

برآورد بیزی تابع تاوان در آزمون همگنی مدل‌های آمیخته

رحمان فرنوش^۱، افشین فلاح^۲، آرزو حاج رجبی^۱

^۱گروه آمار، دانشگاه علم و صنعت

^۲گروه آمار، دانشگاه بین‌المللی امام خمینی

تاریخ دریافت: ۱۳۸۷/۷/۱۷ تاریخ آخرین بازنگری: ۱۳۸۷/۱۲/۲۰

چکیده: برای آزمون فرضیه همگنی مدل‌های آمیخته، معمولاً از آزمون نسبت درست‌نمایی اصلاح شده که مبتنی بر افزودن یک تابع تاوان مناسب به تابع لگ درست‌نمایی می‌باشد، استفاده می‌شود. کارایی این آزمون به شدت تحت تأثیر شکل تابع تاوان انتخابی است. انتخاب تابع تاوان در این نوع آزمون معمولاً بر اساس پرهیز از پیچیدگی و میسر بودن برآورد پارامترها صورت می‌پذیرد، که لزوماً نتایج مطلوبی بدنبال ندارد. در این مقاله یک تابع تاوان جامع در نظر گرفته شده است، که دارای یک پارامتر تعیین کننده شکل است. سپس پارامتر تعیین کننده شکل این تابع تاوان و پارامترهای مدل آمیخته با در نظر گرفتن توزیع‌های پیشین مناسب برای آنها با استفاده از رهیافت بیزی، بصورت پسینی برآورد شده‌اند. نشان داده شده است که رهیافت بیزی پیشنهادی در برآورد پارامترهای مدل، در مقایسه با رهیافت بسامدی، به مراتب کارایی مطلوب‌تری دارد. این کارایی خصوصاً در شرایط شناخت ناپذیری توزیع آمیخته که روش‌های بسامدی کارایی اندکی دارند، بیشتر است.

واژه‌های کلیدی: آزمون نسبت درست‌نمایی، تابع تاوان، الگوریتم EM ، زنجیره‌های مارکف مونت کارلوئی.

آدرس الکترونیک مسئول مقاله: رحمان فرنوش، rfaranoosh@iust.ac.ir
کد موضوع‌بندی ریاضی (۲۰۰۰): ۶۲۴۱۵

توزیع‌های آمیخته به دلیل انعطاف پذیر بودن در توصیف بسیاری از پدیده‌های تصادفی مفید هستند. پس از مقاله کارل پیرسون در دهه هشتاد این توزیع‌ها بطور گسترده‌ای در تحقیقات کاربردی و در زمینه‌های متفاوتی مانند ژنتیک، زیست‌شناسی، اقتصاد و غیره بکار گرفته می‌شوند. کاربرد توزیع‌های آمیخته معمولاً در مواردی است که جامعه‌ی آماری ناهمگن و ترکیبی از چند زیر جامعه است. در چنین حالتی مشاهداتی که از این جامعه بدست می‌آیند، می‌توانند با احتمال معینی به هر یک از این زیر جامعه‌ها تعلق داشته باشند. روش‌های مختلفی برای آزمون تعداد مؤلفه‌های یک توزیع آمیخته وجود دارد، که از آن جمله می‌توان به روش‌هایی بر پایه نظریه اطلاع، فاصله تاوانیده، روش لگ درست‌نمایی تاوانیده و روش بیزی اشاره نمود. آکائیک (۱۹۷۳) به کمک نظریه اطلاع، معیار اطلاع آکائیک را برای انتخاب تعداد مؤلفه‌های مدل آمیخته، از طریق مینیمم ساختن فاصله کولبک-لایبلر بین توزیع جامعه و توزیع مدل‌های پیشنهادی مورد مطالعه قرار داد. چن و کال فلیسک (۱۹۹۶) روشی را بر پایه مینیمم ساختن فاصله تاوانیده بین تابع توزیع تجربی و تابع توزیع تجمعی برازش داده شده، ارائه دادند. چن و خلیلی (۲۰۰۶) روشی را بر پایه لگ درست‌نمایی تاوانیده، برای انتخاب تعداد مؤلفه‌های مدل آمیخته معرفی کردند. ریچاردسون و گرین (۱۹۹۷) انتخاب تعداد مؤلفه‌های مدل آمیخته را از دیدگاه بیزی مورد مطالعه قرار دادند. مارین و همکاران (۲۰۰۵) دشواری‌های تحلیل بیزی مدل‌های آمیخته را بررسی نموده و نشان دادند که وقتی توزیع جامعه آمیخته است، کاربست روشهای معمول زنجیر مارکف مونت کارلویی برای نمونه‌گیری از توزیع پسین پارامترها با دشواری‌های زیادی روبرو است. آزمون فرضیه همگن بودن جامعه توسط محققان بسیاری مورد توجه قرار گرفته است. در حالت خاص، برای انتخاب یک مدل همگن در مقابل یک مدل آمیخته با دو مؤلفه، هارتیگان (۱۹۸۵) استفاده از آزمون نسبت درست‌نمایی را پیشنهاد نمود. ولی استفاده از این روش به دلیل پیچیدگی توزیع حدی آماره آزمون نسبت درست‌نمایی، عملاً امکان پذیر نیست. چن (۱۹۹۸) منابع بی‌نظمی مؤثر بر توزیع حدی آماره نسبت درست‌نمایی را مورد مطالعه قرار داد و بر این اساس آزمون نسبت درست‌نمایی اصلاح شده را معرفی نمود، که در آن با افزودن یک تابع تاوان به تابع لگ درست‌نمایی، توزیع حدی ساده‌ای برای آماره‌ی آزمون نسبت درست‌نمایی بدست می‌آید. یکی از محدودیت‌های این روش آن است که استفاده از آزمون نسبت درست‌نمایی اصلاح شده، منوط به برقراری برخی شرایط نظم است. بعلاوه توان این آزمون به شدت تحت تاثیر تابع تاوان انتخابی است. از اینرو لی و همکاران (۲۰۰۸) آزمون EM را بر اساس شکل دیگری از

تابع توان پیشنهادی توسط چن و همکاران (۲۰۰۱)، که مستقل از فرضیات لازم برای آزمون نسبت درستنمایی اصلاح شده است، مطرح نمودند.

