



## Clustering Based on Nonparanormal Graphical Mixture Models

Haji Aghabozorgi, H. , Eskandari, F. 

Department of Statistics, Allameh Tabatabaei University, Tehran, Iran.

**Corresponding author:** F. Eskandari, askandari@atu.ac.ir

**Received:** 3 December 2021 **Revised:** 25 July 2022 **Accepted and Published Online:** 1 August 2022.

### Introduction

Graphical mixture models provide a powerful tool to visually depict the conditional independence relationships between high-dimensional heterogeneous data. In the study of graphical mixture models, the distribution of the mixture components is mostly considered multivariate normal with different covariance matrices. The resulting model is the Gaussian graphical mixture model (GGMM). The nonparanormal graphical mixture model (NGMM) has been introduced by replacing the normal assumption with a semiparametric Gaussian copula, which extends the nonparanormal graphical model and mixture models. This study proposes clustering based on NGMM under two forms of  $\ell_1$  penalty functions. Its performance is compared with clustering based on GGMM, in terms of cluster reconstruction and parameters estimation.

### Material and Methods

The clustering based on NGMM is performed via a penalized EM algorithm under conventional and unconventional forms of  $\ell_1$  penalty functions (denoted by  $NGMM_0$  and  $NGMM_1$ , respectively) and its performance over Gaussian and non-Gaussian simulated data sets are compared with the Gaussian ones (represented by  $GGMM_0$  and  $GGMM_1$ , respectively). Along with the conventional  $\ell_1$  penalty, an alternative, unconventional penalty term is considered, which depends on the mixture proportions. Thus, the choice of mixture model distribution (Gaussian or nonparanormal) along with the choice of penalty function has emerged as the primary key of comparison. To

better compare the studied methods in terms of robustness against outliers, we considered deterministic and random contamination mechanisms. The proposed methodology is applied to Wisconsin diagnostic breast cancer data set to diagnose malignant or benign cancer patients.

### Results and Discussion

The results of the simulation study on normal and nonparanormal datasets in ideal and noisy settings, as well as the application of breast cancer data set, showed that clustering approaches based on NGMM ( $NGMM_0$  and  $NGMM_1$ ) are more efficient and robust in the recovery of true cluster assignments than the clustering based on GGMM ( $GGMM_0$  and  $GGMM_1$ ), whereas, the unconventional PMLEs ( $GGMM_1$  and  $NGMM_1$ ) are more efficient in estimating the elements of precision matrices than the conventional PMLEs ( $GGMM_0$  and  $NGMM_0$ ).

### Conclusion

The performance of clustering methods depends on the choice of penalty function and model selection, such that the combination of the nonparanormal graphical mixture model and the penalty term depending on the mixing proportions ( $NGMM_1$ ) is more accurate than Gaussian ones in terms of cluster reconstruction and parameters estimation.

**Keywords:** Clustering, Graphical mixture models, Nonparanormal distribution, Penalized log-likelihood.

**Mathematics Subject Classification (2010):** 62H30, 62H20.





مجله علوم آماری، بهار و تابستان ۱۴۰۱

جلد ۱۶، شماره ۱، ص ۶۳ - ۸۹

DOI: 10.29252/jss.16.1.63

مقاله پژوهشی

## خوشه‌بندی مبتنی بر مدل‌های آمیخته گرافی نرمال ناپارامتری

حمید حاجی‌آقابزرگی، فرزاد اسکندری

گروه آمار، دانشگاه علامه طباطبایی (ره)

نویسنده مسئول: فرزاد اسکندری، askandari@atu.ac.ir

تاریخ دریافت: ۱۴۰۰/۰۹/۱۲ تاریخ بازنگری: ۱۴۰۱/۰۵/۰۳ تاریخ پذیرش و انتشار: ۱۴۰۱/۰۵/۱۰

**چکیده:** مدل‌های آمیخته گرافی، ابزاری قدرتمند برای نمایش دیداری روابط استقلال شرطی بین داده‌های ناهمگن بالا بعد فراهم کرده است. در مطالعه این مدل‌ها، اغلب توزیع مولفه‌های آمیخته، نرمال چندمتغیره با ماتریس‌های کواریانس متفاوت در نظر گرفته شده که مدل حاصل، به مدل آمیخته گرافی گاوسی معروف است. با جای‌گزین کردن فرض محدودکننده نرمال با یک مفصل نیمه‌پارامتری نرمال، مدل آمیخته گرافی نرمال ناپارامتری معرفی شده که هم مدل گرافی نرمال ناپارامتری و هم مدل‌های آمیخته را تعمیم داده است. در این مطالعه، خوشه‌بندی مبتنی بر مدل آمیخته گرافی نرمال ناپارامتری با دو فرم تابع تاوان  $l_1$  (متعارف و نامتعارف) پیشنهاد شده است و عملکرد آن با روش خوشه‌بندی مبتنی بر مدل آمیخته گرافی گاوسی مقایسه شده است. نتایج مطالعه شبیه‌سازی روی داده‌های نرمال و غیرنرمال، در حضور و عدم حضور داده‌های دورافتاده و همچنین نتایج کاربردی روی داده‌های سرطان سینه نشان داد که ترکیب مدل آمیخته گرافی نرمال ناپارامتری با تابع تاوان وابسته به نسبت‌های آمیخته، از نظر بازسازی خوشه‌ها و برآورد پارامترهای مدل، نسبت به سایر روش‌های خوشه‌بندی مبتنی بر مدل از دقت بالاتری برخوردار است.

**واژه‌های کلیدی:** توزیع نرمال ناپارامتری، خوشه‌بندی، لگاریتم درست‌نمایی تاوانیده، مدل‌های آمیخته گرافی.

کد موضوع‌بندی ریاضی (۲۰۱۰): 62H30، 62H20.



©نویسندگان). ناشر انجمن آمار ایران است.

این مقاله با دسترسی آزاد تحت شرایط و ضوابط (CC BY-NC 4.0) توزیع شده است.

## ۱ مقدمه

مدل‌های گرافی بی‌سو<sup>۱</sup> برای توصیف روابط وابستگی یا استقلال شرطی بین متغیرها در داده‌های بالابعد<sup>۲</sup> به‌کار می‌رود. در این مدل‌ها، تعیین توزیع توأم بردار تصادفی  $\mathbf{X} = (X_1, \dots, X_p)$  متناظر با راس‌های یک گراف  $G$ ، چالشی مهم است. وقتی متغیرهای تصادفی متناظر با راس‌ها پیوسته باشند، توزیع آن‌ها اغلب به‌عنوان توزیع نرمال چند متغیره در نظر گرفته می‌شود ( $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$ ) و مدل حاصل، به مدل گرافی گاوسی<sup>۳</sup> (GGM) تعبیر می‌شود. در این حالت، می‌توان ماتریس دقت<sup>۴</sup>  $\Omega = \Sigma^{-1}$  را به‌طور مستقیم به یک مدل گرافی گاوسی تعبیر کرد. بنابراین، مساله برآورد مدل گرافی گاوسی به برآوردیابی ماتریس دقت محدود می‌شود. الگوریتم‌های بسیاری نظیر ماکسیمم درستنمایی  $\ell_1$ -منظم شده<sup>۵</sup>، انتخاب همسایگی<sup>۶</sup> و الگوریتم لاسوی گرافی<sup>۷</sup>، برای برآورد ماتریس دقت در مدل گرافی گاوسی تعمیم داده شده است (مینشاوسن و بولمن، ۲۰۰۶؛ فریدمن و همکاران، ۲۰۰۸).

در عمل ممکن است فرض نرمال بودن برقرار نباشد که در این حالت، از تبدیل کردن داده‌ها (با تبدیل‌هایی نظیر تبدیل لگاریتمی) برای بهره‌مندی از ویژگی‌های مدل گرافی گاوسی استفاده می‌شود. بدین منظور، باید همه تبدیل‌ها را بررسی کرد تا تبدیلی که داده‌ها را به نرمال نزدیک می‌کند، حاصل شود. از دیدگاه مدل‌سازی، راه حل ساده‌تر این است که توابع تبدیل، نامعلوم فرض شوند و با روش‌های مناسب برآورد شوند. در این رابطه، فرض نرمال توسط لیو و همکاران (۲۰۰۹) با یک مفصل نیمه‌پارامتری گاوسی<sup>۸</sup> جای‌گزین شده است. به‌طور مشخص، بردار تصادفی  $\mathbf{X} = (X_1, \dots, X_p)$  با بردار تصادفی تبدیل‌شده  $f(\mathbf{X}) = (f_1(X_1), \dots, f_p(X_p))$  جای‌گزین شده است که در آن  $f(\mathbf{X}) \sim N(\mathbf{0}, \Sigma)$ . این رویکرد به یک تعمیم توزیع نرمال منجر شد که توزیع نرمال ناپارامتری<sup>۹</sup> نامیده شده است. مدل گرافی نرمال ناپارامتری<sup>۱۰</sup> (NGM) انعطاف‌پذیرتر از GGM است و این ویژگی GGM را حفظ می‌کند که روابط استقلال شرطی بین متغیرها در ماتریس دقت  $\Omega = \Sigma^{-1}$  در یک فرم گرافی کدگذاری می‌شود. در نتیجه، برآورد NGM به برآورد توابع تبدیل تک متغیره  $\{f_j\}_{j=1}^p$  و ماتریس دقت  $\Omega$  بستگی دارد.

در ادبیات مدل‌های گرافی، غالباً فرض بر آن است که داده‌های مشاهده‌شده از یک منبع همگن آمده‌اند و از مدل گرافی گاوسی پارامتری (یوان و لین، ۲۰۰۷)، مدل گرافی نرمال ناپارامتری (لیو و همکاران، ۲۰۰۹) یا مدل گرافی نرمال ناپارامتری SKEPTIC (لیو و همکاران، ۲۰۱۲) پیروی می‌کنند. با این وجود، مشاهده‌ها معمولاً از منبع‌های متفاوت می‌آیند و در طیف گسترده‌ای از کاربردهای دنیای واقعی، ممکن است وابستگی‌های ساختاری

<sup>1</sup>Undirected graphical models

<sup>2</sup>High-dimensional data

<sup>3</sup>Gaussian Graphical Model

<sup>4</sup>Sparse precision matrix

<sup>5</sup> $\ell_1$ -regularized maximum likelihood

<sup>6</sup>Neighborhood selection

<sup>7</sup>Graphical lasso

<sup>8</sup>Gaussian semiparametric copula

<sup>9</sup>Nonparanormal distribution

<sup>10</sup>Nonparanormal Graphical Model

در کل جامعه وجود داشته باشد. به عبارت دقیق‌تر، ممکن است مشاهده‌ها از خوشه‌ها یا مولفه‌های متفاوت آمده باشند؛ بدون آن‌که اطلاعاتی در مورد عضویت‌شان در خوشه‌ها در دست باشد. بنابراین، کشف زیرجامعه‌های متناهی با وجه اشتراک‌های خاص و برآورد وابستگی‌ها یا استقلال‌های شرطی در این زیرجامعه‌ها، یک چالش اساسی است. خوشه‌بندی مبتنی بر مدل<sup>۱</sup> با استفاده از مدل‌های آمیخته، ابزاری قدرتمند برای مدل‌بندی چنین داده‌هایی فراهم آورده است (لیندسی، ۱۹۹۵؛ مک لاکلان و پیل، ۲۰۰۴). ادبیات غنی پیرامون الگوریتم‌های خوشه‌بندی از جمله K- میانگین<sup>۲</sup>، خوشه‌بندی سلسله‌مراتبی<sup>۳</sup> و خوشه‌بندی مبتنی بر مدل وجود دارد (مک کوئین، ۱۹۶۷).

در روش‌های پیشین، به توزیع آماری داده‌ها توجه نشده است. به عنوان مثال، الگوریتم K- میانگین بر اساس میانگین فاصله اشیاء تا میانگین هر خوشه عمل می‌کند و اشیاء به‌گونه‌ای به خوشه‌ها تخصیص داده می‌شود که میانگین مجموع مربعات فاصله‌ها در خوشه‌ها، کم‌ترین مقدار را داشته باشد. در مقابل، در خوشه‌بندی مبتنی بر مدل، یک مدل آمیخته برای داده‌ها فرض می‌شود و با استفاده از آن، به برآورد برجسب اشیاء که معرف خوشه‌ی آن‌ها است، پرداخته می‌شود. به عبارت دقیق‌تر، خوشه‌بندی مبتنی بر مدل، یک الگوریتم خوشه‌بندی احتمالی است که با آن، میزان شدت باور شخص مبنی بر این که یک نقطه داده شده به یک خوشه خاص تعلق دارد، با تخصیص احتمال‌های پسین به آن نقطه معین ابراز می‌شود.

