



انتخاب مدل در خوشه‌بندی بلوکی به کمک انتگرال درستنمایی رده‌بندی

حمیده سادات فاطمی قمی^۱، موسی گل‌علی‌زاده^۲، منصور رزقی^۳

^۱ گروه آمار دانشگاه تربیت مدرس

^۳ گروه علوم کامپیوتر دانشگاه تربیت مدرس

چکیده: هدف خوشه‌بندی بلوکی، یافتن ساختاری به صورت بلوک‌های همگن است که موجب افزایش همزمان سطر و ستون‌ها در ماتریس داده می‌شود. به کارگیری مدل بلوکی پنهان برای این منظور به عنوان یک مدل خوشه‌بندی مدل مبنا است که در واقع یک گسترشی از مدل آمیخته متناهی است. همان‌طور که در خوشه‌بندی یک سویه، انتخاب مدل از نظر یافتن بهترین تعداد خوشه‌ها یک موضوع حیاتی به شمار می‌رود؛ در خوشه‌بندی همزمان نیز یافتن بهترین تعداد بلوک‌ها باید مورد بررسی قرار بگیرد. در این مقاله، ضابطه مبتنی بر تقریب انتگرال درستنمایی رده‌بندی برای مدل بلوکی تعمیم داده می‌شود و یک معیار اطلاع‌بیزی نیز به کمک آن بدست می‌آید. این ضابطه‌ها روی داده‌های شبیه‌سازی شده و با اندازه‌های مختلف از ماتریس و با چند دسته از تداخل خوشه‌ها مقایسه خواهند شد. واژه‌های کلیدی: خوشه‌بندی بلوکی، مدل بلوکی پنهان، انتخاب مدل، انتگرال درستنمایی رده‌بندی.
کد موضوع‌بندی ریاضی (۲۰۱۰): 62H30، 07 – 62.

۱ مقدمه

در سال‌های اخیر یکی از روش‌های مورد علاقه محققان در زمینه خوشه‌بندی، به کارگیری خوشه‌بندی همزمان^۱ در یک ماتریس داده است که در زمینه‌های مختلفی مانند تشخیص الگو، ژنتیک و تحلیل متن کاربرد دارد. خوشه‌بندی بلوکی^۲ به کمک یک مدل احتمالی به نام مدل

^۱ حمیده سادات فاطمی قمی : h.sadat.fatemi@gmail.com

^۱ Biclustering

^۲ Block clustering

بلوکی پنهان^۳ یکی از این روش‌ها است که در یک ماتریس با افراز همزمان سطر و ستون‌ها، بلوک‌های همگن (خوشه‌ها) را تشخیص می‌دهد (گووارت و ندیف (۲۰۰۸)). اما پیاده‌سازی این نوع خوشه‌بندی با ابعاد مختلفی از بلوک امکان‌پذیر است بنابراین انتخاب تعداد مناسب خوشه‌ها یک مسئله مهم شناخته شده در بحث خوشه‌بندی از جمله خوشه‌بندی همزمان به‌شمار می‌رود. تاکنون روش‌های مختلفی برای انتخاب مدل برای خوشه‌بندی همزمان به‌کار گرفته شده است. اولین رویکرد ایده خاصی از ذات مسئله خوشه‌بندی همزمان نمی‌گیرد و تنها ضابطه ارزیابی برای خوشه‌بندی یک‌سویه (نظیر ضابطه سیل‌هات) را با دوبار محاسبه نمودن برای سطر و ستون‌ها ترکیب می‌کند (اسچپر و همکاران (۲۰۰۸) و چاراد و همکاران (۲۰۱۰)). روسی و ویچی (۲۰۰۸) در روشی دیگر، واگرایی بین بلوک‌ها را به واگرایی خوشه‌بندی بلوکی تقسیم کردند. همچنین برخی ضابطه‌های اطلاع برای انتخاب مدل، برای خوشه‌بندی بلوکی به‌کار گرفته شده اند. ون دیجک و همکاران (۲۰۰۹) یک ضابطه آکائیک اصلاح شده به‌این منظور معرفی نمودند. اما در سال‌های اخیر نیز ضابطه مبتنی بر انتگرال درست‌نمایی رده‌بندی (بیرناکی و همکاران (۲۰۰۰)) تعمیم داده شده است که در این مقاله روی این ضابطه متمرکز خواهیم شد. ساختار مقاله حاضر به‌صورت زیر تدوین شده است. در بخش ۲ مدل بلوکی پنهان تشریح می‌شود. انتخاب مدل براساس انتگرال درست‌نمایی رده‌بندی در بخش ۳ می‌آید و مطالعه شبیه‌سازی تشکیل‌دهنده بخش ۴ است.