در این مقاله آزمون فرض همگنی جامعه با در نظر گرفتن فرم جامعی برای تابع توان که به یک پارامتر شکل وابسته است، از دیدگاه بیزی مورد مطالعه قرار گرفته است. برای این منظور با در نظر گرفتن پیشین‌های مناسب، پارامتر تعیین کننده شکل تابع توان و پارامترهای مدل آمیخته بصورت پسینی برآورد شده‌اند. در بخش ۲ آزمون نسبت درستنمایی و شکل اصلاح شده آن معرفی و مشکلات این نوع آزمونها مورد بررسی قرار گرفته است. در بخش ۳ چگونگی برآورد پارامترهای مدل با استفاده از الگوریتم EM شرح داده شده است. مسئله انتخاب تابع توان در بخش ۴ مورد بحث قرار گرفته و روشی برای برآورد پارامتر تعیین کننده شکل تابع توان به کمک رهیافت بیزی، پیشنهاد شده است. در بخش ۵، کارایی آزمون نسبت درستنمایی و روش بیزی پیشنهادی با استفاده از یک مطالعه شبیه‌سازی مورد ارزیابی قرار گرفته است.

۲ آزمون نسبت درستنمایی

متغیر تصادفی Y دارای توزیع آمیخته است، هرگاه تابع چگالی یا جرم احتمال آن بصورت

$$p_Y(y; \theta, p) = \sum_{j=1}^k p_j f_j(y|\theta_j), \quad y \in \mathcal{Y}, \quad 0 \leq p_j \leq 1, \quad j = 1, \dots, k, \quad \sum_{j=1}^k p_j = 1,$$

باشد، که در آن \mathcal{Y} تکیه‌گاه متغیر تصادفی Y را نشان می‌دهد و $\theta = (\theta_1, \dots, \theta_k)$ که θ_j پارامتر مؤلفه‌ی j ام و $\varphi = (p_1, \dots, p_k)$ که p_j نسبت آمیخته مؤلفه‌ی j ام نامیده می‌شود. هنگامی که از یک توزیع آمیخته نمونه‌ای مشاهده می‌شود، از آنجا که مشخص نیست هر مشاهده مربوط به کدام زیر جامعه می‌باشد، مشاهدات بدست آمده را داده‌های ناقص گویند. در اینصورت تابع لگاریتم درستنمایی داده‌های ناقص بصورت

$$\ell_n(\theta, p) = \sum_{i=1}^n \log\left(\sum_{j=1}^k p_j f_j(y_i|\theta_j)\right),$$

است. داده‌های ناقص را می‌توان با در نظر گرفتن مجموعه‌ای از متغیرهای نشانگر که تعلق مشاهدات به زیر مجموعه‌ها را مشخص می‌سازند، به داده‌های کامل تبدیل نمود. تابع لگاریتم

درست‌نمایی داده‌های کامل که معمولاً برای انجام استنباط مناسبتر می‌باشد، بصورت

$$\ell_n^c(z, \theta, p) = \sum_{i=1}^n \sum_{j=1}^k z_{ij} \{\log p_j + \log f_j(y_i | \theta_j)\},$$

است. درست‌نمایی داده‌های کامل، تابعی از متغیرهای پنهان $Z = (z_1, \dots, z_n)$ و مشاهدات می‌باشد، که در آن $z_i = (z_{i1}, \dots, z_{ik})$ و $z_{ij} = 1$ نشان دهنده تعلق مشاهده i ام به مؤلفه‌ی j ام است. فرض کنید Y_1, \dots, Y_n نمونه‌ای تصادفی از توزیع آمیخته

$$pf_Y(y | \theta_1) + (1-p)f_Y(y | \theta_2), \quad (1)$$

باشد، که در آن $\theta_1, \theta_2 \in \Theta$ و Θ زیر مجموعه‌ای فشرده از خط حقیقی است. هدف آزمون

$$\begin{cases} \mathcal{H}_0 : p(1-p)(\theta_1 - \theta_2) = 0 \\ \mathcal{H}_1 : p(1-p)(\theta_1 - \theta_2) \neq 0, \end{cases} \quad (2)$$

است، که در آن فرض صفر به معنی همگنی جامعه و فرض مقابل به معنی تشکیل جامعه از دو زیر جامعه ناهمگن مطابق رابطه (۱) است. تابع لگاریتم درست‌نمایی را می‌توان بصورت

$$\ell_n(p, \theta_1, \theta_2) = \sum_{i=1}^n \log \{pf(y_i | \theta_1) + (1-p)f(y_i | \theta_2)\},$$

نوشت. اگر $\hat{\theta}_0$ و $(\hat{p}, \hat{\theta}_1, \hat{\theta}_2)$ به ترتیب ماکسیمم‌کننده تابع درست‌نمایی تحت فرض صفر و مقابل باشند، در اینصورت آماره آزمون نسبت درست‌نمایی بصورت

$$R_n = 2 \{ \ell_n(\hat{p}, \hat{\theta}_1, \hat{\theta}_2) - \ell_n(0/5, \hat{\theta}_0, \hat{\theta}_0) \},$$

می‌باشد و مقادیر بزرگ این آماره منجر به رد فرض صفر می‌شوند. چن (۱۹۹۸) نشان داد دو منبع بی‌نظمی که توزیع حدی آماره آزمون نسبت درست‌نمایی را به طور قابل ملاحظه‌ای پیچیده می‌سازند، یکی قرار داشتن نقاط $p = 0$ یا $p = 1$ در محدوده فرض صفر و دیگری شناخت ناپذیر بودن توزیع آمیخته (۱) تحت فرض صفر است. وی پیشنهاد نمود برای برطرف کردن این مشکل، تابع تاوانی برحسب p به صورت $T(p)$ که در شرایط

$$\lim_{p \rightarrow 0 \text{ or } 1} T(p) = -\infty, \quad \arg \max_{p \in [0, 1]} T(p) = 0/5, \quad (3)$$

صدق کند، به آماره آزمون نسبت درست‌نمایی افزوده شود. در اینصورت تابع لگاریتم درست‌نمایی تاوانیده بصورت $T\ell_n(p, \theta_1, \theta_2) = \ell_n(p, \theta_1, \theta_2) + T(p)$ است و آماره آزمون نسبت درست‌نمایی بصورت $M_n = 2 \{ T\ell_n(\hat{p}, \hat{\theta}_1, \hat{\theta}_2) - T\ell_n(0/5, \hat{\theta}_0, \hat{\theta}_0) \}$ اصلاح می‌شود.