مدل‌های آمیخته گرافی یک رویکرد خوشه‌بندی مبتنی بر شبکه است که خوشه‌بندی مبتنی بر مدل و مدل‌های گرافی را گرد هم می‌آورد. این مدل‌ها امکان بازسازی خوشه‌ها و برآورد زیرشبکه‌های خوشه‌ای را به‌طور هم‌زمان فراهم آورده است. مدل آمیخته گرافی گاوسی<sup>۴</sup> (GGMM) توسط لاتسی و ویت (۲۰۱۶) معرفی شد که روش‌شناسی و کاربرد مدل GGMM را تعمیم داد. برای برآورد نسبت‌های آمیخته و ماتریس‌های دقت مدل GGMM، نسخه جدیدی از الگوریتم امید ریاضی-بیشینه‌سازی<sup>۵</sup> (EM) طراحی شد. اگرچه GGMM نسبت به GGMM برای مدل‌بندی داده‌های ناهمگن انعطاف‌پذیرتر و مناسب‌تر است، اما هنوز فرض محدودکننده نرمال را برای مولفه‌های مدل در نظر می‌گیرد. برای تقلیل فرض نرمال بودن مولفه‌های مدل، مدل آمیخته گرافی نرمال ناپارامتری<sup>۶</sup> (NGMM) با استفاده از توزیع نرمال ناپارامتری، توسط خلیلی و همکاران (۲۰۲۱) پیشنهاد شد. برآورد ساختار مدل NGMM، به برآورد نسبت‌های آمیخته، ماتریس‌های دقت و تبدیل‌های حاشیه‌ای تک متغیره بستگی دارد.

هدف این مطالعه، بازسازی خوشه‌ها از طریق یک رویکرد خوشه‌بندی مبتنی بر مدل است که مولفه‌های آن توسط مدل گرافی نرمال ناپارامتری تعریف شده است. بدین منظور، یک لگاریتم درست‌نمایی تاوانیده در نظر گرفته شده است که عبارت تاوان آن به نسبت‌های آمیخته وابسته است. الگوریتم EM مورد استفاده، امکان بازیابی تخصیص خوشه‌ها و برآورد پارامترهای مدل را می‌دهد. در سرتاسر این مطالعه، خوشه‌بندی مبتنی بر مدل توسط

<sup>1</sup>Model-based clustering

<sup>2</sup>K-means

<sup>3</sup>Hierarchical clustering

<sup>4</sup>Gaussian Graphical Mixture Model

<sup>5</sup>Expectation-maximization algorithm

<sup>6</sup>Nonparanormal Graphical Mixture Model

مک لاکلان و پیل (۲۰۰۴) در نظر گرفته شده و فرض بر آن است که داده‌ها از یک مدل متناهی آمیخته نرمال ناپارامتری استخراج شده‌اند که تعداد مولفه‌های آن  $(K)$ ، معلوم است.

در بخش ۲ پیش زمینه مورد نیاز پیرامون مدل گرافی نرمال ناپارامتری از لیو و همکاران (۲۰۰۹) و مدل آمیخته گرافی نرمال ناپارامتری از خلیلی و همکاران (۲۰۲۱) مرور شده است. سپس، در بخش ۳ رویکرد برآورد درست‌نمایی  $\ell_1$ -تاوانیده با دو فرم تابع تاوان (متعارف و نامتعارف) از طریق یک الگوریتم EM مبتنی بر مفصل ارائه شده است. بخش ۴ شامل یک مقایسه تجربی بین چندین روش خوشه‌بندی مبتنی بر مدل روی داده‌های شبیه‌سازی شده با دو نوع سازوکار آلودگی داده‌ها (قطعی و تصادفی) است. در بخش ۵، روش‌های خوشه‌بندی مورد بحث روی یک مجموعه داده واقعی سرطان سینه مقایسه شده‌اند. بخش پایانی به نتیجه‌گیری و بحث در مورد یافته‌ها می‌پردازد.

## ۲ پیش زمینه

یک گراف بی‌سو  $G = (V, E)$  شامل مجموعه‌ای متناهی از راس‌های  $V = \{1, \dots, p\}$  و مجموعه‌ای از یال‌های  $E \subset V \times V$  (که برخی از جفت راس‌های  $(i, j)$  را به هم متصل می‌کند) است. مدل‌های گرافی بی‌سو، توزیع توام  $n$  مشاهده مستقل از یک بردار تصادفی  $p$ -بعدی  $\mathbf{X} = (X_1, \dots, X_p)$  متناظر با راس‌های گراف  $G$  را به تصویر می‌کشد. زمانی که توزیع توام  $\mathbf{X}$ ، توزیع گاوسی، نرمال ناپارامتری، آمیخته گاوسی و آمیخته نرمال ناپارامتری فرض شده باشد، ویژگی‌های مدل گرافی مورد نظر برآورد شده است. توزیع نرمال ناپارامتری و رویکرد برآوردیابی جای‌گذاری<sup>۱</sup> برای برآورد مدل گرافی نرمال ناپارامتری شرح داده شده است. همچنین به مرور الگوریتم لاسوی گرافی که برای برآورد ماتریس دقت به‌کار می‌رود، پرداخته شده است. در نهایت، مدل آمیخته گرافی نرمال ناپارامتری معرفی شده است.

### ۲.۱ مدل نرمال ناپارامتری (NPN)

فرض کنید  $\{f_j\}_{j=1}^p$  مجموعه‌ای از تابع‌های یکنوا و مشتق‌پذیر تک متغیره نامعلوم است که به‌عنوان تبدیل‌های حاشیه‌ای در نظر گرفته شده‌اند. همچنین فرض کنید  $\Sigma \in \mathbb{R}^{p \times p}$  یک ماتریس همبستگی همیشه‌مثبت است که بدون از دست دادن کلیت مساله  $\mathbf{1} = \text{diag}(\Sigma)$  و  $\boldsymbol{\mu} = \mathbf{0}$ .

تعریف ۰.۱ (لیو و همکاران، ۲۰۰۹) بردار تصادفی  $p$ -بعدی  $\mathbf{X} = (X_1, \dots, X_p)$  دارای توزیع  $p$ -بعدی نرمال ناپارامتری<sup>۲</sup> است و آن را با نماد  $\mathbf{X} \sim \text{NPN}(\mathbf{f}, \Sigma)$  نمایش می‌دهند، هرگاه مجموعه‌ای از تبدیلات یکنوا و مشتق‌پذیر تک متغیره  $\{f_j\}_{j=1}^p$  موجود باشد، به‌نحوی که

$$\mathbf{f}(\mathbf{X}) = (f_1(X_1), \dots, f_p(X_p)) := (Z_1, \dots, Z_p) = \mathbf{Z} \sim N(\mathbf{0}, \Sigma). \quad (1)$$

<sup>1</sup>Plug-in estimation approach

<sup>2</sup>Nonparanormal

نرمال بودن حاشیه‌ای همیشه با تبدیلات به دست می‌آید، بنابراین مدل (۱) فرض می‌کند که توزیع توأم آن متغیرهای تبدیل‌شده، نرمال است. از آنجایی که توزیع‌های حاشیه‌ای تک متغیره نرمال هستند، مولفه‌های ماتریس همبستگی  $\Sigma$ ، دقیقاً همبستگی زوجی بین متغیرها را نشان می‌دهند. اگر  $f_j$ ها توابع مشتق‌پذیر باشند، تابع چگالی احتمال توأم  $\mathbf{X}$  به فرم

$$P_{\mathbf{X}}(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{f}(x))^T \Sigma^{-1} (\mathbf{f}(x)) \right\} \prod_{j=1}^p |f'_j(x_j)|. \quad (2)$$

خواهد بود که شناساپذیر نیست. برای این‌که این خانواده شناساپذیر باشد، فرض شده است که

$$\mu_j = E(Z_j) = E(X_j) = 0, \quad \sigma_j^2 = Var(Z_j) = Var(X_j) = 1, \quad j = 1, \dots, p.$$

با فرض این‌که  $\mathbf{Z} = (Z_1, \dots, Z_p) = (f_1(X_1), \dots, f_p(X_p))$  از یک توزیع نرمال پیروی می‌کند، واضح است که  $\omega_{ij} = 0$  اگر و تنها اگر  $Z_i \perp Z_j | \mathbf{Z}_{\setminus \{i,j\}}$ ، که  $\omega_{ij}$  نماد درایه  $(i, j)$  از ماتریس دقت  $\Omega = \Sigma^{-1}$  است. این ایده توسط لیو و همکاران (۲۰۰۹) به صورت

$$\omega_{ij} = 0 \iff Z_i \perp Z_j | \mathbf{Z}_{\setminus \{i,j\}} \iff X_i \perp X_j | \mathbf{X}_{\setminus \{i,j\}} \iff (i, j) \notin E, \quad (3)$$

کامل شد، که در آن  $\mathbf{X}_{\setminus \{i,j\}} := \{X_k : k \neq i, j\}$ . از این رو، برای برآورد گراف نرمال ناپارامتری، کافی است تابع‌های تک متغیره نامعلوم  $\{f_j\}_{j=1}^p$  برآورد شوند و سپس، مکان درایه‌های غیرصفر ماتریس  $\Omega$  مشخص شوند. در نتیجه، روند طبیعی برآورد مدل گرافی نرمال ناپارامتری، یک رویکرد برآوردیابی جای‌گذاری است که شامل دو گام زیر است:

- ۱- برای برآورد تبدیل  $f_j$  روی متغیر  $X_j$ ، از یک برآوردگر طبیعی به صورت  $\tilde{f}_j$  استفاده می‌شود.
- ۲- از یک تعمیم روش برآوردیابی مدل گرافی گاوسی (نظیر لاسوی گرافی) روی داده‌های تبدیل‌شده، برای برآورد ماتریس دقت استفاده می‌شود.

حال فرض کنید  $\mathbf{X}^1, \dots, \mathbf{X}^n$  یک نمونه تصادفی از  $NPN(\mathbf{0}, \Sigma)$  باشد که  $\mathbf{X}^i = (X_1^i, \dots, X_p^i)$ . همچنین  $\tilde{\mathbf{f}}(\mathbf{X}^1), \dots, \tilde{\mathbf{f}}(\mathbf{X}^n)$  نمونه تبدیل‌شده باشد که  $\tilde{\mathbf{f}}(\mathbf{X}^i) = (\tilde{f}_1(X_1^i), \dots, \tilde{f}_p(X_p^i))$ . ماتریس همبستگی نمونه‌ای  $R$  برای متغیرهای تصادفی تبدیل‌شده، به‌عنوان برآورد ماتریس همبستگی نامعلوم  $\Sigma$ ، به صورت

$$\hat{\Sigma} = R = \sum_{i=1}^n (\tilde{\mathbf{f}}(\mathbf{X}^i)) (\tilde{\mathbf{f}}(\mathbf{X}^i))^T,$$

۷۰ ..... خوشه‌بندی مبتنی بر مدل‌های آمیخته گرافی نرمال ناپارامتری

$$r_{jk} = \frac{\sum_{i=1}^n \tilde{f}_j(x_j^i) \tilde{f}_k(x_k^i)}{\sqrt{\sum_{i=1}^n \tilde{f}_j^2(x_j^i)} \sqrt{\sum_{i=1}^n \tilde{f}_k^2(x_k^i)}} \text{ با } R \text{ برابرند}$$

تعریف می‌شود، که در آن مولفه‌های ماتریس  $R$  برابرند با  $\hat{\Omega} = R^{-1} = \Sigma^{-1}$ . زمانی که تعداد متغیرها بسیار بیشتر از تعداد مشاهدات باشد ( $n \ll p$ )،  $R$  یک ماتریس تکین است و نمی‌توان از آن برای برآورد  $\Omega$  استفاده کرد. در این حالت، برآوردهای تنگ را می‌توان با اعمال یک تاوان  $\ell_1$  روی مولفه‌های  $\Omega$  به دست آورد. لگاریتم درست‌نمایی تاوانیده به صورت

$$\ell_p(\Omega) = \log |\Omega| - \text{tr}(R\Omega) - \lambda \|\Omega\|_1, \quad (4)$$

است، که در آن  $\lambda$  یک پارامتر تنظیم‌کننده<sup>۱</sup> برای کنترل سطح تنگی است و  $\|\Omega\|_1 = \sum_{i \neq j} \omega_{ij}$  یکی از روش‌های یافتن برآورد ماکسیمم درست‌نمایی تاوانیده<sup>۲</sup> (PMLE) برای ماتریس دقت  $\Omega$  در (۴)، استفاده از نسخه گرافی الگوریتم لاسو<sup>۳</sup> (glasso) است که توسط فریدمن و همکاران (۲۰۰۸) به فرم

$$\hat{\Omega}_{glasso}^\lambda = \arg \max_{\Omega > 0} \left\{ \log |\Omega| - \text{tr}(R\Omega) - \lambda \|\Omega\|_1 \right\}, \quad (5)$$

پیشنهاد شد، که در آن  $\Omega > 0$  بیان‌گر این است که  $\Omega$  یک ماتریس همیشه مثبت است.