۲ مدل بلوکی پنهان

در خوشه‌بندی همزمان یا بلوکی، هدف یافتن افراز همزمان سطرها و ستون‌ها (z, w) است. در این روش z یک افرازی است که موجب خوشه‌بندی سطرها (I) به خوشه می‌شود و w نیز یک افراز روی مجموعه ستون‌ها (J) به خوشه است. در این حالت توزیع داده‌ها را می‌توان با تابع چگالی احتمال $f(x, \theta) = \sum_{u \in U} p(u) f(x|u, \theta)$ نشان داد که در آن U نشان‌دهنده کلیه افرازه‌های ممکن روی $I \times J$ است. اگر مجموعه افرازه‌ها روی I و J مستقل از هم فرض شوند، در این حالت $p(u) = p(z)p(w)$ و لذا

$$f(x; \theta) = \sum_{(z,w) \in Z \times W} p(z)p(w) f(x|z, w; \alpha), \quad (1.2)$$

که در آن Z و W به‌عنوان متغیرهای پنهان، مجموعه همه افرازه‌های ممکن I به خوشه J و خوشه J به خوشه m تشکیل می‌دهند. با توجه به فرض استقلال مکانی^۴ مربوط به مدل کلاس پنهان^۵ تابع چگالی با شرطی شدن روی z و w به‌صورت

$$f(x|z, w; \theta) = \prod_{i,j,k,l} f(x_{ij}; \alpha_{kl})^{z_{ik}w_{jl}}.$$

بدست می‌آید. فرض می‌شود $\rho = (\rho_1, \dots, \rho_m)$ و $\pi = (\pi_1, \dots, \pi_g)$ بردارهای احتمال‌های ρ_k و π_l باشند و این احتمال پیشین را بیان کنند که سطر و ستون به‌ترتیب متعلق به k امین مولفه و l امین مولفه باشد. با توجه به (۱.۲) تابع درست‌نمایی برای مدل بلوکی به‌صورت

$$f(x; \theta) = \sum_{(z,w) \in Z \times W} \prod_{i,k} \pi_k^{z_{ik}} \prod_{j,l} \rho_l^{w_{jl}} \prod_{i,j,k,l} f(x_{ij}; \alpha_{kl})^{z_{ik}w_{jl}} \quad (2.2)$$

بدست می‌آید. از آنجایی که مدل شامل متغیرهایی پنهان برای تبیین افرازه‌ها هستند؛ روند الگوریتم EM که شامل فرآیندی متناوب بین گام امیدگیری (گام E) و گام ماکسیمم‌سازی (گام M) است؛ به ذهن می‌رسد. لگاریتم درست‌نمایی رده‌بندی با استفاده از داده کامل می

^۳ Latent Block Model (LBM)

^۴ Local independence

^۵ Latent class model

تواند به صورت زیر نوشته شود:

$$L_C(x, z, w, \theta) = \sum_{i,k} z_{ik} \log \pi_k + \sum_{j,l} w_{jl} \log \rho_l + \sum_{i,j,k,l} z_{ik} w_{jl} \log f(x_{ij}; \alpha_{kl}).$$