استفاده از این تابع تاوان موجب می شود مقادیر برازش داده شده p تحت درستنمایی اصلاح شده دور از صفر و یک واقع شوند. چن (۱۹۹۸) نشان داد که اگر برخی شرایط نظم روی هسته چگالی برقرار باشند، آنگاه توزیع مجانبی آماره آزمون نسبت درستنمایی اصلاح شده، آمیخته‌ای از دو توزیع با وزن‌های یکسان بصورت $\frac{1}{2}\chi^2 + \frac{1}{2}\chi^2_0$ می باشد، که در آن χ^2_0 توزیع تباهیده در نقطه صفر و χ^2 توزیع کای دو با درجه آزادی ۱ را نشان می دهند. نکته‌ی قابل توجه آن است که این توزیع حدی به فرم تابع تاوان انتخابی وابسته نیست.

۳ برآورد پارامترها با الگوریتم EM

برای یافتن برآورد ماکسیمم درستنمایی پارامترها در توزیع‌های آمیخته، معمولاً از الگوریتم تکراری EM استفاده می شود. که شامل دو مرحله امیدگیری E و ماکسیمم سازی M است. فرض کنید Y_1, \dots, Y_n نمونه‌ای تصادفی از توزیع آمیخته (۱) باشد، در این صورت تابع لگاریتم درستنمایی داده‌های کامل بصورت

$$\ell_n((\theta_1, \theta_2), p) = \sum_{i=1}^n [z_{i1} \{\log(p) + \log(f(y_i|\theta_1))\} + (\lambda - z_{i1}) \{\log(1-p) + \log(f(y_i|\theta_2))\}], \quad (4)$$

است. در مرحله E ، امید ریاضی

$$z_{i1}^{(t)} = \frac{E(z_{i1}|y_i, \theta^{(t-1)}, p^{(t-1)})}{p^{(t-1)}f(y_i|\theta_1^{(t-1)}) + (1-p^{(t-1)})f(y_i|\theta_2^{(t-1)})}, \quad i = 1, \dots, n, \quad (5)$$

محاسبه و جایگزین z_{ij} , $i = 1, \dots, n$ می شود. در مرحله M با قرار دادن مقادیر (۵) در تابع لگ درستنمایی داده‌های کامل و ماکسیمم کردن این تابع نسبت به پارامترهای مدل، داریم

$$p^{(t)} = \frac{1}{n} \sum_{i=1}^n z_{i1}^{(t)}$$

$$\theta_1^{(t)} = \arg \max_{\theta_1 \in \Theta} \sum_{i=1}^n \{z_{i1}^{(t)} \log(f(y_i|\theta_1))\},$$

$$\theta_2^{(t)} = \arg \max_{\theta_2 \in \Theta} \sum_{i=1}^n \{(1 - z_{i1}^{(t)}) \log(f(y_i|\theta_2))\}.$$

در این صورت، تکرار مراحل E و M تا حصول همگرایی، برآوردهای ماکسیمم درستنمایی پارامترها را بدست می دهد. برآورد p به انتخاب تابع تاوان بستگی دارد، چن و همکاران

(۲۰۰۱) تابع تاوانی به صورت

$$T(p) = C \log(4p(1-p)), \quad (6)$$

را پیشنهاد نمودند، که در آن C ثابتی مثبت و تأثیر گذار بر میزان اصلاح تابع درست‌نمایی است. لی و همکاران (۲۰۰۸) نیز تابع تاوانی به فرم

$$T(p) = C^* \log(1 - |1 - 2p|), \quad (7)$$

را مورد استفاده قرار دادند. توابع تاوان (۶) و (۷) در شرایط (۳) صدق می‌کنند و به سادگی می‌توان نشان داد که مقدار p در مرحله M از تکرار t الگوریتم EM ، برای این توابع به

$$p^{(t+1)} = \frac{\sum_{i=1}^n z_i^{(t)} + C}{2C+n}$$

ترتیب بصورت

$$p^{(t)} = \begin{cases} \min\left\{\frac{\sum_{i=1}^n z_i^{(t)} + C^*}{n+C^*}, 0.5\right\} & \frac{\sum_{i=1}^n z_i^{(t)}}{n} < 0.5, \\ 0.5 & \frac{\sum_{i=1}^n z_i^{(t)}}{n} = 0.5, \\ \max\left\{\frac{\sum_{i=1}^n z_i^{(t)}}{n}, 0.5\right\} & \frac{\sum_{i=1}^n z_i^{(t)}}{n} > 0.5, \end{cases}$$

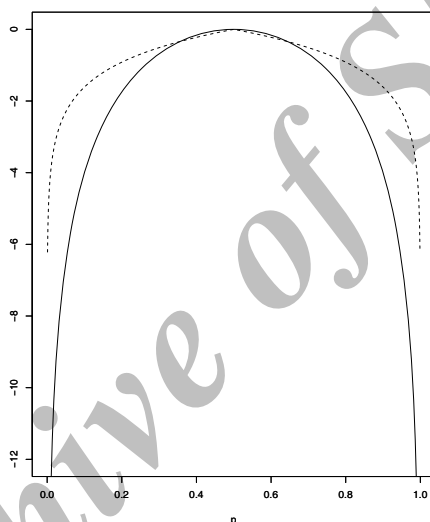
بدست می‌آیند.