## ۲.۲ مدل آمیخته گرافی نرمال ناپارامتری (NGMM)

فرض کنید  $X^1, \dots, X^n$  یک نمونه داده شده به اندازه  $n$  از  $K$  مولفه آمیخته باشد؛ به نحوی که  $X^i = (X_1^i, \dots, X_p^i) \in \mathbb{R}^p$ ،  $i = 1, \dots, n$ . همچنین فرض کنید که هر  $X^i$  از یکی از  $K$  مولفه با انتخاب تصادفی  $Y_i$  از مجموعه  $\{1, \dots, K\}$  با احتمال  $\pi_k$  استخراج شده باشد. نماد نسبت‌های آمیخته است که  $0 < \pi_k < 1$  و مجموع آن‌ها در شرط  $\sum_k \pi_k = 1$  صدق می‌کند. بنابراین  $Y_i$  یک متغیر پنهان چندجمله‌ای با پارامترهای  $\pi_k$  است؛ به طوری که  $P(Y_i = k) = \pi_k$  برای  $k = 1, \dots, K$ . هدف، مدل‌بندی داده‌های داده شده با تعیین توزیع توام  $P(X^i, Y_i) = P(X^i | Y_i) P(Y_i)$  است. هر خوشه، به طور جداگانه با در نظر گرفتن یک توزیع نرمال ناپارامتری مدل‌بندی می‌شود، به طوری که  $X^i | Y_i = k \sim NPN(\mathbf{f}_k, \Sigma_k)$  تابع چگالی  $X^i$  در خوشه  $k$  ام باشد.

تعریف ۲. (خلیلی و همکاران، ۲۰۲۱) بردار تصادفی  $p$ -بعدی  $X^i = (X_1^i, \dots, X_p^i)$  در مدل آمیخته نرمال ناپارامتری صدق می‌کند و آن را با نماد  $X^i \sim MNP N\{(\pi_k, \mathbf{f}_k, \Sigma_k)\}_{k=1}^K$  نمایش می‌دهند، هرگاه تابع

<sup>1</sup>Regularization tuning parameter

<sup>2</sup>Penalized Maximum likelihood estimate

<sup>3</sup>Graphical lasso



چگالی توام  $\mathbf{X}^i = (X_1^i, \dots, X_p^i)$  به فرم

$$f(\mathbf{x}^i | \Psi) = \sum_{k=1}^K \pi_k NPN(\mathbf{x}^i | \xi_k) = \sum_{k=1}^K \pi_k NPN(\mathbf{x}^i | \mathbf{f}_k, \Sigma_k), \quad (6)$$

باشد، که در آن  $NPN(\cdot | \xi_k)$  نماد چگالی نرمال ناپارامتری به فرم (۲) است،  $\Psi = (\Pi, \Xi)$  بردار همه پارامترهای نامعلوم است و  $\xi_k = (\mathbf{f}_k, \Sigma_k)$  پارامترهای مولفه  $k$ ام مدل برای  $k = 1, \dots, K$  است.

در مدل (۶)، فضای پارامتر را می‌توان به فرم

$$\Psi = \Pi \times \Xi, \quad (7)$$

نوشت، که در آن  $\Pi = \left\{ \{\pi_k\}_{k=1}^K \mid 0 < \pi_k < 1, \sum_k \pi_k = 1, k = 1, \dots, K \right\}$  و  $\Xi = \left\{ \{\mathbf{f}_k\}_{k=1}^K, \{\Omega_k\}_{k=1}^K \mid \Omega_k > 0, k = 1, \dots, K \right\}$

مدل آمیخته نرمال ناپارامتری که در تعریف ۲ ارائه شده است، در صورتی شناساپذیر است که برای هر دو پارامتر  $\psi, \psi^* \in \Psi$  و برای همه مقادیرهای ممکن  $\mathbf{x}$ ، رابطه زیر برقرار باشد

$$\sum_{k=1}^K \pi_k NPN(\mathbf{x} | \xi_k) = \sum_{k=1}^{K^*} \pi_k^* NPN(\mathbf{x} | \xi_k^*) \iff$$

$$K = K^*, \quad \pi_k = \pi_k^*, \quad NPN(\mathbf{x} | \xi_k) = NPN(\mathbf{x} | \xi_k^*).$$

تساوی  $NPN(\mathbf{x} | \xi_k) = NPN(\mathbf{x} | \xi_k^*)$  زمانی برای همه مقادیرهای  $\mathbf{x}$  برقرار است که  $f_{j(k)}$ ها میانگین و واریانس را حفظ کنند؛ یعنی

$$\mu_{j(k)} = E(Z_{j(k)}) = E(X_j) = 0, \quad \sigma_{j(k)}^2 = Var(Z_{j(k)}) = Var(X_j) = 1 \quad (8)$$

به‌ازای هر  $k \in \{1, \dots, K\}$  و  $j \in \{1, \dots, p\}$ .

### ۳ برآورد و خوشه‌بندی همزمان برای مدل آمیخته گرافی

در این بخش، ابتدا بردارهای تبدیل برآورد شده‌اند. سپس برآورد PMLE و به‌طور همزمان برآورد خوشه‌ها با استفاده از یک الگوریتم EM بر مبنای مفصل ارائه شده است.

### ۳.۱ برآورد بردارهای تبدیل حاشیه‌ای

فرض کنید  $\mathbf{X}^1, \dots, \mathbf{X}^n$  نمونه‌ای به اندازه  $n$  از آمیخته متناهی از مدل نرمال ناپارامتری باشد؛ یعنی  $\mathbf{X}^i = (X_1^i, \dots, X_p^i)^{i.i.d.} \sim MNP N\{(\pi_k, \mathbf{f}_k, \Sigma_k)\}_{k=1}^K$ . همچنین فرض کنید  $F_j(x)$  نماد تابع توزیع حاشیه‌ای  $X_j$ ،  $j = 1, \dots, p$ ، باشد. بردار تبدیل متناظر با مولفه  $k$ ام را به صورت

$$\mathbf{f}_k(\mathbf{X}) = (f_{1(k)}(X_1), \dots, f_{p(k)}(X_p)) := (Z_{1(k)}, \dots, Z_{p(k)}) = \mathbf{Z}_k \sim N(\mathbf{0}, \Sigma_k)$$

در نظر بگیرید که مولفه‌های قطری ماتریس همبستگی برابر با یک هستند؛  $\sigma_{jj(k)} = 1$ ،  $j = 1, \dots, p$ ، و  $k = 1, \dots, K$  آن‌گاه

$$F_j(x) = \sum_{k=1}^K \pi_k P(X_j \leq x) = \sum_{k=1}^K \pi_k P(f_{j(k)}(X_j) \leq f_{j(k)}(x)) = \sum_{k=1}^K \pi_k \Phi(f_{j(k)}(x)),$$

که در آن  $\Phi$  نماد تابع توزیع تجمعی نرمال استاندارد تک متغیره است. تحت شرایط شناساپذیری در (۸) و با فرض این‌که بردارهای تبدیلات برای همه مولفه‌ها یکسان‌اند  $f_{j(k)} = f_{j(k')}$ ،  $k, k' \in \{1, \dots, K\}$ ، که می‌توان از اندیس  $k$  صرف‌نظر کرد، تبدیل حاشیه‌ای متغیر  $j$ ام در همه مولفه‌ها برابر  $f_j(x) = \Phi^{-1}(F_j(x))$  است و برآوردگر تبدیل حاشیه‌ای  $f_j$  به صورت

$$\tilde{f}_j(x) = \Phi^{-1}(\tilde{F}_j(x)), \quad (9)$$

تعریف می‌شود، که در آن  $\tilde{F}_j$  یک برآوردگر  $F_j$  است. مشابه لیو و همکاران (۲۰۰۹)، برای اجتناب از بزرگی واریانس برآوردگرها در ابعاد بالا، از برآوردگر  $(n+1)$   $n\hat{F}_j(x) / (n+1)$  استفاده می‌کنیم، که در آن  $\hat{F}_j(t) = \frac{1}{n} \sum_{i=1}^n I(X_j^i \leq t)$ ، تابع توزیع تجربی حاشیه‌ای  $X_j$  است و  $I(\cdot)$  تابع نشان‌گر و  $t$  یک مقدار ثابت است.

### ۳.۲ درست‌نمایی تاوانیده مبتنی بر مدل

فرض کنید  $\mathbf{X}^1, \dots, \mathbf{X}^n$  نمونه‌ای تصادفی از  $MNP N\{(\pi_k, \mathbf{f}_k, \Sigma_k)\}_{k=1}^K$  باشد و همچنین فرض کنید داده‌های تبدیل شده متناظر با این نمونه با استفاده از (۹) باشد. طبق رویکرد برآوردیابی جای‌گذاری، برای پیاده‌سازی PMLE، مشاهدات با داده‌های تبدیل شده جای‌گزین می‌شوند. بنابراین مساله برآورد مدل NGMM، به مساله برآورد ماتریس‌های دقت  $\Omega_k = \Sigma_k^{-1}$  و نسبت‌های آمیخته  $\pi_k$  برای  $k = 1, \dots, K$

تقلیل می‌یابد و به دنبال آن، فضای پارامتری (۷) با فضای پارامتری  $\Theta = \Pi \times \Omega$  جای‌گزین می‌شود، که در آن

$$\begin{aligned} \Pi &= \left\{ \{\pi_k\}_{k=1}^K \mid 0 < \pi_k < 1, \sum_k \pi_k = 1, k = 1, \dots, K \right\} \\ \Omega &= \left\{ \{\Omega_k\}_{k=1}^K \mid \Omega_k > 0, k = 1, \dots, K \right\}. \end{aligned}$$

حال می‌توان تابع درست‌نمایی داده‌های ناتمام را به صورت

$$L(\Theta) = \prod_{i=1}^n \left( \sum_{k=1}^K \pi_k NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right),$$

نوشت، که در آن  $\Theta = \{(\pi_k, \Omega_k) : k = 1, \dots, K\}$  مجموعه تمام پارامترهای نامعلوم است.

تابع لگاریتم درست‌نمایی برابر است با

$$\ell(\Theta) = \sum_{i=1}^n \log \left( \sum_{k=1}^K \pi_k NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right). \quad (10)$$

برای افزایش سطح تنگی، یک تکنیک ماکسیمم درست‌نمایی تاوانیده با اعمال یک تاوان  $\ell_1$  روی هریک از  $K$  ماتریس دقت توسط ژو و همکاران (۲۰۰۹) و موخرجی و هیل (۲۰۱۱)، برای خوشه‌بندی با استفاده از مدل گرافی گاوسی ارائه شده است. همچنین یک تابع تاوان مشابه توسط خلیلی و چن (۲۰۰۷) و استدلر و همکاران (۲۰۱۰) برای آمیخته‌ی متناهی از مدل‌های رگرسیونی تاوانیده ارائه شده است. با اعمال این تاوان بر لگاریتم درست‌نمایی (۱۰)، لگاریتم درست‌نمایی تاوانیده به صورت

$$\ell_p(\Theta) = \sum_{i=1}^n \log \left( \sum_{k=1}^K \pi_k NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right) - \frac{n}{\gamma} \lambda \sum_{k=1}^K \pi_k^\gamma \|\Omega_k\|_1, \quad (11)$$

تعریف می‌شود، که در آن  $\lambda$  یک پارامتر تنظیم‌کننده است،  $\gamma$  یک پارامتر دودویی با مقادیر ۰ و ۱ است که فرم عبارت تاوان را کنترل می‌کند،  $\|\cdot\|_1$  نرم  $\ell_1$  است و  $\|\Omega_k\|_1 = \sum_{i \neq j} |w_{ij(k)}|$  وقتی  $\gamma = 0$  باشد، عبارت تاوان به نسبت‌های آمیخته  $\pi_k$  وابسته نیست و بنابراین، لگاریتم درست‌نمایی (۱۱)، به فرم متعارف لگاریتم درست‌نمایی تاوانیده که توسط لاتسی و ویت (۲۰۱۶) برای مدل آمیخته گرافی گاوسی و همچنین توسط خلیلی و همکاران (۲۰۲۱) برای مدل آمیخته گرافی نرمال ناپارامتری مورد استفاده قرار گرفته، تبدیل می‌شود. اما وقتی  $\gamma = 1$  باشد، لگاریتم درست‌نمایی (۱۱)، به عبارت تاوان از هر  $K$  خوشه به نسبت  $\pi_k$  وزن می‌دهد. بنابراین، تاوان  $\ell_1$  (۱۱) فرم کلی‌تر دارد؛ زیرا می‌تواند به نسبت‌های آمیخته بستگی داشته باشد یا نداشته باشد. در این مطالعه، این دو فرم تابع تاوان،

۷۴ ..... خوشه‌بندی مبتنی بر مدل‌های آمیخته گرافی نرمال ناپارامتری

برای خوشه‌بندی و برآورد پارامترهای مدل آمیخته گرافی گاوسی و مدل آمیخته گرافی نرمال ناپارامتری به‌کار رفته است.