اما روند الگوریتم EM برای LBM به دلیل ضرورت محاسبه $P(Z_{ik}W_{jl} = 1 | \theta^{(c)}, X = x)$ غیرقابل اجراست. برای رفع این مشکل **گووارت و ندیف (۲۰۰۸)** یک تقریب به نام الگوریتم EM تغییرگرا^۶ را پیشنهاد دادند. بنا به این الگوریتم در گام E در این الگوریتم توزیع توأم $p(z, w | X, \theta)$ برابر با حاصلضرب $p(z | X, \theta)p(w | X, \theta)$ خواهد شد و ضابطه

$$F(P, \theta) = E_P(L_C(z, w, \theta)) + H(P) \quad (۳.۲)$$

مورد ماکسیم سازی قرار می‌گیرد. به پیشنهاد **کریبین و همکاران (۲۰۱۲)** بعد از همگرایی، مقدار F که یک کران پایین لگاریتم درست‌نمایی است؛ می‌تواند به عنوان تقریبی از آن در نظر گرفته شود. لازم به ذکر است پیاده‌سازی این الگوریتم مبتنی بر LBM مستلزم اجرای گام‌هایی متناوب برای دو الگوریتم خوشه‌بندی روی سطرها ستون‌ها مبتنی بر مدل‌های آمیخته متناهی است.

۳ انتخاب مدل

موضوع انتخاب مدل در LBM به چند دلیل مسئله‌ساز است. اولین دلیل این است که یک زوج از تعداد خوشه‌ها یعنی (g, m) به جای یک عدد تنها باید انتخاب شود. دومین دلیل این است که ضابطه‌های درست‌نمایی جریمه‌ای موجود نظیر AIC و BIC که برای انتخاب مدل در خوشه‌بندی مدل مینا به کار گرفته می‌شوند؛ نمی‌توانند به طور مستقیم مورد استفاده قرار گیرند. سوم اینکه انتخاب مدل معمولاً به تعداد واحدهای آماری مدل وابسته است، لذا گزینش این موضوع در LBM سؤال برانگیز است. به زبانی دیگر انتخاب تعداد واحد آماری مدل باید از تعداد مشاهدات n یا تعداد متغیرها d یا تعداد سلول‌ها صورت گیرد. در ادامه با در نظر گرفتن این مسئله، انتخاب مدل بررسی می‌شود. به طور کلی، هدفی که در بحث انتخاب مدل دنبال می‌شود؛ رسیدن به بهترین مدل است. در چارچوب بیزی بهترین مدل، محتمل‌ترین آنهاست. درحقیقت در میان مجموعه متناهی مدل، مدلی انتخاب می‌شود که بیشترین احتمال پسین را دارا باشد. به زبانی ساده

$$M_h = \arg \max_{M \in \mathcal{M}} p(M|X).$$

اگر فرض شود که احتمال پیشین $p(M)$ برای هر مدل، ثابت باشد؛ آنگاه ماکسیم نمودن $p(M|X)$ معادل با ماکسیم نمودن $p(X|M)$ است و سپس هدف یافتن مدلی می‌شود که انتگرال درست‌نمایی را ماکسیم کند. لذا

$$M_h = \arg \max_M p(X|M),$$

طوری که $p(X|M) = \int_{\Theta} p(X|M; \theta)p(\theta|M)d\theta$. در اینجا Θ فضای پارامتر و $p(\theta|M)$ احتمال پیشین پارامتر θ برای مدل M است. در حالاتی که مشاهدات مستقل و هم‌توزیع باشند؛ ضابطه BIC یکی راهی برای تقریب مجانبی از لگاریتم انتگرال درست‌نمایی است. فرم استاندارد این ضابطه به صورت

$$BIC(M) = \log L(x; \hat{\theta}_{ML}) - \frac{1}{n} D_M \log(n)$$

^۶ Variational EM (VEM)