۴ انتخاب تابع تاوان

گرچه توزیع حدی آماره آزمون نسبت درست‌نمایی به فرم تابع تاوان انتخابی وابسته نیست، اما فرم تابع تاوان انتخاب شده بر توان آزمون تأثیر گزار است. آزمون اصلاح شده‌ای که از افزودن تابع تاوان (۶) به تابع لگاریتم درست‌نمایی حاصل می‌شود، مشکلات آزمون نسبت درست‌نمایی معمول را تا حدودی مرتفع می‌سازد. با این وجود گاهی حتی با وجود مشاهدات مشهود از توزیع آمیخته، آزمون نسبت درست‌نمایی اصلاح شده قادر به رد فرض H_0 نیست. دلیل این امر آن است که تابع تاوان (۶) به ازای نسبت‌های آمیخته نزدیک صفر یا یک، تاوان زیادی را روی تابع لگاریتم درست‌نمایی اعمال می‌کند. به همین دلیل به تابع تاوان معقول‌تری نیاز است، به نحوی که برای آن توان آزمون نسبت درست‌نمایی اصلاح شده، حتی در حالتی که نسبت‌های آمیخته به صفر و یک نزدیک هستند، افزایش یابد. تابع تاوان (۷) به این منظور پیشنهاد شده است. به سادگی ملاحظه می‌شود که بین دو تابع تاوان (۶) و (۷) نابرابری

$$\log(1 - |1 - 2p|) \leq \log(1 - |1 - 2p|^2) = \log(4p(1-p)),$$

برقرار است، که به ازای مقادیر نزدیک به $0/5$ نسبت آمیخته، این نابرابری به برابری تقریبی $\log(1 - |1 - 2p|) \approx -|1 - 2p|$ تبدیل می‌شود. بنابراین توابع تاوان (۶) و (۷) به ازای مقادیر p نزدیک به $0/5$ تقریباً معادل‌اند، اما به ازای مقادیر p نزدیک به صفر یا یک، تفاوت قابل ملاحظه‌ای بین آنها وجود دارد. شکل ۱ نمودار این دو تابع تاوان را به ازای $C^* = 1$ نشان می‌دهد. ملاحظه می‌شود که تابع تاوان (۶) در نقطه $p = 0/5$ تقریباً مسطح می‌باشد و توانی را اعمال نمی‌کند، در حالی که تابع تاوان (۷) علاوه بر مرتفع‌تر بودن، به ازای مقادیر نزدیک به صفر و یک نیز تاوان بسیار کمتری را لحاظ می‌نماید.



شکل ۱: نمودار توابع تاوان (۶) (خط) و (۷) (نقطه چین).

با وجود اینکه تابع تاوان پیشنهادی لی و همکاران (۲۰۰۸)، دارای مزیت‌هایی می‌باشد و کارایی آزمون نسبت درستی را افزایش می‌دهد، اما دلیلی مبنی بر بهینه بودن هیچ یک از این توابع تاوان وجود ندارد. بسادگی ملاحظه می‌شود که توابع تاوان (۶) و (۷) حالات خاصی از تابع کلی

$$g(p, h) = C \log(1 - |1 - 2p|^h), \quad 0 < h \leq 2 \quad (8)$$

هستند، که به ازای مقادیر $h = 1$ و $h = 2$ حاصل می‌شوند. انتخاب مقدار h به شکل تابع

تاوان و در نتیجه به استنباط‌های حاصل از آن به شدت تأثیر گذار است. در رهیافت بسامدی انتخاب مقادیر دیگری برای h به پیچیدگی تابع تاوان (۸) نتیجه می‌شود، که آن نیز به نوبه خود منجر به دشواری ماکسیمم سازی در مرحله M از الگوریتم EM می‌شود. با توجه به مشکلات رهیافت بسامدی در تعیین پارامترهای مدل و تابع تاوان، در این بخش این پارامترها در چارچوب رهیافت بیزی برآورد می‌شوند. برای این منظور ابتدا برای پارامترهای مورد نظر، توزیع‌های پیشین مناسب در نظر گرفته و سپس برآورد این پارامترها بصورت پسینی محاسبه می‌شوند. در حالت کلی، دامنه تغییرات پارامتر h در تابع تاوان (۸) مجموعه اعداد حقیقی مثبت است. با این وجود میزان تاوانی که این تابع اعمال می‌کند با بزرگ شدن h به صفر میل می‌نماید، به نحوی که برای مقادیر $h > 1$ میزان تاوان اعمال شده توسط این تابع اندک بوده و برای مقادیر $h > 2$ عملاً هیچ تاوانی اعمال نمی‌شود. از این رو برخی محققان مانند لی و همکاران (۲۰۰۸) مقداری در بازه $(0, 1]$ را برای این پارامتر پیشنهاد کرده‌اند. در این مقاله از توزیع $U(0, 1)$ به عنوان توزیع پیشین برای پارامتر h استفاده شده است. به منظور حفظ بی طرفی، توزیع پیشین برای نسبت آمیخته p ، $U(0, 1)$ در نظر گرفته شده است. چون توابع هسته مختلف تشکیل دهنده یک توزیع آمیخته، دارای پارامترهای متفاوتی هستند، برای انجام تحلیل بیزی لازم است برای هر یک از این پارامترها نیز توزیع پیشین مناسبی لحاظ شوند. اگر Y_1, \dots, Y_n نمونه‌ای تصادفی از توزیع آمیخته متناهی دو مؤلفه‌ای نرمال باشد، در این صورت تابع درست‌نمایی تاوانیده که از ضرب تابع درست‌نمایی در تابع تاوان حاصل می‌شود، بصورت

$$L(y|\mu_1, \mu_2, p, h) \propto p^{n_1} (1-p)^{n_2} e^{-\frac{n_1}{2}(\bar{y}_1 - \mu_1)^2} e^{-\frac{n_2}{2}(\bar{y}_2 - \mu_2)^2} (1 - |1 - 2p|^h)$$