### ۳.۳ الگوریتم EM مبتنی بر مفصل

برای داده‌های مستقل  $C = \{(\mathbf{x}^i, y_i), i = 1, \dots, n\}$ ، لگاریتم درست‌نمایی تاوانیده به صورت

$$\ell_{p,c}(\Theta) = \sum_{i=1}^n \sum_{k=1}^K 1_{\{Y_i=k\}} \log \left( \pi_k NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right) - \frac{n}{\gamma} \lambda \sum_{k=1}^K \pi_k^\gamma \|\Omega_k\|_1,$$

است، که در آن اگر مشاهده  $\mathbf{x}^i$  به خوشه  $k$  متعلق باشد،  $Y_i = k$  است. بنابراین،  $1_{\{Y_i=k\}}$  نشان می‌دهد که اگر مشاهده  $\mathbf{x}^i$  به خوشه  $k$  متعلق باشد، از  $\Theta_k$  استفاده می‌شود. گام E از الگوریتم EM، احتمال‌های شرطی این‌که  $\mathbf{X}^i$  از خوشه  $k$  آمده باشد، به شرط آن‌که برآوردهای جاری پارامترهای  $\Theta^{(t)}$  و داده‌های تبدیل شده، داده شده باشند، محاسبه می‌شود. ابتدا کمیت

$$\begin{aligned} Q(\Theta \mid \Theta^{(t)}) &= E_{Y_i} \left[ \ell_{p,c}(\Theta) \mid \{\tilde{\mathbf{f}}(\mathbf{x}^i)\}_{i=1}^n, \Theta^{(t)} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K E_{Y_i} \left[ 1_{\{Y_i=k\}} \tilde{\mathbf{f}}(\mathbf{x}^i), \Theta^{(t)} \right] \\ &\quad \times \left[ \log \pi_k + \log \left( NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right) \right] - \frac{n}{\gamma} \lambda \sum_{k=1}^K \pi_k^\gamma \|\Omega_k\|_1 \\ &= \sum_{i=1}^n \sum_{k=1}^K \tau_k^{i(t)} \left[ \log \pi_k + \log \left( NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1}) \right) \right] - \frac{n}{\gamma} \lambda \sum_{k=1}^K \pi_k^\gamma \|\Omega_k\|_1, \end{aligned} \quad (12)$$

محاسبه می‌شود، که در آن  $\tau_k^{i(t)}$  احتمال پسین تعلق مشاهده  $\mathbf{x}^i$  به خوشه  $k$  است. طبق قضیه بیز

$$\tau_k^{i(t)} = P(Y_i = k \mid \tilde{\mathbf{f}}(\mathbf{x}^i), \Theta^{(t)}) = \frac{\pi_k^{(t)} NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_k^{-1(t)})}{\sum_{j=1}^K \pi_j^{(t)} NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) \mid \Omega_j^{-1(t)})}. \quad (13)$$

برای اجرای خوشه‌بندی مبتنی بر مدل، هر مشاهده  $\mathbf{x}^i$  را به خوشه  $k$  تخصیص می‌دهند، هرگاه

$$k' = 1, \dots, K \text{ و } \forall k' \neq k, \tau_k^{i(t)} > \tau_{k'}^{i(t)}$$

گام M. برای به‌دست آوردن برآوردهای به‌روز رسانی شده  $\Theta^{(t+1)}$ ، کمیت  $Q(\Theta \mid \Theta^{(t)})$  نسبت به  $(\pi_k, \Omega_k)$  است.

$k = 1, \dots, K$  ماکسیم می‌شود. بنابراین، گام  $M$  به دو مساله بهینه‌سازی تقسیم می‌شود؛ ماکسیم‌سازی نسبت به  $\pi_k$  و ماکسیم‌سازی نسبت به  $\Omega_k$ .

۱- گام  $M$  برای  $\pi_k$ : وقتی  $\gamma = 0$ ، نسبت‌های آمیخته  $\pi_k$  در عبارت تاوان ظاهر نمی‌شوند، لذا قید  $\pi_K = 1 - \sum_{k=1}^{K-1} \pi_k$  و  $\pi_k > 0$  به کار برده می‌شود. در این حالت، با ماکسیم‌سازی (۱۲) نسبت به  $\pi_k$ ، به‌روز رسانی استاندارد EM مشابه خلیلی و همکاران (۲۰۲۱) به‌دست می‌آید

$$\pi_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_k^{i(t)}}{n}. \quad (14)$$

به ازای  $\gamma = 1$ ، نسبت‌های آمیخته  $\pi_k$  در عبارت تاوان ظاهر می‌شوند. در این حالت، از به‌روز رسانی استاندارد (۱۴) استفاده می‌کنیم. اگر این به‌روز رسانی استاندارد،  $Q(\Theta | \Theta^{(t)})$  را بهبود بخشد، آن‌گاه برای به‌دست آوردن ماکسیم (۱۱) کافی است. اگرچه بهبود در این مورد تضمین نشده است، اما این روش در کاربرد به‌خوبی عمل می‌کند، هم‌چنان‌که در خلیلی و چن (۲۰۰۷) نشان داده شده است.

۲- گام  $M$  برای  $\Omega_k$ : برای ماکسیم کردن (۱۲) نسبت به ماتریس دقت  $\Omega_k$ ، عباراتی که به  $\Omega_k$  بستگی ندارند، حذف می‌شوند. بنابراین، کمیت  $Q(\Theta | \Theta^{(t)})$  به کمیت  $Q(\Omega_k)$  تقلیل می‌یابد

$$\begin{aligned} Q(\Omega_k) &= \sum_{i=1}^n \tau_k^{i(t)} \log(NPN(\tilde{\mathbf{f}}(\mathbf{x}^i) | \Omega_k^{-1})) - \frac{n}{\gamma} \lambda \pi_k^\gamma \|\Omega_k\|_1 \\ &= \sum_{i=1}^n \tau_k^{i(t)} \left[ \frac{1}{\gamma} \log |\Omega_k| - \frac{1}{\gamma} (\tilde{\mathbf{f}}(\mathbf{x}^i))^T \Omega_k (\tilde{\mathbf{f}}(\mathbf{x}^i)) \right] - \frac{n}{\gamma} \lambda \pi_k^\gamma \|\Omega_k\|_1 \\ &= \frac{1}{\gamma} \sum_{i=1}^n \tau_k^{i(t)} \left[ \log |\Omega_k| - tr(R_k \Omega_k) - \tilde{\lambda}_k \|\Omega_k\|_1 \right]. \end{aligned}$$

بنابراین، به‌روز رسانی  $\Omega_k$  به‌صورت

$$\Omega_k^{(t+1)} = \arg \max_{\Omega_k > 0} \left\{ \log |\Omega_k| - tr(R_k^{(t)} \Omega_k) - \tilde{\lambda}_k^{(t)} \|\Omega_k\|_1 \right\}, \quad (15)$$

به‌دست می‌آید، که در آن  $R_k^{(t)}$  به‌روز رسانی ماتریس همبستگی تجربی وزنی است که برابر است با

$$\begin{aligned} R_k^{(t)} &= \frac{\sum_{i=1}^n \tau_k^{i(t)} (\tilde{\mathbf{f}}(\mathbf{x}^i)) (\tilde{\mathbf{f}}(\mathbf{x}^i))^T}{\sum_{i=1}^n \tau_k^{i(t)}}, \\ \tilde{\lambda}_k^{(t)} &= n \lambda \frac{(\pi_k^{(t+1)})^\gamma}{\sum_{i=1}^n \tau_k^{i(t)}}. \end{aligned} \quad (16)$$

از (۱۴) و (۱۶) نتیجه می‌شود

$$\tilde{\lambda}_k^{(t)} = \begin{cases} \frac{\lambda}{\pi_k^{(t+1)}} & \gamma = 0 \\ \lambda & \gamma = 1 \end{cases} \quad (17)$$

رابطه (۱۵) به فرم لاسوی گرافی در (۵) است که در آن، ماتریس همبستگی نمونه  $R$  با یک ماتریس همبستگی تجربی وزنی  $R_k$  و پارامتر تنظیم‌کننده  $\lambda$  با یک پارامتر تنظیم‌کننده مقیاس  $\tilde{\lambda}_k$  جای‌گزین شده است. بنابراین، می‌توان از الگوریتم لاسوی گرافی برای حل این مساله بهینه‌سازی استفاده کرد. توجه داشته باشید که وقتی  $\gamma = 0$ ، فرم متعارف لگاریتم درست‌نمایی تاوانیده (۱۱) از یک پارامتر تنظیم‌کننده  $\lambda$  استفاده می‌کند، درحالی‌که به‌روز رسانی EM (۱۵) یک پارامتر تنظیم‌کننده متناسب با عکس نسبت آمیخته خوشه را ارائه می‌دهد. برای  $\gamma = 1$ ، عکس این مطلب برقرار است.

### ۳.۴ انتخاب پارامتر تنظیم‌کننده

میزان سطح تنکی در ماتریس دقت و مدل گرافی مربوطه، توسط پارامتر تنظیم‌کننده  $\lambda$  کنترل می‌شود. روش‌های مختلف انتخاب پارامتر تنظیم‌کننده همراه با فرم‌های متفاوت تاوان  $\ell_1$ ، منجر به رویکردهای برآوردیابی متفاوت می‌شود. فرم کلی معیار اطلاع‌بیزی<sup>۱</sup> (BIC) به‌صورت

$$BIC(\lambda) = -2\ell(\hat{\Theta}^\lambda) + \log(n)df^\lambda, \quad (18)$$

است، که در آن  $\ell(\cdot)$  لگاریتم درست‌نمایی غیرتاوانیده (۱۰)،  $\hat{\Theta}^\lambda$  برآورد درست‌نمایی تاوانیده با پارامتر تنظیم‌کننده  $\lambda$  و  $df$  درجه آزادی برای اندازه‌گیری پیچیدگی مدل است. برای تعیین درجه آزادی در مدل گرافی گاوسی، کمیت

$$df^\lambda = K(p+1) - 1 + \sum_{k=1}^K \# \{ (i, j) : i \leq j, \hat{\omega}_{ij}^\lambda \neq 0 \}, \quad (19)$$

توسط **یوان و لین (۲۰۰۷)** پیشنهاد شده است، که در آن  $\hat{\omega}_{ij}^\lambda$  درایه  $(i, j)$  در  $\hat{\Omega}_k^\lambda$  است.  $\lambda$  با مینیمم‌سازی نمره‌های BIC موجود برای هر  $\lambda$ ، انتخاب می‌شود؛ یعنی  $\hat{\lambda} = \arg \min_{\lambda} BIC(\lambda)$ . توجه کنید که مقدار بهینه‌ی پارامتر تنظیم‌کننده کمتر تحت تاثیر تخصیص خوشه‌ها است، ولی به شدت به ویژگی‌های کلی داده‌ها نظیر  $n$ ،  $p$  و  $K$  بستگی دارد.

<sup>1</sup>Bayesian Information Criterion

### ۳.۵ الگوریتم برآورد تخصیص خوشه‌ها و پارامترهای مدل

الگوریتم ۱. تخصیص خوشه‌ها و برآورد پارامترها

ورودی:  $n$  مشاهده مستقل  $\mathbf{X}^1, \dots, \mathbf{X}^n$  که  $\mathbf{X}^i = (X_1^i, \dots, X_p^i)$ ،  $i = 1, \dots, n$

گام ۱- جای‌گزینی مشاهده‌های هر متغیر  $X_j$ ،  $j = 1, \dots, p$ ، با مقدار استاندارد شده آن‌ها.

گام ۲- برآورد تبدیل‌های  $f_j$  از رابطه  $\tilde{f}_j(x) = \Phi^{-1}\left(\frac{1}{n+1} \sum_{i=1}^n I(X_j^i \leq x)\right)$ ،  $j = 1, \dots, p$

گام ۳- مقدار دهی اولیه  $\Theta^{(0)} = \{(\pi_k^{(0)}, \Omega_k^{(0)})\}_{k=1}^K$

۱.۳- تخصیص تصادفی هر مشاهده  $\mathbf{x}^i$  به یکی از  $K$  خوشه منوط به اندازه معین  $n_k^{(0)} = n/2$

۲.۳- تعیین  $\pi_k^{(0)} = n_k^{(0)}/n = 0.5$  برای  $k = 1, \dots, K$  که  $n_k^{(0)}$  تعداد مشاهده‌های تخصیص داده شده

اولیه به خوشه  $k$  است.

۳.۳- تعیین  $\Omega_k^{(0)}$ ،  $k = 1, \dots, K$ ، برای ماتریس دقت خوشه  $k$  با استفاده از لاسوی گرافی

$$\Omega_k^{(0)} = \arg \max_{\Omega_k > 0} \left\{ \log |\Omega_k| - \text{tr}(R_k^{(0)} \Omega_k) - \lambda \|\Omega_k\|_1 \right\},$$

که  $R_k^{(0)} = \sum_{i=1}^n (\tilde{f}_k(\mathbf{X}^i)) (\tilde{f}_k(\mathbf{X}^i))^T$ ، ماتریس همبستگی نمونه اولیه در خوشه  $k$  است.