است که در آن D_M بعد مدل، n تعداد مشاهده و $L(y; \hat{\theta}_{ML})$ ماکسیمم درست‌نمایی است. اما این تقریب به دلیل ساختار وابستگی مشاهدات بین سطرها و ستون‌ها نمی‌تواند در مدل بلوکی، به‌کار گرفته شود. برای حل این مشکل، تعریف دیگری از انتخاب مدل،

$$\text{به صورت } M_c = \arg \max_M P(X, z, w|M) \text{ ارائه شد که در آن}$$

$$P(X, z, w|M) = \int_{\Theta} p(X, z, w|\theta; M) p(\theta|M) d\theta.$$

در نهایت ضابطه مدنظر، با لگاریتم گرفتن به صورت $ICL(M) = \log P(X, z, w|M)$ تعریف شده و به اختصار انتگرال درست‌نمایی کامل (یا رده‌بندی) ^۲ نامیده می‌شود (لومت و همکاران (۲۰۱۲)). همچنین تحت فرض استقلال برچسب‌ها داریم:

$$ICL(M) = \log p(X|z, w, M) + \log p(z|M) + \log p(w|M). \quad (۱.۳)$$

در ادامه یک تقریب مجانبی از ICL معرفی می‌شود. مشاهده می‌شود که این تقریب صریحاً به توزیع پیشین پارامترها بستگی ندارد. همچنین با الهام از آن ضابطه مشابه BIC با نماد BIC_{like} برای مدل بلوکی ارائه می‌شود.

۱.۳ تقریب ICL و بدست آوردن ضابطه BIC در مدل بلوکی

تحت فرض استقلال پیشین پارامترها، یعنی $p(\theta) = p(\alpha)p(\pi)p(\rho)$ ، بخش اول ICL به صورت

$$\log p(X|z, w, M) = \int_A p(X|z, w, M, \alpha) p(\alpha|M) d\alpha \underbrace{\int_{P \times R} p(\pi|M) p(\rho|M) d\pi d\rho}_{=1}$$

محاسبه می‌شود. این بخش می‌تواند با استفاده از BIC_{like} تقریب زده شود. یعنی

$$\log p(X|z, w, M) = \max \log p(X|z, w, M, \alpha) - \frac{\lambda}{4} \log(nd)$$

که در آن λ بعد بردار پارامتر α در فضای A است. طبق $p(z|\pi) = \prod_{i,k} \pi_k^{z_{ik}}$ و در نظر گرفتن توزیع ناآگاهی بخش دیریکله $D(a, \dots, a)$ به‌عنوان توزیع پیشین برای π ، و با استفاده از قانون بیز، $p(z)$ به صورت

$$p(z) = \frac{p(z|\pi)p(\pi)}{p(\pi|z)} = \frac{\Gamma(ga) \prod_k \Gamma(z_{.k} + a)}{\Gamma(a)^g \Gamma(n + ga)}$$

بدست می‌آید و با همین روند، $p(w)$ محاسبه می‌شود. با قرار دادن $a = \frac{1}{4}$ برای یک توزیع ناآگاهی بخش جفریز (رابرت (۲۰۰۱))، لگاریتم توزیع‌های پیشین به صورت

$$\log p(z|M) = \log \Gamma\left(\frac{g}{4}\right) + \sum_{k=1}^g \log \Gamma\left(n_k + \frac{1}{4}\right) - g \log \Gamma\left(\frac{1}{4}\right) - \log \Gamma\left(n + \frac{g}{4}\right),$$

$$\log p(w|M) = \log \Gamma\left(\frac{m}{4}\right) + \sum_{l=1}^m \log \Gamma\left(d_l + \frac{1}{4}\right) - m \log \Gamma\left(\frac{1}{4}\right) - \log \Gamma\left(d + \frac{m}{4}\right).$$

خواهند بود. اگر n_k و d_l به اندازه کافی بزرگ باشند؛ با به‌کارگیری تقریب استرلینگ $\Gamma(t+1) \approx t^{t+1/2} \exp(-t) (\sqrt{\pi})^{1/2}$ تابع گاما و نادیده گرفتن بخش‌هایی که از $O(1)$ هستند؛ لگاریتم توزیع پیشین به صورت زیر حاصل خواهند شد:

$$\log p(z|M) \approx \sum_{k=1}^g n_k \log n_k - n \log n - \frac{1}{4}(g-1) \log n$$

^۲ Integrated Complete(Classification) Likelihood (ICL)

$$\log p(w|M) \approx \sum_{l=1}^m d_l \log d_l - d \log d - \frac{1}{\gamma} (m-1) \log m.$$