است، که در آن تعداد مشاهدات متناسب به مؤلفه اول، n_2 تعداد مشاهدات متناسب به مؤلفه دوم، \bar{y}_1 میانگین مشاهدات متناسب به مؤلفه اول و \bar{y}_2 میانگین مشاهدات متناسب به مؤلفه دوم است. با فرض اینکه h و p بصورت پیشینی مستقل هستند و با در نظر گرفتن توزیع‌های پیشین نرمال $\mu_j \sim N(\delta, \frac{1}{\lambda}), j = 1, 2$ با $\delta \in R, \lambda > 0$ برای پارامترهای مکانی دو مؤلفه توزیع آمیخته نرمال، می‌توان نوشت

$$\begin{aligned} \pi(\mu_1, \mu_2) &\propto e^{-\frac{\lambda}{2}\{(\mu_1 - \delta)^2 + (\mu_2 - \delta)^2\}}, \\ \pi(h, p) &= \pi(h)\pi(p) = 1, \\ \pi(\mu_1, \mu_2, h, p) &\propto e^{-\frac{\lambda}{2}\{(\mu_1 - \delta)^2 + (\mu_2 - \delta)^2\}}, \\ \pi(\mu_1, \mu_2, h, p|y) &= f(y|\mu_1, \mu_2, p, h)\pi(\mu_1, \mu_2, h, p). \end{aligned}$$

از این رو توزیع پسین را می توان بصورت

$$\pi(\mu_1, \mu_2, h, p|y) \propto p^{n_1}(1-p)^{n_2} e^{-\frac{n_1}{p}(\bar{y}_1 - \mu_1)^2} e^{-\frac{n_2}{1-p}(\bar{y}_2 - \mu_2)^2} (1 - |1 - 2p|^h) e^{-\frac{\delta}{p}\{(\mu_1 - \delta)^2 + (\mu_2 - \delta)^2\}}$$

نوشت. به همین ترتیب، اگر Y_1, \dots, Y_n یک نمونه تصادفی از توزیع آمیخته متناهی دو مؤلفه‌ای پواسن باشد، در این صورت تابع درستنمایی توانیده از ضرب تابع درستنمایی در تابع تاوان بصورت

$$L(y|\lambda_1, \lambda_2, p, h) \propto p^{n_1}(1-p)^{n_2} \exp\{-n_1\lambda_1 - n_2\lambda_2 + n_1\bar{y}_1 \ln \lambda_1 + n_2\bar{y}_2 \ln \lambda_2\} (1 - |1 - 2p|^h),$$

حاصل می شود، که در آن n_1 تعداد مشاهدات متناسب به مؤلفه اول، n_2 تعداد مشاهدات متناسب به مؤلفه دوم، \bar{y}_1 میانگین مشاهدات متناسب به مؤلفه اول و \bar{y}_2 میانگین مشاهدات متناسب به مؤلفه دوم است. با فرض اینکه h و p بصورت پیشینی مستقل هستند و با در نظر گرفتن توزیع‌های پیشین گاما $\lambda_j \sim \Gamma(\alpha_j, \beta_j)$ ، $j = 1, 2$ با $\alpha_j > 0, \beta_j > 0$ برای پارامترهای λ در دو مؤلفه‌ای توزیع آمیخته‌ی پواسن، می توان نوشت

$$\begin{aligned} \pi(\lambda_1, \lambda_2) &\propto \lambda_1^{(\alpha_1-1)} e^{-\frac{1}{\beta_1}\lambda_1} \lambda_2^{(\alpha_2-1)} e^{-\frac{1}{\beta_2}\lambda_2}, \\ \pi(h, p) &= \pi(h)\pi(p) = 1, \\ \pi(\lambda_1, \lambda_2, h, p) &\propto \lambda_1^{(\alpha_1-1)} e^{-\frac{1}{\beta_1}\lambda_1} \lambda_2^{(\alpha_2-1)} e^{-\frac{1}{\beta_2}\lambda_2}, \\ \pi(\lambda_1, \lambda_2, h, p|y) &= f(y|\lambda_1, \lambda_2, p, h)\pi(\lambda_1, \lambda_2, h, p), \end{aligned}$$

از این رو توزیع پسین را می توان بصورت

$$\pi(\lambda_1, \lambda_2, h, p|y) \propto p^{n_1}(1-p)^{n_2} e^{\{-n_1\lambda_1 - n_2\lambda_2 + n_1\bar{y}_1 \ln \lambda_1 + n_2\bar{y}_2 \ln \lambda_2\}} (1 - |1 - 2p|^h) \lambda_1^{(\alpha_1-1)} e^{-\frac{1}{\beta_1}\lambda_1} \lambda_2^{(\alpha_2-1)} e^{-\frac{1}{\beta_2}\lambda_2},$$

نوشت. تحلیل بیزی مدل‌های آمیخته به دلیل پیچیدگی ذاتی این مدل‌ها، با دشواری‌های زیادی همراه است (دایبولت و رابرت، ۱۹۹۰). معمولاً در این حالت توزیع پسین فاقد یک فرم بسته می باشد. از این رو برای انجام استنباط، به کمک الگوریتم‌های زنجیر مارکوف مونت کارلویی (MCMC)، از توزیع پسین نمونه‌گیری می شود. الگوریتم‌های گیبز و مترو پولیس-هاستینگ، از جمله مهمترین الگوریتم‌ها در زمینه نمونه‌گیری از توزیع پسین هستند. اما، برای استفاده از الگوریتم گیبز لازم است توزیع‌های شرطی کامل پارامترها در دسترس باشند. بعلاوه نشان داده شده است که وقتی این الگوریتم برای برآورد پارامترهای توزیع‌های آمیخته به کار گرفته می شود، همگرایی آن به مقادیر اولیه وابسته است. از این رو اگر مقادیر اولیه به ماکسیمم موضعی نزدیک باشند، ممکن است حتی برای تکرارهای بسیار زیاد