گام ۴-  $E$ : محاسبه احتمال‌های پسین  $\tau_k^{i(t)} = \frac{\pi_k^{(t)} NPN(\mathbf{f}(\mathbf{x}^i) | \Omega_k^{-1(t)})}{\sum_{j=1}^K \pi_j^{(t)} NPN(\mathbf{f}(\mathbf{x}^i) | \Omega_j^{-1(t)})}$

گام ۵-  $M$ : محاسبه برآوردهای به‌روز رسانی شده پارامترها  $\Theta^{(t+1)} = \{(\pi_k^{(t+1)}, \Omega_k^{(t+1)})\}_{k=1}^K$

۱.۵- محاسبه  $\pi_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_k^{i(t)}}{n}$  برای  $k = 1, \dots, K$

۲.۵- تعیین به‌روز رسانی  $EM$  برای ماتریس دقت خوشه  $k$ ،  $k = 1, \dots, K$ ، از رابطه

$$\Omega_k^{(t+1)} = \arg \max_{\Omega_k > 0} \left\{ \log |\Omega_k| - \text{tr}(R_k^{(t)} \Omega_k) - \tilde{\lambda}_k^{(t)} \|\Omega_k\|_1 \right\},$$

که  $\tilde{\lambda}_k^{(t)} = \frac{\lambda}{\pi_k^{(t+1)}} I_{\{\gamma=0\}} + \lambda I_{\{\gamma=1\}}$  و  $R_k^{(t)} = \frac{\sum_{i=1}^n \tau_k^{i(t)} (\tilde{f}(\mathbf{x}^i)) (\tilde{f}(\mathbf{x}^i))^T}{\sum_{i=1}^n \tau_k^{i(t)}}$

گام ۶- تخصیص هر مشاهده  $\mathbf{x}^i$  به خوشه  $k$  با بزرگترین احتمال  $\tau_k^{i(t)}$ ،  $i = 1, \dots, n$

گام ۷- تکرار گام‌های ۴ تا ۶ تا توقف آن‌ها در صورت برآورده شدن یکی از شرایط زیر:

(i) ماکسیمم تعداد تکرارها ( $T$ ) حاصل شود:  $t > T$

(ii) برای خوشه  $k$ ، مینیمم اندازه خوشه ( $n_{\min}$ ) حاصل شود:  $\sum_{i=1}^n \tau_k^{i(t)} < n_{\min}$

(iii) تغییر نسبی در لگاریتم درست‌نمایی، ناچیز باشد:  $\left| \ell_p(\Theta)^{(t)} / \ell_p(\Theta)^{(t-1)} - 1 \right| \leq \epsilon$ . خروجی.

برآورد تخصیص خوشه‌ها، نسبت‌های آمیخته و ماتریس‌های دقت خوشه‌ای.

## ۴ مطالعه شبیه‌سازی

نتایج تجربی خوشه‌بندی مبتنی بر مدل آمیخته گرافی نرمال ناپارامتری و خوشه‌بندی مبتنی بر مدل آمیخته گرافی گاوسی، برای دو فرم کلی تابع تاوان  $\ell_1$  (متعارف و نامتعارف)، مقایسه شده است. بدین منظور، ابتدا فرایند تولید داده‌ها و روند مقایسه کمی شرح داده شده است. سپس مقدار پارامتر تنظیم‌کننده انتخاب شده برای هر روش، نمایش داده شده است. در ادامه، میزان تطبیق خوشه‌های بازسازی‌شده در هر روش با خوشه‌های اصلی و همچنین میزان تطبیق پارامترهای برآورد شده از خوشه‌های بازسازی‌شده در هر روش با پارامترهای واقعی از خوشه‌های اصلی، مقایسه شده است.

### ۴.۱ تولید داده

برای تولید داده از مدل آمیخته نرمال ناپارامتری، داده‌های  $p$ -بعدی شامل دو خوشه ( $K = 2$ ) با دو ماتریس دقت تنک معلوم در نظر گرفته شده است. بدین منظور، ابتدا دو گراف تنک تولید می‌شود و سپس، ماتریس‌های دقت  $\Omega_1$  و  $\Omega_2$  متناظر با این دو گراف به دست می‌آیند. در ادامه،  $n$  نمونه مستقل از  $N(\mathbf{0}, \Sigma_k)$  که  $\Sigma_k = \Omega_k^{-1}$  و  $k = 1, 2$ ، استخراج می‌شود. بعد از آن، با استفاده از تبدیل‌های  $\{f_{j(k)}\}_{j=1}^p$ ، داده‌های نرمال  $N(\mathbf{0}, \Sigma_k)$  به داده‌های نرمال ناپارامتری  $NPN(f_k, \Sigma_k)$  تبدیل می‌شوند. در نهایت، مقادیر متغیر برنولی  $Y_i \sim Ber(1, \delta)$ ،  $i = 1, \dots, n$ ، تولید و براساس مقادیر آن، داده‌های  $\mathbf{X}^1, \dots, \mathbf{X}^n$  از مدل آمیخته نرمال ناپارامتری دو مولفه‌ای با احتمال‌های برابر  $\pi_k = \delta$ ، تولید می‌شوند. جزئیات الگوهای تولید گراف در ادامه شرح داده شده است. الگوی **I** (بسته نرم‌افزاری Huge<sup>۱</sup>). الگوی گراف بدین ترتیب تولید می‌شود که احتمال آن که بین هر زوج راس‌های  $(i, j)$ ،  $(i \neq j)$ ، یالی وجود داشته باشد، برابر است با

$$P((i, j) \in E) = \sqrt{\delta} I(\delta_p < \delta) + (1 - \sqrt{\delta(1 - \delta_p)}) I(\delta_p \geq \delta),$$

که  $\delta_p = \min(1, 3/p)$  و  $p$  تعداد راس‌ها است. مجموعه یال  $E$  شامل هر زوج راس  $(i, j)$  می‌شود، اگر و تنها اگر، مقدار احتمال  $P((i, j) \in E)$ ، از احتمال متناظر با یک توزیع یکنواخت دو متغیره در مربع  $[0, 1]^2$  کمتر باشد. معکوس ماتریس همبستگی  $\Omega_1 = (\omega_{ij(1)})_{p \times p}$  برای خوشه ۱ به صورت زیر ساخته می‌شود

$$\omega_{ij(1)} = \begin{cases} |e| + \delta + u & i = j \\ \delta & (i, j) \in E \\ 0 & \text{سایر نقاط} \end{cases}$$

<sup>۱</sup><http://cran.r-project.org/web/packages/huge>.



که در آن،  $u$  عددی مثبت است که برای کنترل بزرگی همبستگی‌های جزئی، به عناصر قطری ماتریس دقت اضافه شده است (مقدار پیش‌فرض  $0/1$  است). برای تضمین همیشه مثبت بودن معکوس ماتریس همبستگی، عناصر قطری  $\Omega_1$ ، با مقدار قدر مطلق کوچکترین مقدار ویژه  $\Omega_1$  (که با نماد  $|e|$  نشان داده شده است)، جای‌گزین شده‌اند.

**الگوی II (میشاوسن و بولمن ، ۲۰۰۶).** الگوی گراف این‌گونه ایجاد می‌شود که مکان هر گره به‌طور یکسان و یکنواخت در یک مربع دوبعدی  $[0, 1]^2$  توزیع شده است. به عبارت دقیق‌تر، هر اندیس  $j \in \{1, \dots, p\}$  با یک نقطه داده دومتغیره  $(U_j^{(1)}, U_j^{(2)}) \in [0, 1]^2$  در تناظر است، به طوری که  $U_1^{(m)}, \dots, U_p^{(m)} \sim uniform[0, 1]$  برای  $m = 1, 2$ . هر جفت از راس‌های  $(i, j)$  ابتدا با احتمال

$$P((i, j) \in E) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\|u_i - u_j\|^2}{2s}\right),$$

در مجموعه راس  $E$  قرار می‌گیرد، که در آن  $u_i = (u_i^{(1)}, u_i^{(2)})$  مقدار مشاهده شده  $(U_i^{(1)}, U_i^{(2)})$  و  $\|\cdot\|$  نماد فاصله اقلیدسی بین زوج متغیرها است.  $s = 0/125$  پارامتری است که سطح تنکی گراف تولید شده را کنترل می‌کند. برای دستیابی به سطح تنکی مطلوب، ماکسیمم درجه گراف، یعنی حداکثر تعداد یال‌های متصل به هر گره، به  $4$  محدود شده است و یال‌های اضافی متصل به هر گره، به‌طور تصادفی حذف می‌شوند. بنابراین، معکوس ماتریس همبستگی  $\Omega_2 = (\omega_{ij(2)})_{p \times p}$  برای خوشه  $2$  به صورت

$$\omega_{ij(2)} = \begin{cases} 1 & i = j \\ 0/245 & (i, j) \in E \\ 0 & \text{سایر نقاط} \end{cases}$$

تشکیل می‌شود که در آن همبستگی جزئی بین راس‌های مجاور،  $0/245$  در نظر گرفته شده است. مقادیر مطلق کمتر از  $0/245$ ، همیشه مثبت بودن معکوس ماتریس همبستگی را تضمین می‌کند. برای به‌دست آوردن ماتریس همبستگی، ماتریس کواریانس  $\Sigma = \Omega^{-1}$  حاصل، به ماتریس همبستگی تبدیل می‌شود؛ به نحوی که همه عناصر قطری آن برابر با  $1$  شود.

داده‌های نرمال و داده‌های نرمال ناپارامتری از مدل‌های زیر تولید می‌شود:

- مدل (۱) - مدل آمیخته گاوسی چند متغیره با دو مولفه با استفاده از تبدیل خطی روی داده‌ها.
- مدل (۲) - مدل آمیخته نرمال ناپارامتری چند متغیره با دو مولفه با استفاده از تبدیل‌های توانی تک متغیره یکسان روی هر بعد، ارائه شده توسط لیو و همکاران (۲۰۰۹).
- مدل (۳) - مدل آمیخته نرمال ناپارامتری چند متغیره با دو مولفه با استفاده از تبدیل‌های تابع توزیع (CDF) تک متغیره یکسان روی هر بعد، ارائه شده توسط لیو و همکاران (۲۰۰۹).

در مدل‌های (۲) و (۳)، برای هر مولفه  $k = 1, 2$ ، ابتدا  $n$  داده مستقل از  $N(\mathbf{0}, \Sigma_k)$  که  $\Sigma_k = \Omega_k^{-1}$ ، تولید می‌شود و سپس، با استفاده از بردارهای تبدیل  $(f_{1(k)}^{-1}, \dots, f_{p(k)}^{-1})$ ،  $\mathbf{g}_k := \mathbf{f}_k^{-1}$ ، داده‌ها به  $NPN(\mathbf{f}_k, \Sigma_k)$  تبدیل می‌شود. همان‌طور که در بخش ۳.۱ فرض شد، بردارهای تبدیل روی هر مولفه، تابع‌های یکسان هستند؛ یعنی  $f_1 = f_2$  یا  $f_{j(1)} = f_{j(2)}$ . بنابراین از اندیس  $(k)$  برای تابع‌های تبدیل، صرف‌نظر شده است. همچنین در مدل‌های (۲) و (۳) تابع تبدیل‌های تک متغیره یکسان روی هر متغیر  $X_j$ ، برای  $j = 1, \dots, p$ ، در نظر گرفته شده است؛ یعنی  $f_1 = f_2 = \dots = f_p = f$ . در مدل‌های (۲) و (۳)، دو نسخه مختلف از  $f_j^{-1} := g_j$  به کار برده شده است که این دو تبدیل، در تعریف ۹ و تعریف ۱۰ از مقاله لیو و همکاران (۲۰۰۹) تعریف شده‌اند.

همچنین دو نوع سازوکار آلودگی داده‌ها در سطح  $r \in (0, 1)$  در نظر گرفته شد؛ قطعی و تصادفی. برای آلودگی قطعی (سطری)، تعداد  $[nr]$  مشاهده با بردار قطعی  $(+5, -5, +5, -5, \dots) \in \mathbb{R}^p$  جای‌گزین می‌شود، به‌نحوی که اعداد  $+5$  و  $-5$  به‌صورت متناوب تکرار شوند. برای آلودگی تصادفی (ستونی)،  $[nr]$  مولفه از هر بعد براساس یک توزیع یکنواخت، با اعداد  $+5$  یا  $-5$  به‌صورت تصادفی جای‌گزین می‌شوند. این دو نوع آلودگی، برای تجزیه و تحلیل داده‌های علمی مدرن، واقع‌بینانه‌تر هستند.

برای تولید داده، دو بُعد  $p = 15, 20$  و دو اندازه نمونه  $n = 30, 150$  در نظر گرفته شده که با آلودگی قطعی، ۵ درصد داده دورافتاده ( $r = 0/05$ ) و با آلودگی تصادفی، ۲۰ درصد داده دورافتاده ( $r = 0/20$ ) جای‌گزین شده است. همچنین، معیار توقف یکسان  $T = 100$ ،  $n_{min} = 4$  و  $\epsilon = 10^{-4}$ ، برای الگوریتم EM اعمال شده است. همه محاسبات با نرم‌افزار R انجام شده است.

به‌طور کلی، روش‌های زیر با هم مقایسه شده‌اند:

(i)  $GGMM$ : خوشه‌بندی مبتنی بر مدل GGMM با استفاده از فرم متعارف لگاریتم درست‌نمایی  $\ell_1$ -تاوانیده از لاتسی و ویت (۲۰۱۶).

(ii)  $GGMM_1$ : خوشه‌بندی مبتنی بر مدل GGMM با استفاده از فرم نامتعارف لگاریتم درست‌نمایی  $\ell_1$ -تاوانیده از موخرجی و هیل (۲۰۱۱).