علاوه بر این، به کمک تساوی $p(z|\pi) = \prod_k \pi_k^{n_k}$ و شرط $\sum_k \pi_k = 1$ رابطه $\sum_{k=1}^g n_k \log \frac{n_k}{n} = \max_{\pi} \log p(z|\pi, M)$ و به طور مشابه $\sum_{l=1}^m d_l \log \frac{d_l}{d} = \max_{\rho} \log p(w|\rho, M)$ بدست خواهند آمد. با جای‌گذاری پارامتر θ و بردار برچسب‌های (z, w) از طریق الگوریتم *VEM* سرانجام ضابطه *ICL* به صورت زیر حاصل خواهد شد:

$$ICL_{BIC}(M) = L_c(X, \hat{z}, \hat{w}|M, \hat{\theta}) - \frac{\lambda}{\gamma} \log(nd) - \frac{g-1}{\gamma} \log n - \frac{m-1}{\gamma} \log d.$$

این ضابطه یک درست‌نمایی رده‌بندی جریمه‌ای محسوب می‌شود که در آن بخش جریمه جمع پارامترهای مرتبط با توزیع‌ها و همچنین پارامترهای مربوط به نسبت‌های سطری و نسبت‌های ستونی هستند. با الهام از رابطه $ICL = BIC + H$ در مدل آمیخته کلاسیک، که H تنها یک بخش آنتروپی مربوط به احتمال‌های شرطی برچسب‌هاست؛ ضابطه BIC_{like} به صورت

$$BIC_{like} = \log p(X|M, \theta^*) - \frac{\lambda}{\gamma} \log(nd) - \frac{g-1}{\gamma} \log n - \frac{m-1}{\gamma} \log d \quad (2.3)$$

حاصل می‌شود.

۲.۳ محاسبه ضابطه دقیق *ICL*

در این بخش ضابطه دقیقی برای تعریف (۱.۳) بدست می‌آید. این ضابطه نیازی به هیچ تعریف مجانبی ندارد؛ اما توزیع پیشین برای پارامترها باید به‌طور مناسبی در نظر گرفته شوند. این موضوع برای دو حالت دنبال می‌شوند:

الف) ضابطه دقیق *ICL* برای داده دودویی

توزیع‌های پیشین برای پارامترهای مربوط به متغیر برنولی به صورت $\alpha_{kl} \sim B(b, b)$ و $\rho \sim D(a, \dots, a)$ ، $\pi \sim D(a, \dots, a)$ در نظر گرفته می‌شود. لذا،

$$\begin{aligned} ICL_{(a,b)}(g, m) &= \log p(x, z, w) \\ &= \log \Gamma(ga) - (m+g) \log \Gamma(a) + mg(\log \Gamma(2b) - 2 \log \Gamma(b)) \\ &\quad - \log \Gamma(n+ga) - \log \Gamma(d+ma) + \sum_k \log \Gamma(z_{.k} + a) + \sum_l \log \Gamma(w_{.l} + a) \\ &\quad + \sum_{k,l} \left[\left(\sum_{h=1}^2 \log \Gamma(N_{kl}^h + b) \right) - \log \Gamma(z_{.k} w_{.l} + rb) \right]. \end{aligned}$$

طبق کریبین و همکاران (۲۰۱۴) در نظر گرفتن ابرپارامترهای $a = 4$ برای تعداد خوشه‌های کمتر از ۸ و نیز مقدار $a = 16$ برای ابعاد بزرگتر و همچنین در نظر گرفتن ابرپارامتر $b = 1$ نتایج بهتری را حاصل خواهد کرد.