الگوریتم توانایی‌رهایی از جذب ماکسیمم موضعی را نداشته باشد. به همین دلیل از الگوریتم مترو پولیس-هاستینگ که در هنگام کار با توزیع‌های آمیخته فاقد معایب الگوریتم گیبز می‌باشد، برای برآورد پارامترهای مدل استفاده شده است (مارین و همکاران، ۲۰۰۵). در این الگوریتم با ساختن زنجیر مارکوف $\{\theta^{(t)}\}_{t=1}^N$ که دارای توزیع مانای $\pi(\theta|y)$ است، مقدار $E(g(\theta)|y)$ توسط برآوردگر سازگار $\frac{1}{N} \sum_{t=1}^N g(\theta^{(t)})$ تخمین زده می‌شود. فرض کنید Y_1, \dots, Y_n یک نمونه تصادفی از توزیع آمیخته‌ی دو مؤلفه‌ای باشد، شکل کلی الگوریتم متروپولیس - هاستینگ برای برآورد پارامترهای مدل آمیخته و نیز پارامتر h بصورت زیر است:

- آ. مقادیر اولیه‌ی $p^{(0)}$ ، $\theta^{(0)}$ و $h^{(0)}$ را در نظر بگیرید.
- ب. مراحل زیر را تا رسیدن به توزیع مانای زنجیر مارکوف $\{\theta^{(t)}\}_{t=1}^N$ تکرار کنید.
- پ. $(\tilde{\theta}, \tilde{p}, \tilde{h})$ را از توزیع پیشنهادی $(\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)})$ تولید کنید.
- ت. مقدار r را محاسبه کنید،

$$r = \frac{f(y|\tilde{\theta}, \tilde{p}, \tilde{h})\Pi(\tilde{\theta}, \tilde{p}, \tilde{h})q(\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)}|\tilde{\theta}, \tilde{p}, \tilde{h})}{f(y|\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)})\Pi(\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)})q(\tilde{\theta}, \tilde{p}, \tilde{h}|\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)})}$$

ث. مشاهده‌ی u را از توزیع $U(0, 1)$ در نظر بگیرید. اگر $r > u$ قرار دهید $(\theta^{(t)}, p^{(t)}, h^{(t)}) = (\tilde{\theta}, \tilde{p}, \tilde{h})$ و در غیر اینصورت قرار دهید $(\theta^{(t)}, p^{(t)}, h^{(t)}) = (\theta^{(t-1)}, p^{(t-1)}, h^{(t-1)})$.

۵ شبیه‌سازی

به منظور ارزیابی و مقایسه‌ی تأثیر توابع تاوان (۶) و (۷) بر آزمون نسبت درستمایی اصلاح شده، یک مطالعه‌ی شبیه‌سازی اجرا شده است. برای این منظور آمیخته‌ای از دو توزیع پواسن با پارامترهای مختلف به نحوی در نظر گرفته شده است، که میانگین توزیع تحت فرض‌های \mathcal{H}_0 و \mathcal{H}_1 با یکدیگر مساوی و برابر با ۵ باشد. برای این منظور از روابط

$$\begin{aligned} E_{\mathcal{H}_0}(Y) &= \text{Var}_{\mathcal{H}_0}(Y), \\ E_{\mathcal{H}_1}(Y) &= p\lambda_1 + (1-p)\lambda_2, \\ \text{Var}_{\mathcal{H}_1}(Y) &= E_{\mathcal{H}_1}(Y^2) - E_{\mathcal{H}_1}^2(Y) \\ &= p(1-p)(\lambda_1 - \lambda_2)^2 + 5, \end{aligned}$$

استفاده شده است. بصورت مشابه آمیخته‌ای از توزیع نرمال با پارامترهای متفاوت نیز به نحوی در نظر گرفته شده است، که میانگین توزیع تحت فرض‌های \mathcal{H}_0 و \mathcal{H}_1 با یکدیگر

مساوی و برابر با ۰ باشد. برای این منظور روابط

$$\begin{aligned} E_{\mathcal{H}_0}(Y) &= 0, \quad \text{Var}_{\mathcal{H}_0}(Y) = 1, \\ E_{\mathcal{H}_1}(Y) &= p\mu_1 + (1-p)\mu_2, \\ \text{Var}_{\mathcal{H}_1}(Y) &= E_{\mathcal{H}_1}(Y^2) - E_{\mathcal{H}_1}^2(Y) \\ &= p(1-p)(\mu_1 - \mu_2)^2 + 1, \end{aligned}$$

مورد استفاده قرار گرفته‌اند. مدل‌های حاصل از این انتخاب‌ها در جدول ۱ نشان داده شده‌اند.

جدول ۱: توزیع‌های آمیخته‌ی پواسن و نرمال با میانگین‌های برابر تحت فرض‌های

توزیع نرمال		توزیع پواسن		p	مدل
μ_2	μ_1	λ_2	λ_1		
۰/۱۱۵	-۲/۱۷۹	۵/۲۵۶	۰/۱۲۷	۰/۰۵	۱
۰/۱۶۷	-۱/۵۰۰	۵/۳۷۳	۱/۶۴۶	۰/۱۰	۲
۰/۲۸۹	-۰/۸۶۶	۵/۶۴۵	۲/۰۶۴	۰/۲۵	۳
۰/۵۰۰	-۰/۵۰۰	۶/۱۱۸	۲/۸۸۲	۰/۵۰	۴