(iii)  $NGMM$ : خوشه‌بندی مبتنی بر مدل NGMM با استفاده از فرم متعارف لگاریتم درست‌نمایی  $\ell_1$ -تاوانیده از خلیلی و همکاران (۲۰۲۱).

(iv)  $NGMM_1$ : خوشه‌بندی مبتنی بر مدل NGMM با استفاده از فرم نامتعارف لگاریتم درست‌نمایی  $\ell_1$ -تاوانیده پیشنهاد شده در این مقاله.

شایان ذکر است لاتسی و ویت (۲۰۱۶) و خلیلی و همکاران (۲۰۲۱) از فرم متعارف لگاریتم درست‌نمایی  $\ell_1$ -تاوانیده تنها برای برآورد پارامترهای GGMM و NGMM استفاده کردند.

## ۴.۲ مقادیر پارامتر تنظیم‌کننده

- بعد از تولید داده از مدل‌های (۱) تا (۳)، گام‌های زیر برای انتخاب پارامتر تنظیم‌کننده طی شده‌اند:
- ۱- تولید  $K = 2$  خوشه‌نما با تخصیص تصادفی مشاهدات به هر خوشه.
  - ۲- محاسبه مقادیر اولیه پارامتر برای هر دو خوشه‌نما:
  - $\hat{\pi}_k^{(e)}$  به‌عنوان نسبت اندازه‌های نمونه در خوشه‌نمای  $k$  تعیین می‌شود.
  - برای  $10$  مقدار متفاوت  $\lambda$ ،  $\hat{\Omega}_k^{(e)\lambda}$  به‌طور همزمان با استفاده از الگوریتم لاسوی گرافی (۵) در خوشه‌نمای  $k$  به‌دست می‌آید. این کار با استفاده از تابع `huge.glasso` در بسته نرم‌افزاری Huge انجام شده است.
  - ۳- با استفاده از این برآوردها، نمره‌های BIC محاسبه می‌شوند و برای انتخاب یک پارامتر  $\lambda$  مشترک بین هر دو خوشه، مینیمم می‌شوند.
  - ۴- گام‌های ۱ تا ۳،  $50$  بار تکرار می‌شوند و از مقدارهای  $\lambda$  به‌دست آمده، میانگین گرفته می‌شود تا یک مقدار نهایی حاصل شود.

میانگین این مقادیر برای هر روش و به ازای دو مقدار متفاوت برای  $n$  و  $p$ ، در جداول ۱ و ۲ منعکس شده است. طبق (۲۰)، برای روش‌های  $NGMM_1$  و  $GMM_1$  ( $\gamma = 1$ )، به‌روز رسانی‌های EM برای پارامترهای تنظیم‌کننده  $\lambda_k$  برابر با  $\lambda$  است که در جداول ۱ و ۲ نیز این برابری برقرار است. همان‌طور که از فرم BIC در (۱۸) و (۱۹) انتظار می‌رفت، مقدارهای  $\lambda$  به دست آمده با افزایش  $n$ ، کاهش می‌یابند. همچنین با افزایش سطح آلودگی داده‌ها ( $r$ )، انحراف استاندارد مقدارهای  $\lambda$  افزایش می‌یابد، اما میانگین مقدارهای  $\lambda$  در برابر افزایش آلودگی، تقریباً استوار است.

## ۴.۳ انتساب خوشه

روش‌های خوشه‌بندی دارای دو بخش معیار و تکنیک است. معیار، به هر خوشه یک مقدار عددی اختصاص می‌دهد که گویای مطلوبیت نسبی روش مورد نظر است. شاخص رند<sup>۱</sup> (RI)، یک معیار عینی برای سنجش میزان تشابه بین تخصیص‌های خوشه‌های واقعی و خوشه‌های برآورد شده است که در بازه (۰، ۱) مقدار می‌پذیرد.  $RI = 0$  بدین معنی است که الگوریتم خوشه‌بندی مورد نظر در برآورد ساختار جامعه شکست خورده است و در مقابل،  $RI = 1$  به معنی انطباق کامل خوشه‌های برآورد شده با خوشه‌های واقعی است (رند، ۱۹۷۱). شکل‌های ۱ و ۲ شاخص‌های رند به‌دست آمده در هر روش از خوشه‌بندی  $50$  مجموعه داده شبیه‌سازی شده با سه تبدیل داده (خطی، توانی و تابع توزیع) و دو سازوکار آلودگی (قطعی با ۵ درصد آلودگی و تصادفی با ۲۰ درصد آلودگی)، برای دو مقدار متفاوت  $n$  و  $p$  را نشان می‌دهد. تفسیر نتایج بدین شرح است:

الف- داده‌های نرمال و بدون داده دورافتاده ( $r = 0$ ). زمانی که تبدیل خطی روی داده‌های نرمال اعمال می‌شود، روش‌های مورد بررسی در برآورد خوشه‌های واقعی تفاوت معنی‌داری ندارند و هر چهار روش عملکرد مشابهی دارند.

<sup>1</sup>Rand index

جدول ۱. متوسط مقادیرهای پارامتر تنظیم‌کننده برای هر روش روی ۵۰ مجموعه داده شبیه‌سازی شده در بُعد  $p = ۱۵$  و اندازه نمونه  $n = ۳۰$ . انحراف استانداردها داخل پرانتزها نوشته شده‌اند.

$NGMM_1$		$NGMM_0$		$GGM_1$		$GGM_0$		$r$	تبدیل
$\lambda = \bar{\lambda}_k$	$\bar{\lambda}_r$	$\bar{\lambda}_1$	$\lambda$	$\lambda = \bar{\lambda}_k$	$\bar{\lambda}_r$	$\bar{\lambda}_1$	$\lambda$		
۰.۶۳ (۰.۰۵)	۱.۳۰ (۰.۱۰)	۱.۲۳ (۰.۱۶)	۰.۶۳ (۰.۰۵)	۰.۶۵ (۰.۰۴)	۱.۳ (۰.۰۸)	۱.۲۶ (۰.۱۳)	۰.۶۴ (۰.۰۴)	۰.۰۰	خطی
۰.۶۴ (۰.۰۵)	۱.۳۱ (۰.۱۰)	۱.۳۱ (۰.۱۱)	۰.۶۶ (۰.۰۵)	۰.۷۲ (۰.۰۸)	۱.۷۱ (۰.۳۹)	۱.۳۷ (۰.۴۸)	۰.۷۰ (۰.۰۹)	۰.۰۵	
۰.۶۰ (۰.۶۷)	۱.۷۴ (۰.۶۷)	۱.۱۰ (۰.۷۴)	۰.۶۰ (۰.۵۳)	۰.۶۲ (۰.۷۳)	۱.۳۲ (۰.۴۳)	۱.۴۵ (۰.۶۵)	۰.۶۲ (۰.۶۳)	۰.۲۰	
۰.۶۲ (۰.۰۶)	۱.۳۰ (۰.۱۰)	۱.۲۲ (۰.۱۷)	۰.۶۳ (۰.۰۵)	۰.۷۶ (۰.۰۵)	۱.۵۶ (۰.۰۹)	۱.۵۱ (۰.۱۸)	۰.۷۶ (۰.۰۶)	۰.۰۰	توانی
۰.۶۴ (۰.۰۵)	۱.۳۰ (۰.۱۰)	۱.۳۱ (۰.۱۱)	۰.۶۵ (۰.۰۵)	۰.۸۰ (۰.۰۶)	۱.۶۹ (۰.۲۱)	۱.۶۵ (۰.۳۲)	۰.۸۲ (۰.۰۶)	۰.۰۵	
۰.۶۰ (۰.۶۵)	۱.۳۱ (۰.۲۸)	۱.۳۴ (۰.۳۱)	۰.۵۹ (۰.۵۸)	۰.۶۲ (۰.۸۲)	۱.۲۸ (۰.۳۵)	۱.۳۸ (۰.۳۵)	۰.۶۳ (۰.۸۵)	۰.۲۰	
۰.۶۳ (۰.۰۵)	۱.۳۰ (۰.۰۹)	۱.۲۵ (۰.۱۵)	۰.۶۴ (۰.۰۵)	۰.۶۵ (۰.۰۶)	۱.۲۵ (۰.۱۲)	۱.۲۱ (۰.۱۶)	۰.۶۱ (۰.۰۷)	۰.۰۰	تابع توزیع
۰.۶۴ (۰.۰۵)	۱.۳۱ (۰.۱۱)	۱.۳۲ (۰.۱۲)	۰.۶۶ (۰.۰۶)	۰.۷۲ (۰.۰۹)	۱.۷۲ (۰.۳۸)	۱.۳۵ (۰.۴۸)	۰.۷۰ (۰.۱۰)	۰.۰۵	
۰.۶۰ (۰.۶۷)	۱.۱۴ (۰.۴۵)	۱.۲۰ (۰.۴۹)	۰.۶۰ (۰.۵۲)	۰.۶۱ (۰.۷۳)	۱.۳۲ (۰.۴۰)	۱.۳۸ (۰.۵۷)	۰.۶۲ (۰.۶۷)	۰.۲۰	

ب- داده‌های غیرنرمال و بدون داده دورافتاده ( $r = 0$ ). زمانی که با تبدیل توانی توزیع داده‌ها با توزیع نرمال اختلاف پیدا می‌کند، عملکرد روش  $GGM_0$  افت می‌کند، ولی همچنان سه روش دیگر عملکرد مشابه دارند. در مورد تبدیل تابع توزیع، هر چهار روش عملکرد نزدیک به هم در خوشه‌بندی داده‌ها دارند.

ج- داده‌های نرمال با داده دورافتاده ( $r \in \{0/0.5, 0/20\}$ ). زمانی که با سازوکار آلودگی قطعی، ۵ درصد از داده‌های نرمال یا با سازوکار آلودگی تصادفی، ۲۰ درصد از داده‌ها با داده‌های دورافتاده جای‌گزین می‌شوند، عملکرد روش متعارف  $GGM_0$  (و در یک مورد، روش  $GGM_1$ ) افت می‌کند، ولی روش‌های نرمال ناپارامتری  $NGMM_0$  و  $NGMM_1$ ، عملکرد بهتری دارند.

د- داده‌های غیرنرمال با داده دورافتاده ( $r \in \{0/0.5, 0/20\}$ ). در این حالت‌ها، روش‌های مورد بررسی بهتر مقایسه می‌شوند. زمانی که اندازه نمونه کوچک است، روش‌های نامتعارف  $GGM_1$  و  $NGMM_1$  عملکرد بهتری نسبت به روش‌های متعارف  $GGM_0$  و  $NGMM_0$  دارند و در برخی موارد،  $NGMM_1$  بهتر از  $GGM_1$  است. با افزایش اندازه نمونه، روش‌های مبتنی بر نرمال ناپارامتری  $NGMM_0$  و  $NGMM_1$ ، نسبت به روش‌های مبتنی بر نرمال  $GGM_0$  و  $GGM_1$ ، در مقابل آلودگی داده‌ها استوارتر هستند و عملکرد بهتری دارند.

جدول ۲. متوسط مقادیرهای پارامتر تنظیم‌کننده برای هر روش روی ۵۰ مجموعه داده شبیه‌سازی شده در بُعد  $p = ۲۰$  و اندازه نمونه  $n = ۱۵۰$ . انحراف استانداردها داخل پرانتزها نوشته شده‌اند.