ب) محاسبه ضابطه دقیق *ICL* برای داده پیوسته

با فرض متفاوت بودن واریانس‌های هر بلوک، یک توزیع پیشین نرمال با پارامترهای $(\mu_{\circ}, \frac{\sigma_{kl}^2}{\kappa_{\circ}})$ برای میانگین شرطی هر بلوک یعنی $\mu_{kl} | \sigma_{kl}^2$ در نظر گرفته می‌شود. همچنین یک توزیع پیشین معکوس‌خی‌دوی مقیاس شده برای پارامتر واریانس هر بلوک (σ_{kl}^2) ، لحاظ می‌شود. به‌زبانی دقیق‌تر، $\sigma_{kl}^2 \sim \text{Scale-Inv} - \chi^2(v_{\circ}, \sigma_{\circ}^2)$. آنگاه بخش اول معادله (۱.۳) به صورت

$$\log p(X|z, w, M) = -\frac{nd}{\gamma} \log \pi + \frac{gm v_{\circ}}{\gamma} \log(v_{\circ} \sigma_{\circ}^2) - gm \log \Gamma\left(\frac{v_{\circ}}{\gamma}\right) + \frac{gm}{\gamma} \log \kappa_{\circ}$$

$$-\sum_{k,l} \frac{v_0 + n_k d_l}{\nu} \log(v_0 \sigma_0^2 + (n_k d_l - 1) s_{kl}^{*2} + \frac{\kappa_0 n_k d_l}{\kappa_0 + n_k d_l} (\bar{x}_{kl} - \mu_0)^2) + \sum_{k,l} \log \Gamma\left(\frac{v_0 + n_k d_l}{\nu}\right) - \frac{1}{\nu} \log(\kappa_0 + n_k d_l).$$

خواهد شد که در آن n_k و d_l به ترتیب تعداد سطرها و ستون‌ها در بلوک (k, l) $\bar{x}_{kl} = 1/(n_k d_l) \sum_{i,j} z_{ik} w_{jl} x_{ij}$ میانگین نمونه‌ای، $s_{kl}^{*2} = 1/(n_k d_l - 1) \sum_{i,j} z_{ik} w_{jl} (x_{ij} - \bar{x}_{kl})^2$ واریانس نمونه‌ای نارایب هستند. برای پارامترهای نسبت (π, ρ) ، با فرض اینکه درایه‌های آن از همدیگر متفاوت هستند، یک توزیع پیشین دیریکله با پارامترهای (a, \dots, a) در نظر گرفته می‌شود. بنابراین دویخس آخر معادله (۱.۳) به ترتیب به صورت زیر تغییر می‌یابند:

$$\log p(z|M) = \log \Gamma(ga) + \sum_k \log \Gamma(n_k + a) - g \log \Gamma(a) - \log \Gamma(n + ga)$$

$$\log p(w|M) = \log \Gamma(ma) + \sum_l \log \Gamma(d_l + a) - m \log \Gamma(a) - \log \Gamma(d + ma).$$

برخلاف ضابطه‌های تقریبی BIC و ICL ، ضابطه دقیق ICL به توزیع پیشین پارامترها بستگی دارد. اگرچه در نظر گرفتن این اطلاعات پیشین می‌تواند در استنباط \hat{z} و \hat{w} و بنابراین تخمین پارامترهای مدل اثر بگذارد؛ با این حال در شبیه‌سازی‌ها ترجیح داده شده که از برآورد ماکسیمم درست‌مابایی استفاده می‌شود. گرچه این انتخاب غیر متعارف است و ممکن است برای ضابطه دقیق ICL مناسب نباشد؛ ولی بنا به لومت و همکاران (۲۰۱۴). اجازه می‌دهد تا با ضابطه‌های انتخاب مدل با برآورد پارامترهای یکسان مقایسه شود. ضابطه ICL به ۵ ابرپارامتر $(a, \mu_0, \kappa_0, \nu_0, \sigma_0^2)$ نیاز دارد. برای توزیع دیریکله مربوط به پارامترهای نسبت‌ها، یک پیشین ناآگاهی بخش جفریز یعنی $a = \frac{1}{\nu}$ انتخاب خواهد شد. برای پارامتر مربوط به میانگین‌های هر بلوک یعنی μ_0 ، میانگین کلی ماتریس داده در نظر گرفته می‌شود. ابرپارامتر κ_0 نیز برابر $\sqrt{nd/gm}$ لحاظ می‌شود. واضح است که این مقدار با بالارفتن تعداد درایه‌ها در هر خوشه، افزایش می‌یابد که با این انتخاب، پیچیدگی مسئله لحاظ شده است و در واقع گویای این نکته نیز هست که با افزایش بعد داده نسبت خطا کاهش می‌یابد. طبق لومت و همکاران (۲۰۱۴) به صورت تجربی نشان داده شده گرفتن ریشه دوم از این نسبت، نتیجه بهتری را حاصل خواهد کرد. نسبتی که برای واریانس پیشین میانگین هر بلوک انتخاب می‌شود؛ برای پیشین مربوط به واریانس هر بلوک، ابرپارامتر ν_0 و σ_0^2 را طوری انتخاب می‌کند که مد توزیع معکوس‌خی‌دو مقیاس شده، یعنی $(\frac{\nu_0 \sigma_0^2}{\nu_0 + \frac{1}{\nu}})$ برابر واریانس نمونه‌ای ماتریس داده باشد. در نتیجه با در نظر گرفتن مقدار 10^{-1} به منظور ایجاد واگرایی بیشتر، مقدار σ_0^2 به صورت 2×10^{-2} در نظر گرفته می‌شود که در اینجا s^{*2} واریانس نمونه‌ای نارایب است.