این چهار مدل از دو توزیع آمیخته‌ی پواسن و نرمال بصورتی در نظر گرفته می‌شود که $p = 0/05, 0/25, 0/1, 0/05$ اتخاذ شود و واریانس برای مدل‌های آمیخته $1/25$ برابر واریانس تحت مدل همگن باشد و میانگین نیز برای مدل‌های آمیخته برابر با میانگین تحت مدل همگن باشد. سپس از هر یک از این مدل‌ها نمونه‌ای به حجم $n = 200$ شبیه‌سازی شده و فرض همگنی (۲) به کمک آزمون نسبت درست‌نمایی اصلاح شده و با در نظر گرفتن توابع تاوان (۶) و (۷) در سطح $0/05$ آزمون شده است. نتایج حاصل که در جدول ۲ خلاصه شده‌اند، نشان می‌دهند که وقتی مدل آمیخته به سمت شناخت ناپذیری میل می‌کند، استفاده از تابع تاوان (۷) به جای تابع تاوان (۶) در آزمون نسبت درست‌نمایی اصلاح شده، منجر به افزایش توان آزمون و دقت برآوردها با استناد به معیار میانگین توانهای دوم خطا، می‌شود. همچنین، با استفاده از رهیافت بی‌زی، پارامترهای مدل آمیخته و نیز پارامتر شکل تابع تاوان برای توزیع آمیخته نرمال و پواسن برآورد شده‌اند، جداول ۳ و ۴ مقادیر میانگین توان‌های دوم خطای برآوردهای حاصل از دو رهیافت بسامدی و بی‌زی را نشان می‌دهند. در رهیافت بسامدی از الگوریتم EM برای برآورد پارامترهای مدل آمیخته و از تابع تاوان (۷)

۲۴۰ برآورد بیزی تابع تاوان در آزمون همگنی مدل‌های آمیخته

جدول ۲: برآورد پارامترها، میانگین تاوان دوم خطا، مقدار آماره و توان آزمون نسبت درستیابی اصلاح شده، برای توابع تاوان چن و لی.

مدل	تابع تاوان	برآورد			MSE			توان آزمون
		$\hat{\lambda}_1$	$\hat{\lambda}_2$	\hat{p}	$\hat{\lambda}_1$	$\hat{\lambda}_2$	\hat{p}	
۱	چن	۰/۱۷۵۰	۵/۳۸۵	۰/۱۲۱	۱/۸۹۰	۰/۰۸	۰/۰۱۶	۸۸/۲
		۰/۳۳۳	۵/۳۰۳	۰/۰۷	۰/۴۳۹	۰/۴۳۱	۰/۰۰۳	
۲	لی	۲/۹۳۹	۵/۸۱	۰/۳۴۵	۳/۱۰۳	۰/۲۸	۰/۰۷۳	۹۲/۱
		۲/۵۷۱	۵/۶۶۴	۰/۲۶۱	۲/۱۶۶	۰/۲۰۱	۰/۰۴۸	
۳	چن	۳/۷۱۷	۶/۰۱۱	۰/۴۵۲	۰/۶۷۴	۰/۲۵۲	۰/۰۴۵	۹۳/۸
		۳/۷۱۷	۶/۰۱۱	۰/۴۵۲	۰/۶۷۴	۰/۲۵۲	۰/۰۴۵	
۴	لی	۳/۹۱۹	۶/۰۸۲	۰/۵	۰/۱۴۵	۰/۱۶۵	۰/۰۰۳	۹۳/۵
		۳/۹۱۱	۶/۰۹۰	۰/۴۹۹	۰/۱۸۲	۰/۲۱۶	۰/۰۰۵	

جدول ۳: میانگین تاوان دوم خطای برآوردگرهای بیزی و بسامدی.

مدل پواسن	برآوردگر	پارامترهای مدل		
		$\hat{\lambda}_1$	$\hat{\lambda}_2$	\hat{p}
۱	$\hat{\theta}_{ML}$	۱/۵۶۷۰	۰/۱۰۱۰	۰/۰۰۵۰
	$\hat{\theta}_B$	۰/۱۰۷۰	۰/۰۲۴۰	۰/۰۴۸۹
۲	$\hat{\theta}_{ML}$	۱/۹۱۵۰	۰/۱۸۵۰	۰/۰۳۴۰
	$\hat{\theta}_B$	۰/۰۹۹۰	۰/۰۲۱۰	۰/۰۵۴۰
۳	$\hat{\theta}_{ML}$	۰/۰۷۲۰	۰/۰۲۴۰	۰/۰۰۲۲
	$\hat{\theta}_B$	۰/۰۴۵۰	۰/۰۲۳۰	۰/۰۷۳۰
۴	$\hat{\theta}_{ML}$	۰/۰۴۰۶	۰/۰۰۳۱۰	۰
	$\hat{\theta}_B$	۰/۰۷۱۰	۰/۰۴۳۰	۰/۰۵۵۰

استفاده شده، در صورتی که در رهیافت بیزی از الگوریتم متروپولیس - هاستینگ برای برآورد پارامترهای مدل آمیخته و نیز پارامتر h استفاده شده است. همانطور که مشاهده می‌شود، زمانی که مدل آمیخته به سمت شناخت ناپذیری میل می‌کند، استفاده از رهیافت بیزی در برآورد پارامترهای مدل به مراتب از رهیافت بسامدی مطلوب‌تر است. به منظور ارزیابی میزان

جدول ۴: میانگین توان دوم خطای برآوردگرهای بیزی و بسامدی.