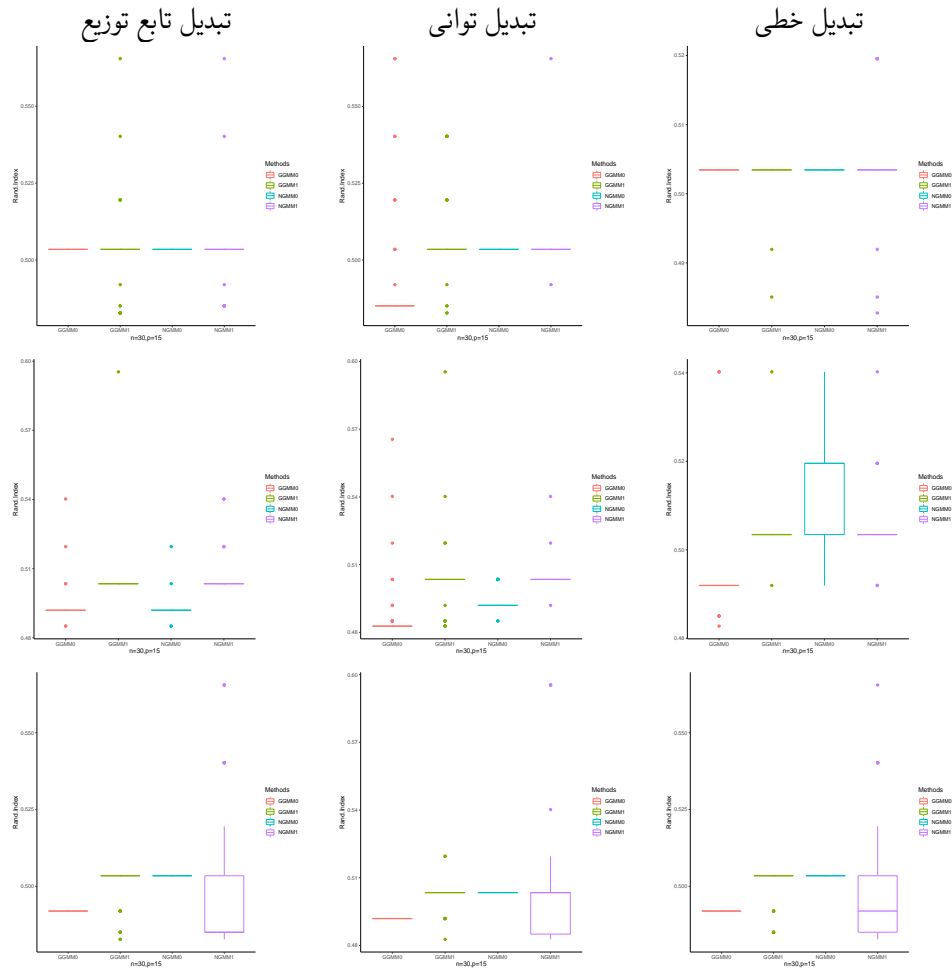
$NGMM_1$	$NGMM_0$			$GGMM_1$	$GGMM_0$				تبدیل
$\lambda = \bar{\lambda}_k$	$\bar{\lambda}_2$	$\bar{\lambda}_1$	$\lambda$	$\lambda = \bar{\lambda}_k$	$\bar{\lambda}_2$	$\bar{\lambda}_1$	$\lambda$	$r$	
۰.۳۶ (۰.۰۷)	۰.۷۱ (۰.۱۵)	۰.۷۱ (۰.۱۵)	۰.۳۵ (۰.۰۷)	۰.۳۷ (۰.۰۷)	۰.۷۶ (۰.۱۱)	۰.۷۶ (۰.۱۲)	۰.۳۸ (۰.۰۶)	۰.۰۰	
۰.۵۱ (۰.۰۵)	۱.۳۶ (۰.۶۸)	۱.۴۰ (۰.۶۷)	۰.۵۳ (۰.۰۵)	۰.۶۴ (۰.۰۶)	۶.۶۲ (۶.۲۵)	۷.۸۶ (۷.۰۲)	۰.۶۴ (۰.۰۵)	۰.۰۵	خطی
۰.۲۹ (۰.۳۰)	۳.۲۲ (۲.۴)	۳.۷۷ (۱.۵۵)	۰.۳۰ (۰.۳۷)	۰.۳۰ (۰.۴۰)	۲.۱۲ (۱.۵۱)	۱.۴۷ (۱.۴۸)	۰.۳۱ (۰.۴۴)	۰.۲۰	
۰.۳۴ (۰.۰۷)	۰.۶۷ (۰.۱۵)	۰.۶۶ (۰.۱۷)	۰.۳۳ (۰.۰۸)	۰.۳۴ (۰.۱۱)	۱.۶۳ (۲.۴۵)	۰.۹۷ (۰.۹۲)	۰.۳۴ (۰.۰۹)	۰.۰۰	
۰.۵۱ (۰.۰۵)	۱.۱۲ (۰.۵۵)	۱.۴۹ (۰.۵۶)	۰.۵۳ (۰.۰۵)	۰.۷۳ (۰.۰۷)	۵.۱۵ (۴.۵۲)	۵.۷۳ (۵.۸۱)	۰.۷۲ (۰.۰۶)	۰.۰۵	توانی
۰.۳۰ (۰.۰۴)	۲.۵۱ (۲.۶۳)	۶.۸۷ (۳.۴۵)	۰.۳۰ (۰.۰۴)	۰.۲۹ (۰.۰۳)	۱.۵۸ (۱.۲۳)	۱.۴۸ (۱.۲۲)	۰.۳۱ (۰.۰۴)	۰.۲۰	
۰.۳۵ (۰.۰۷)	۰.۶۸ (۰.۱۵)	۰.۶۷ (۰.۱۶)	۰.۳۴ (۰.۰۸)	۰.۳۵ (۰.۰۵)	۰.۷۲ (۰.۱۱)	۰.۷۲ (۰.۱۱)	۰.۳۶ (۰.۰۵)	۰.۰۰	
۰.۵۱ (۰.۰۵)	۱.۳۴ (۰.۶۲)	۱.۳۴ (۰.۶۲)	۰.۵۳ (۰.۰۵)	۰.۶۳ (۰.۰۶)	۵.۹۷ (۵.۹۹)	۷.۷۸ (۶.۹۰)	۰.۶۳ (۰.۰۵)	۰.۰۵	تابع توزیع
۰.۲۹ (۰.۰۳)	۱.۰۸ (۱.۹۹)	۲.۵۰ (۱.۷۷)	۰.۳۰ (۰.۰۴)	۴.۰۳۰ (۰.۰۴)	۱.۸۲ (۱.۵۳)	۱.۷۶ (۱.۵۵)	۰.۳۱ (۰.۰۵)	۰.۲۰	

#### ۴.۴ مقایسه برآورد ماتریس‌های دقت

در این بخش، دقت روش‌ها در برآورد درایه‌های ماتریس‌های دقت خوشه‌ها مقایسه شده است. دقت با استفاده از نُرم  $\sum_{k=1}^K \|\hat{\Omega}_k - \Omega_k\|_1$  محاسبه می‌شود که در آن  $\hat{\Omega}_k$ ، ماتریس دقت خوشه برآورد شده  $k$  و  $\Omega_k$ ، ماتریس دقت خوشه واقعی  $k$  است. نتایج در جدول ۳ آورده شده است. عملکرد ضعیف روش متعارف  $GGMM_0$  در خوشه‌بندی داده‌ها موجب شد که این روش، در برآورد ماتریس‌های دقت خوشه‌ها نیز ضعیف‌ترین عملکرد را داشته باشد. در مقابل، روش نامتعارف  $NGMM_1$  بالاترین دقت را نشان داده است. به‌طور کلی، روش‌های نامتعارف  $NGMM_1$  و  $GGMM_1$ ، دقت بالاتری نسبت به روش‌های متعارف  $NGMM_0$  و  $GGMM_0$  در برآورد پارامترهای مدل‌های آمیخته مورد بررسی داشته‌اند.

#### ۵ کاربرد داده واقعی

سرطان سینه، شایع‌ترین سرطان و دومین عامل مرگ و میر ناشی از سرطان در بین زنان است. از این‌رو، نظارت بر شیوع سرطان سینه برای دولت‌ها و سازمان‌های بهداشتی، از اهمیت بالایی برخوردار است. در تحقیقات علمی، داده‌های سرطان سینه به عنوان داده‌های ناهمگن شناخته شده است. روش‌های خوشه‌بندی روی داده‌های تشخیص

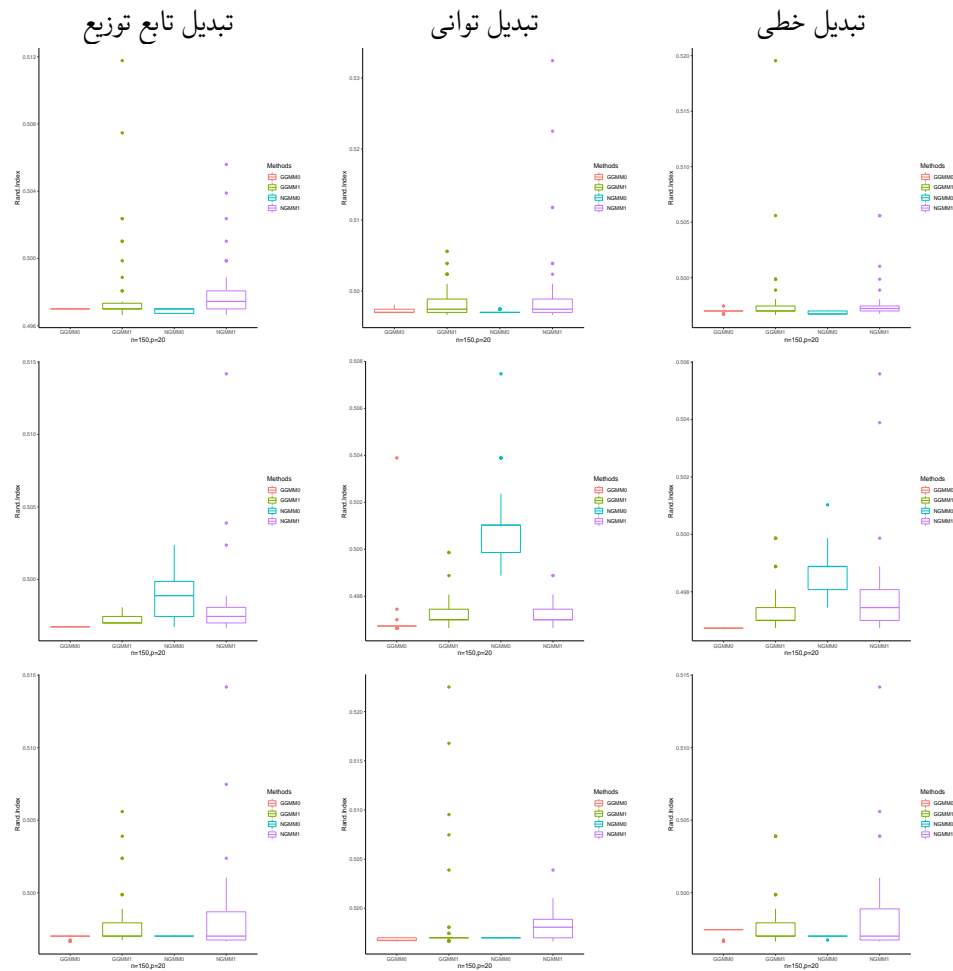


شکل ۱. مقایسه نتایج انتساب خوشه با استفاده از شاخص رند، روی ۵۰ مجموعه داده شبیه‌سازی شده با اندازه نمونه  $n = 30$  و بُعد  $p = 15$  برای روش‌های (به ترتیب از چپ به راست)  $GGMM_0$ ،  $GGMM_1$ ،  $NGMM_0$  و  $NGMM_1$ . داده‌ها فاقد آلودگی (سطر اول)، داده‌ها با سازوکار آلودگی قطعی ۵ درصد (سطر دوم) و داده‌ها با سازوکار آلودگی تصادفی ۲۰ درصد (سطر سوم).

سرطان سینه ویسکانسین<sup>۱</sup> (WDBC) (منگاسریان و همکاران، ۱۹۹۵) اعمال شده‌اند.<sup>۲</sup> این مجموعه داده شامل

<sup>1</sup>Wisconsin Diagnostic Breast Cancer

<sup>2</sup>[http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)).



شکل ۲. مقایسه نتایج انتساب خوشه با استفاده از شاخص رند، روی ۵۰ مجموعه داده شبیه‌سازی شده با اندازه نمونه  $n = 150$  و بُعد  $p = 20$  برای روش‌های (به ترتیب از چپ به راست)  $GGMM_0$ ،  $GGMM_1$ ،  $NGMM_0$  و  $NGMM_1$ . داده‌ها فاقد آلودگی (سطر اول)، داده‌ها با سازوکار آلودگی قطعی ۵ درصد (سطر دوم) و داده‌ها با سازوکار آلودگی تصادفی ۲۰ درصد (سطر سوم).

هسته‌های سلولی به‌دست آمده از تصویرهای دیجیتالی با سوزن ظریف آسپیراسیون<sup>۱</sup> (FNA) از هر توده سرطان سینه است. داده‌ها در ۱۰ ویژگی روی ۵۶۹ بیمار اندازه‌گیری شده که برای هر ویژگی، ۳ عدد به عنوان میانگین، خطای استاندارد و میانگین سه مقدار بزرگ‌تر ثبت شده است. بنابراین، این مجموعه داده شامل ۵۶۹ سطر، ۳۰ ستون اصلی

<sup>1</sup>Fine Needle Aspirate

جدول ۳. مقایسه دقت برآورد ماتریس‌های دقت هر دو خوشه برای هر روش روی ۵۰ مجموعه داده شبیه‌سازی شده. انحراف استانداردها داخل پرانتزها نوشته شده‌اند.