۴ مطالعه شبیه‌سازی

برای g و m های مختلف، الگوریتم خوشه‌بندی بلوکی روی ۲۰ داده شبیه‌سازی شده با ابعاد بلوک صحیح 2×3 از نوع دودویی اجرا شدند تا عملکرد ضابطه‌های ICL_{BIC} ، BIC_{like} و ICL واقعی مورد مقایسه قرار بگیرند. در جدول ۱، اندازه ماتریس برای سه بعد مختلف و دو دسته از تداخل خوشه‌ها روی داده‌های دودویی در نظر گرفته شده است. بدیهی است تعیین پارامترها در میزان تداخل خوشه‌ها موثر خواهد بود و این تداخل براساس نسبت خطا بین دو افراز برآورد شده تعیین می‌شود. نسبت خطای حاصل از الگوریتم خوشه‌بندی، یعنی $u = (z, w)$ و افراز صحیح $u' = (z', w')$ به صورت $u' = (z', w')$ به صورت $\delta(u, u') = 1 - \frac{1}{nd} \sum_{i,j,k,l} z_{ik} w_{jl} z'_{ik} w'_{jl}$ و به عنوان

میانگین تعداد خطای رده بندی های اشتباه[^] تعریف می شود. بازه \bar{d} برای این دو دسته از داده با تداخل نسبتاً تفکیک شده $(++)$ و تفکیک شده ضعیف $(+++)$ به ترتیب بین $[0/10, 0/10]$ و $[0/12, 0/22]$ قرار گرفته اند. اعداد داخل جدول تعداد تشخیص ابعاد بلوک در حالت های مختلف و تحت ضابطه های مربوطه را نشان می دهند.

جدول ۱: نتایج حاصل از مطالعه شبیه سازی

			+++			++																		
			۲۰۰ × ۱۲۰			۷۰ × ۳۰			۵۰۰ × ۲۰۰			۲۰۰ × ۱۲۰			۷۰ × ۳۰									
<i>ICL_{BIC}</i>	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲						
		۴	۲			۲			۱۱	۲			۲			۱۹	۳			۱	۴	۹	۳	
			۱۶	۳			۲	۱۳	۳			۹	۳			۲۰	۳			۱	۴	۹	۳	
			۴			۱	۱	۳	۴				۴				۴				۱	۵	۴	
<i>BIC_{like}</i>	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲						
			۱۴	۲			۲			۱۸	۲			۲			۱۶	۲			۲	۲	۲	
			۶	۳			۳	۱۳	۳			۲	۳			۲۰	۳			۱	۴	۹	۳	
			۴			۱	۳	۴				۴				۴				۴	۴	۴		
<i>ICL دقیق</i>	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲	۴	۳	۲						
			۱	۲			۲			۹	۲			۲			۲			۱	۴	۹	۳	
			۱۶	۳			۱	۲	۱۳	۳			۱	۱۰	۳			۲۰	۳			۱	۴	۹
			۴			۱	۳	۴				۴				۴				۱	۵	۴		