پارامترهای مدل				برآوردگر	مدل نرمال
\hat{h}	\hat{p}	$\hat{\mu}_2$	$\hat{\mu}_1$		
-	۰/۰۵۷۰	۰/۰۶۸۰	۲/۵۲۳۰	$\hat{\theta}_{ML}$	۱
۰/۰۰۹۴	۰/۰۰۴۴	۰/۰۰۴۵	۰/۲۴۳۰	$\hat{\theta}_B$	
-	۰/۰۱۶۷	۰/۰۲۱۰	۰/۱۵۰۵	$\hat{\theta}_{ML}$	۲
۰/۱۹۴۳	۰/۰۰۲۵	۰/۰۰۰۲	۰/۰۰۳۱	$\hat{\theta}_B$	
-	۰/۰۶۱۰	۰/۰۰۲۱	۰/۳۴۱۲	$\hat{\theta}_{ML}$	۳
۰/۰۲۹۰	۰/۰۰۲۲	۰/۰۰۰۶	۰/۰۰۸۱	$\hat{\theta}_B$	
-	۰	۰/۰۳۸۰	۰/۰۰۰۲	$\hat{\theta}_{ML}$	۴
۰/۰۱۶۵	۰/۰۰۰۹	۰/۰۲۶۰	۰/۰۰۰۴	$\hat{\theta}_B$	

برآزش مدل‌هایی که از رهیافت بیزی برآورد شده‌اند، معیار نیکویی برآزش کیش و آکائیک محاسبه و نتایج آن در جدول ۵ ارائه شده است. نکتهٔ حائز اهمیت آن است که گرچه مقادیر کوچک این معیارهای نیکویی برآزش به معنی برآزش دقیق مدل برآورد شده به داده‌ها نیست، اما مقادیر بزرگ آن گواهی روشن بر نامناسب بودن مدل برآورد شده برای داده‌ها است (هاسمر و لِمِشاو، ۲۰۰۰). معیار نیکویی برآزش کیش و آکائیک به ترتیب از روابط

$$D = -2\ell_n(\hat{\theta}, \hat{p}),$$

$$AIC = -2\ell_n(\hat{\theta}, \hat{p}) + 2t,$$

محاسبه می‌شوند، که در آن $\ell_n(\hat{\theta}, \hat{p})$ تابع لگ درست‌نمایی و t تعداد پارامترهای مدل می‌باشد. نتایج ارائه شده در جدول ۵ نشان می‌دهند که در هر دو مورد توزیع‌های آمیخته‌ی پواسن و نرمال با دو مؤلفه ($k = 2$) در مقابل مدل همگن ($k = 1$)، برآزش به مراتب بهتری به داده‌هایی دارند که از همان توزیع‌ها شبیه‌سازی شده‌اند.

جدول ۵: مقادیر معیارهای نیکویی برازش آکائیک و کیش.

توزیع پواسن		توزیع نرمال		k	مدل
$D(\theta)$	AIC	$D(\theta)$	AIC		
۱۴۴۶/۸۰۰	۱۴۴۸/۸۰۰	۹۷۵/۹۰۰	۹۷۷/۰۰۰	۱	۱
۱۳۹۶/۰۸۰	۱۴۰۴/۰۸۰	۹۵۶/۸۰۰	۹۶۴/۸۰۰	۲	
۱۳۹۸/۲۹۵	۱۴۰۰/۲۹۵	۹۳۴/۰۷۱	۹۳۹/۰۸۴	۱	۴
۱۳۹۲/۰۰۵	۱۳۰۰/۰۰۵	۹۲۷/۴۸۳	۹۳۵/۴۱۸	۲	

۶ بحث و نتیجه‌گیری

کارایی آزمون نسبت در سننمایی اصلاح شده به شدت تحت تأثیر شکل تابع تاوان انتخابی است. از طرفی هیچ دلیلی مبنی بر بهینه بودن هیچ یک از این توابع تاوانی که در منابع مختلف ارائه شده‌اند، وجود ندارد. استفاده از رهیافت بیزی هم در برآورد پارامترهای مدل آمیخته و هم برای تعیین شکل بهینه تابع تاوان خصوصاً در شرایطی که مدل آمیخته به سمت شناخت ناپذیری میل می‌کند، به مراتب از رهیافت بسامدی مطلوب‌تر است. آزمون همگنی توزیعهای آمیخته، عمدتاً به منظور تشخیص آمیخته بودن یا نبودن توزیع جامعه صورت می‌پذیرد. تعیین تعداد مؤلفه‌های توزیع آمیخته در مرحله بعدی قرار دارد.

مراجع

- Akaike, H. (1973), Information Theory and an Extension of the Maximum Likelihood Principle, *In Second International Symposium on Information Theory*, B. N. Petrov and F. Csaki (Eds). Budapest: Akademiai Kiado, 261-281.
- Chen, H., Chen, J. and Kalbfleish, D. (2001), The Likelihood Ratio Test for Homogeneity in Finite Mixture Models, *Journal of the Royal Statistical Society*, **63**, 19-29.
- Chen, J., (1998), Penalized Likelihood Ratio Test for Finite Mixture Models with Multinomial Observations, *Canadian Journal of Statistics*, **26**, 583-599.

- Chen, J. and Kalbfleisch, J. D. (1996), Penalized Minimum-Distance Estimates in Finite Mixture Models, *Canadian Journal of Statistics*, **24**, 167-175.
- Chen, J. and Khalili, A. (2006), Order Selection in Finite Mixture Models, *Working Paper 2006-03*, Department of Statistics and Actuarial Science, University of Waterloo.
- Diebolt, J. and Robert, C. (1990), Bayesian Estimation of Finite Mixture Distribution, Part i: Theoretical Aspects. *Technical Report 110*, LSTA, Université Paris VI, Paris.
- Hartigan, J. A. (1985), A Failure of Likelihood Asymptotics for Normal Mixtures, *Proceedings of conference in Honor of J. Neyman and Kiefer*, Volume 2, eds L. LeCam and R. A. Olshen, 807-810.
- Hosmer, D. W. and Lemeshow, S. (2000), *Applied Logistic Regression*, Wiley, Inc., New York.
- Li, P. Chen, J. and Marrriott, P. (2008), Non-Finite Fisher Information and Homogeneity: The EM Approach, *Biometrika*. **96**, 411-426.
- Marin, J.-M., Mengerson, K. and Robert, C. (2005), Bayesian Modelling and Inference on Mixture of Distributions, *Handbook of Statistics*, **25**, 459-507, DOI: 10.1016/S0169-7161(05)25016-2.
- Richardson, S. and Green, P. J. (1997), On Bayesian Analysis of Mixtures with an Unknown Number of Components, *Journal of the Royal Statistical Society*, B, **59**, 731-792.