu stabdymo kriterijumi ir imties tūrio reguliavimu. Atlikus skaičiavimus, gaunami nežinomų parametų įverčiai ir tikimybės. Sudarytas algoritmas gali būti taikomas socialinių ir medicininių duomenų analizei.

Literatūra

- Abramowitz, M.; Stegun, I. A. (1964). Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. New York: Dover Publications.
- Altaleb, Anas; Chauveau, Didier (2002). Bayesian analysis of the Logit model and comparison of two. Metropolis–Hastings strategies. Computational Statistics and Data Analysis, 39: 137–152.
- Bradley, P. C.; Thomas, A. L. (2000). Bayes and Empirical Bayes Methods for Data Analysis. New York: Chapman and Hall.
- Chen, Fang (2009). Bayesian modeling using the MCMC procedure [interaktyvus]. [žiūrėta 2013 m. gegužės 4 d.]. Prieiga per internetą: <<http://support.sas.com/resources/papers/proceedings09/257-2009.pdf>>.
- Clayton, David; Kaldor, John (1987). Empirical Bayes estimates of age-standardized relative risks for use in disease mapping, Biometrics, 43(3): 671–681.
- Dennis J. E.; Schnabel, R. B. (1996). Numerical methods for unconstrained optimization

- and nonlinear equations. Philadelphia: Classics in Applied Mathematics.
- Kantorovich, L. V.; Akilov, G. P. (1982). *Functional Analysis*. Oxford: Pergamon Press.
- Liseo, Brunero; Loperfido, Nicola (2003). A Bayesian interpretation of the multivariate skew-normal distribution. *Statistics & Probability Letters*, 61(4): 395–401.
- Pearce, D. W. (2006). *Aiškinamasis ekonomikos anglų–lietuvių kalbų žodynas*. Lithuania: Vilnius, TEV.
- Sakalauskas, Leonidas (2000). Nonlinear stochastic optimization by Monte–Carlo estimators. *Informatica*, 11(4): 455–468.
- Tsutakawa, Robert K.; Shoop, Gary L., Marienfeld, Carl J. (1985). Empirical Bayes estimation of cancer mortality rates. *Statistics in Medicine*, 4(2): 201–212.

MULTIDIMENSIONAL RARE EVENT PROBABILITY ESTIMATION ALGORITHM

Leonidas Sakalauskas, Ingrida Vaiciulyte

Summary

This work contains Monte–Carlo Markov Chain algorithm for estimation of multi-dimensional rare events frequencies. Logits of rare event likelihood we are modeling with Poisson distribution, which parameters are distributed by multivariate normal law with unknown parameters – mean vector and covariance matrix. The estimations of unknown parameters are calculated by the maximum likelihood method. There are equations derived, those must be satisfied with model’s maximum likelihood parameters estimations. Positive definition of evaluated covariance matrixes are controlled by calculating ratio between matrix maximum and minimum eigenvalues.

Key words: Monte–Carlo Markov Chain, Bayesian method, likelihood method, Poisson–Gaussian model.