همانطور که ملاحظه می شود؛ کیفیت عملکرد ضابطه های با افزایش اندازه ماتریس و کاهش میزان تفکیک پذیری مولفه ها افزایش می یابد. در بیشتر موارد ضابطه *ICL_{BIC}* قادر به انتخاب مدل درست است؛ درحالی که *BIC_{like}* تمایل به کم برآورد تعداد خوشه ها دارد؛ به خصوص هنگامی که ابعاد ماتریس کوچک تر باشد. دلیلی که می توان برای این امر ارائه کرد؛ این است که بعد ماتریس نمی تواند برای اعتبار استفاده از تقریب مجانبی *BIC_{like}* کافی باشد. به بیانی دیگر در این ضابطه به دلیل غیر دسترس بودن ماکسیمم درست نمایی از کران پایین آن به عنوان تقریبی از ماکسیمم انرژی آزاد استفاده شد که بر تفاوت رفتار درست نمایی و کران پایین آن با افزایش تعداد مولفه ها تاکید می کند.

بحث و نتیجه گیری

مدل بلوکی پنهان یک مدل آماری به صرفه از نظر برآورد تعداد پارامترها برای سازماندهی یک ماتریس داده به بلوک های همگن است. مسئله انتخاب تعداد مناسب خوشه های سطری و ستونی یک چالشی برای این ساختارها است که مقاله حاضر بر روی تعمیم ضابطه مبتنی بر انتگرال درست نمایی رده بندی متمرکز شده است. اما در این میان یکی از مشکلات برای داده دودویی و به طور کلی تر داده های رده ای، امکان روبرو شدن با پدیده خالی بودن برخی خوشه ها به خصوص در ابعاد کوچک تر است. به نظر می رسد رهیافت بیزی می تواند تا حدی از این مسئله جلوگیری کند که به عنوان مطالعات آینده مدنظر نویسندگان مقاله است.

[^]Misclassification

مراجع

- Biernacki, C., Celeux, G. and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(7), 719-725.
- Charrad, M., Lechevallier, Y., Saporta, G. and Ahmed, M. B. (2010). Détermination du nombre de classes dans les méthodes de bipartitionnement, 17ème Rencontres de la Société Francophone de Classification, 119-122.
- Govaert, G., and Nadif, M. (2008). Block clustering with Bernoulli mixture models: Comparison of different approaches, *Computational Statistics and Data Analysis*, 52(6), 3233-3245.
- Keribin, C., Brault, V., Celeux, G. and Govaert, G. (2012). Model selection for the binary latent block model, In COMPSTAT 2012 20th International Conference on Computational Statistics.
- Keribin, C., Brault, V., Celeux, G. and Govaert, G. (2014). Estimation and selection for the latent block model on categorical data, *Statistics and Computing*, 1-16.
- Lomet, A., Govaert, G. and Grandvalet, Y. (2012). Model selection in block clustering by the integrated classification likelihood, In 20th International Conference on COMPUTATIONAL STATISTICS (COMPSTAT 2012), 519-530.
- Lomet, A., Govaert, G. and Grandvalet, Y. (2014). Model selection for Gaussian latent block clustering with the integrated classification likelihood, *Advances in Data Analysis and Classification*, 1-20.
- Robert C (2001). *The Bayesian choice*, Springer, Berlin.
- Rocci, R. and Vichi, M. (2008). Two-mode multi-partitioning, *Computational Statistics and Data Analysis*, 52(4), 1984-2003.
- Schepers, J., Ceulemans, E. and Van Mechelen, I. (2008). Selecting among multi-mode partitioning models of different complexities: A comparison of four model selection criteria, *Journal of Classification*, 25(1), 67-85.
- van Dijk, B., van Rosmalen, J. and Paap, R. (2009). A Bayesian approach to two-mode clustering, Technical Report, Econometric Institute Research Papers.