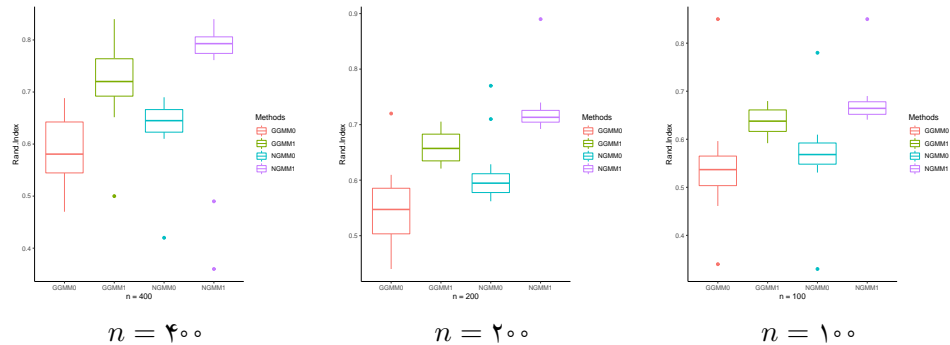
$NGMM_1$	$NGMM_0$	$GGMM_1$	$GGMM_0$	$r$	تبدیل	$(p, n)$
۸۷۸(۰٫۱۱)	۹۷۹(۰٫۱۰)	۸۸۳(۰٫۰۹)	۹۸۰(۰٫۰۸)	۰٫۵۰	خطی	
۸۷۹(۰٫۱۲)	۹۸۴(۰٫۱۲)	۸۹۴(۰٫۱۸)	۱۰۰۱(۰٫۱۳)	۰٫۵۵		
۸۷۲(۰٫۱۴)	۱۰۵۸(۰٫۲۹)	۸۷۵(۰٫۱۶)	۹۸۳(۰٫۱۳)	۰٫۲۰		
۸۷۶(۰٫۱۵)	۹۷۹(۰٫۱۱)	۹۰۳(۰٫۱۱)	۱۰۰۸(۰٫۱۱)	۰٫۵۰	توانی	(۱۵, ۳۰)
۸۸۰(۰٫۱۲)	۹۸۴(۰٫۱۲)	۹۱۰(۰٫۱۵)	۱۰۱۹(۰٫۱۰)	۰٫۵۵		
۸۷۱(۰٫۱۴)	۹۸۲(۰٫۱۷)	۸۷۵(۰٫۱۷)	۱۰۶۱(۰٫۲۰)	۰٫۲۰		
۸۷۸(۰٫۱۰)	۹۸۰(۰٫۱۰)	۸۸۲(۰٫۱۳)	۹۷۴(۰٫۱۷)	۰٫۵۰	تابع توزیع	
۸۸۰(۰٫۱۱)	۹۸۴(۰٫۱۳)	۸۹۳(۰٫۲۱)	۱۰۰۰(۰٫۱۶)	۰٫۵۵		
۸۷۲(۰٫۱۴)	۹۸۱(۰٫۱۴)	۸۷۴(۰٫۱۵)	۱۰۵۷(۰٫۲۹)	۰٫۲۰		
۸۶۴(۰٫۲۸)	۹۵۸(۰٫۳۳)	۸۶۴(۰٫۲۸)	۹۷۱(۰٫۲۶)	۰٫۵۰	خطی	
۹۱۰(۰٫۱۵)	۱۰۵۱(۰٫۱۱)	۹۳۹(۰٫۲۱)	۱۱۴۵(۰٫۰۶)	۰٫۵۵		
۸۴۳(۰٫۱۱)	۱۰۸۸(۰٫۵۸)	۸۴۵(۰٫۱۳)	۱۰۳۳(۰٫۲۴)	۰٫۲۰		
۸۵۴(۰٫۳۳)	۹۴۸(۰٫۳۷)	۸۵۱(۰٫۴۲)	۹۹۲(۰٫۵۵)	۰٫۵۰	توانی	(۲۰, ۱۵۰)
۹۰۹(۰٫۱۵)	۱۰۴۵(۰٫۱۰)	۹۶۲(۰٫۲۱)	۱۱۴۳(۰٫۲۱)	۰٫۵۵		
۸۴۴(۰٫۱۴)	۱۰۲۲(۰٫۲۰)	۸۴۳(۰٫۱۲)	۱۰۹۵(۰٫۵۱)	۰٫۲۰		
۸۵۸(۰٫۳۳)	۹۵۱(۰٫۳۶)	۸۶۰(۰٫۱۹)	۹۶۲(۰٫۲۵)	۰٫۵۰	تابع توزیع	
۹۱۰(۰٫۱۵)	۱۰۴۹(۰٫۱۰)	۹۳۵(۰٫۲۳)	۱۱۴۰(۰٫۱۶)	۰٫۵۵		
۸۴۳(۰٫۱۱)	۱۰۳۲(۰٫۳۰)	۸۴۸(۰٫۱۳)	۱۰۹۰(۰٫۵۷)	۰٫۲۰		

و یک ستون «تشخیص» با دو وضعیت تومور بدخیم و خوش‌خیم است.

هدف، پیش‌بینی متغیر «تشخیص» برای ارزیابی بدخیم یا خوش‌خیم بودن تومور هر بیمار است. بدین منظور، ابتدا مشاهدات به‌طور تصادفی به دو خوشه تخصیص داده شد و سپس، روش‌های خوشه‌بندی مبتنی بر مدل  $NGMM_1$ ،  $NGMM_0$ ،  $GGMM_1$ ،  $GGMM_0$ ، برای ارزیابی تاثیر اندازه نمونه بر عملکرد روش‌های خوشه‌بندی مورد مطالعه، ۵۰ نمونه بوت استرپ با اندازه‌های متفاوت  $n = 100, 200, 400$ ، بدون جای‌گذاری از جامعه استخراج شده است.

شکل ۳ نتایج بازیابی خوشه‌ها برای هر روش را منعکس می‌کند. طبق این شکل، روش  $NGMM_1$  نتایج متفاوتی ارائه می‌دهد و به‌طور قابل توجهی بهتر از سایر روش‌های مورد بررسی عمل می‌کند. همچنین روش  $GGMM_1$  اندکی بهتر از روش  $NGMM_0$  عمل می‌کند و روش متعارف  $GGMM_0$  بدترین عملکرد را دارد. شکل ۳ نشان می‌دهد که با افزایش اندازه نمونه، عملکرد روش‌های مورد نظر ارتقا می‌یابد.





شکل ۳. مقایسه نتایج انتساب خوشه با استفاده از شاخص رند، روی ۵۰ نمونه از مجموعه داده سرطان سینه WDBC با اندازه‌های نمونه  $n = 100, 200, 400$  و بُعد  $p = 30$  برای روش‌های (به ترتیب از چپ به راست)  $NGMM_0, NGMM_1, GGMM_0, GGMM_1$ .

## بحث و نتیجه‌گیری

خوشه‌بندی مبتنی بر مدل آمیخته گرافی نرمال ناپارامتری با دو فرم تابع تاوان  $\ell_1$  متعارف و نامتعارف ارائه شد و عملکرد آن با روش خوشه‌بندی مبتنی بر مدل آمیخته گرافی گاوسی با دو فرم تابع تاوان  $\ell_1$  متعارف و نامتعارف، هم از نظر بازسازی خوشه‌ها و هم از نظر برآورد پارامترهای مدل، مقایسه شد. همراه با تاوان  $\ell_1$  استاندارد، یک تاوان جای‌گزین در نظر گرفته شد که به نسبت‌های آمیخته  $\pi_k$  وابسته است، بنابراین، انتخاب تابع تاوان به‌عنوان کلید اصلی مقایسه در این مطالعه ظاهر شده است، ولی انتخاب پارامتر تنظیم‌کننده در مقایسه بی‌تاثیر بوده است.

در مطالعه شبیه‌سازی، عملکرد روش‌های مورد بررسی زمانی‌که توزیع داده‌ها نرمال یا غیرنرمال است، بررسی شد. برای مقایسه روش‌ها از نظر استحکام در برابر داده‌های دورافتاده، دو نوع سازوکار آلودگی قطعی و تصادفی در نظر گرفته شد. در برآورد ماتریس‌های دقت خوشه‌ها، روش‌های نامتعارف  $GGMM_1$  و  $NGMM_1$  دقت بالاتری نسبت به روش‌های متعارف  $GGMM_0$  و  $NGMM_0$  داشتند. در بازسازی خوشه‌ها، روش‌های خوشه‌بندی مبتنی بر مدل آمیخته گرافی نرمال ناپارامتری ( $NGMM_1$  و  $NGMM_0$ ) نسبت به روش‌های مبتنی بر مدل آمیخته گرافی گاوسی، عملکرد بهتری از خود بروز دادند، به خصوص در حضور داده‌های دورافتاده، استوارتر بودند. با این وجود، عملکرد هر چهار روش در بازیابی خوشه‌ها، نزدیک به هم بود و از کارایی مناسب برخوردار نبود ( $0.45 < RI < 0.55$ )، زیرا داده‌ها از مدل‌های آمیخته گاوسی و آمیخته نرمال ناپارامتری با دو مولفه با بردارهای میانگین صفر برای هر مولفه شبیه‌سازی شده‌اند. همچنین الگوهای I و II، ماتریس‌های همبستگی با درایه‌های بین  $[0, 1]$  تولید می‌کنند. از آنجایی‌که داده‌های تولید شده برای هر دو مولفه بسیار نزدیک به هم است، بازیابی خوشه‌ها دشوار است. بنابراین، شاخص‌های رند در محدوده ۵٪ محاسبه شده‌اند.

در نهایت، روش‌های خوشه‌بندی مورد نظر روی داده‌های سرطان سینه برای تشخیص توده‌های سرطانی خوش‌خیم

و بدخیم بین بیماران، به کار برده شد. نتایج حاکی از آن است که روش‌های خوشه‌بندی نامتعارف  $GMM_1$  و  $NGMM_1$ ، به‌طور قابل توجهی عملکرد بهتری در تشخیص وضعیت توده‌های سرطانی داشته‌اند و با افزایش اندازه نمونه، کارایی آن‌ها نیز افزایش می‌یابد ( $0.8 < RI < 0.9$ ).

به‌طور کلی، نتایج نشان داد که عملکرد روش‌های خوشه‌بندی مبتنی بر مدل، به انتخاب عبارت تاوان و انتخاب مدل بستگی دارد؛ به‌نحوی که ترکیبی از مدل آمیخته گرافی نرمال ناپارامتری با عبارت تاوان وابسته به نسبت‌های آمیخته (که در خلال این مطالعه به عنوان روش  $NGMM_1$  معرفی شد)، دقیق‌ترین خوشه‌بندی و برآورد پارامترها را ارائه می‌دهد که با افزایش اندازه نمونه، دقت آن بیشتر نمود پیدا می‌کند.

## تقدیر و تشکر

نویسندگان از نظرات روشن‌گرانه داوران محترم و پیشنهادهای ارزنده سردبیر محترم، هیئت تحریریه و ویراستار مجله علوم آماری که در بهبود کیفیت مقاله نقش مهمی داشته‌اند، کمال قدردانی و تشکر را می‌نمایند.

## مراجع

- Friedman, J. H., Hastie, T. and Tibshirani, R. (2008), Sparse Inverse Covariance Estimation with the Graphical Lasso, *Biostatistics*, **9**, 432-441.
- Khalili, A. and Chen, J. (2007), Variable Selection in Finite Mixture of Regression Models, *Journal of the American Statistical Association*, **102**, 1025-1038.
- Khalili, A., Eskandari, F. and Nematollahi, N. (2021), Estimation of Undirected Graph with Finite Mixture of Nonparanormal Distribution, *Journal of Statistical Theory and Practice*, **15**, 1-18.
- Lindsay, B. G. (1995), Mixture Models: Theory, Geometry and Applications, *In NSF-CBMS Regional Conference Series in Probability and Statistics*, 1-163.
- Liu, H., Han, F., Yuan, M., Lafferty, J. and Wasserman, L. (2012), High-dimensional Semiparametric Gaussian Copula Graphical Models, *The Annals of Statistics*, **40**, 2293-2326.

- Liu, H., Lafferty, J. and Wasserman, L. (2009), The Nonparanormal: Semiparametric Estimation of High Dimensional Undirected Graphs, *Journal of Machine Learning Research*, **10**, 2295-2328.
- Lotsi, A. and Wit, E. (2016), Sparse Gaussian Graphical Mixture Model, *Afrika Statistika*, **11**, 1041-1059.
- MacQueen, J. (1967), Some Methods for Classification and Analysis of Multivariate Observations, In *Fifth Berkeley Symposium on Mathematical Statistics and Probability*, **1**, 281-297.
- Mangasarian, O. L., Street, W. N. and Wolberg, W. H. (1995), Breast Cancer Diagnosis and Prognosis via Linear Programming, *Operations Research*, **43**, 570-577.
- McLachlan, G., and Peel, D. (2004), *Finite Mixture Models*, Wiley Series in Probabilities and Statistics, New York, NY: Wiley-Interscience.
- Meinshausen, N. and Bühlmann, P. (2006), High-dimensional Graphs and Variable Selection with the Lasso, *The Annals of Statistics*, **34**, 1436-1462.
- Mukherjee, S. and Hill, S. M. (2011), Network Clustering: Probing Biological Heterogeneity by Sparse Graphical Models, *Bioinformatics*, **27**, 994-1000.
- Rand, W. M. (1971), Objective Criteria for the Evaluation of Clustering Methods, *Journal of the American Statistical Association*, **66**, 846-850.
- Städler, N., Bühlmann, P. and Van de Geer, S. (2010),  $\ell_1$ -Penalization for Mixture Regression Models, *TEST*, **19**, 209-256.
- Yuan, M. and Lin, Y. (2007), Model Selection and Estimation in the Gaussian Graphical Model, *Biometrika*, **94**, 19-35.
- Zhou, H., Pan, W. and Shen, X. (2009), Penalized Model-based Clustering with Unconstrained Covariance Matrices, *Electronic Journal of Statistics*, **3**, 1473-1496.