

تعیین تعداد مؤلفه‌ها در توزیع آمیخته با مؤلفه‌های تی-نرمال چوله

محمد بهرامی^۱، فهیمه طورانی فرانی^۲

تاریخ دریافت: ۱۳۹۱/۸/۲۸

تاریخ پذیرش: ۱۳۹۶/۴/۲۵

چکیده:

تعیین تعداد مؤلفه‌ها در یک توزیع آمیخته، مسئله‌ای دشوار و حائز اهمیت است. برای تعیین تعداد بهینه مؤلفه‌ها در توزیع‌های آمیخته، روش‌های مختلفی وجود دارد که در این مقاله به ذکر چند مورد از آنها خواهیم پرداخت. روش اول که تحت عنوان الگوریتم greedy EM بیان شده، بر اساس الگوریتمی است که طی هر مرحله آن مؤلفه‌ای جدید به مدل اضافه می‌شود و این روند تا زمانی که منجر به تعیین تعداد بهینه مؤلفه‌ها در توزیع آمیخته شود، ادامه می‌یابد. روش دوم بر اساس ماکسیمم آنتروپی ادغام در تکرار زیر کلاس‌های روی هم افتاده تا زمانی است که در نتیجه ادغام این مؤلفه‌ها، توزیع آمیخته مورد بررسی دارای یک مؤلفه شود. این روش با عنوان ادغام آمیختگی شرح داده شده است و روش سوم نیز توسط تعریف متغیرهای نشانگر، به صورت ناپارامتری تعداد مؤلفه‌های توزیع آمیخته را تعیین می‌کند. شایان ذکر است که مؤلفه‌های توزیع آمیخته مورد نظر در این مقاله توزیع تی-نرمال چوله در نظر گرفته شده‌اند.

واژه‌های کلیدی: توزیع آمیخته متناهی، تعداد مؤلفه‌ها، الگوریتم greedy EM، آنتروپی، ادغام آمیختگی، توزیع تی-نرمال چوله.

تصادفی X در زیرجامعه‌های P_1 و P_2 و \dots و P_k باشند، آنگاه مدل آمیخته بر حسب توابع چگالی عبارت‌اند از:

$$f(x) = \pi_1 f_1(x) + \pi_2 f_2(x) + \dots + \pi_k f_k(x) \\ = \sum_{i=1}^k \pi_i f_i(x); \quad x \in S. \quad (2)$$

در حالت کلی لزومی ندارد که مؤلفه‌های f_1 و f_2 و \dots و f_k متعلق به یک خانواده از توابع چگالی باشند، اما در اغلب موارد، این اتفاق رخ می‌دهد. توابع چگالی f_i برای $i = 1, 2, \dots, k$ می‌توانند شامل بردار پارامتر $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{ih})^T$ باشند. بنا بر این یک چگالی آمیخته متناهی پارامتری با k مؤلفه به شکل کلی

$$f(x|\theta) = \sum_{i=1}^k \pi_i f(x|\theta_i); \quad \theta \in \Theta, \quad x \in S \quad (3)$$

پارامترهای مدل آمیخته می‌نامیم که در فضای پارامتر Θ تغییر

۱ معرفی توزیع آمیخته

فرض کنید جامعه آماری P شامل k ($k \geq 2$) زیر جامعه P_1 و P_2 و \dots و P_k باشد. هدف، بررسی خصوصیتی در این جامعه است که متغیر تصادفی X نشان‌دهنده اندازه آن خصوصیت است. اگر F_1 و F_2 و \dots و F_k تابع توزیع متغیر تصادفی X در زیر جامعه‌های P_1 و P_2 و \dots و P_k باشند، آنگاه تابع توزیع متغیر تصادفی X با تکیه‌گاه S در جامعه P به صورت

$$F(x) = \pi_1 F_1(x) + \pi_2 F_2(x) + \dots + \pi_k F_k(x), \quad x \in S$$

است، که در آن احتمال تعلق X به زیرجامعه P_i و دارای توزیع F_i است. برای $i = 1, 2, \dots, k$ $0 < \pi_i < 1$ و $\sum_{i=1}^k \pi_i = 1$ مقادیر π_i را نسبت‌های آمیختگی می‌نامند. مدل (۳) یک توزیع آمیخته متناهی با k مؤلفه بر حسب تابع توزیع است. چنانچه f_1 و f_2 و \dots و f_k به ترتیب توابع چگالی متغیر

^۱ عضو هیئت علمی گروه آمار، دانشگاه اصفهان، ایران

^۲ دانشجوی دکتری، دانشگاه اصفهان، ایران

توزیع آمیخته دلخواه به کار برد. دو روش دیگر در حالت کلی و برای هر توزیع آمیخته بیان شده‌اند.

۲ روش الگوریتم gEM

روش gEM^۳ اولین بار توسط ولاسیس و لیکاس [۴] معرفی شد. در این روش کمترین مقدار مؤلفه برای یک توزیع آمیخته در نظر گرفته می‌شود و سپس طی چند مرحله با اضافه کردن مؤلفه جدید به آن، مدل مناسب برای توصیف داده‌ها و در نتیجه تعداد بهینه مؤلفه‌ها برای توزیع آمیخته اولیه تعیین می‌شود. ابتدا توسط الگوریتم امید ریاضی ماکسیم‌سازی (EM)^۴ و یا امید ریاضی ماکسیم‌سازی شرطی (EMC)^۵ پارامترهای توزیع آمیخته اولیه با کمترین مؤلفه (k مؤلفه) برآورد می‌شود. توجه کنید که بردار پارامترهای هر مؤلفه به صورت

$$\theta_i = (\mu_i, \sigma_i^2, \alpha_i, v_i)^T$$

است.

فرض کنید مؤلفه جدید $g(y_j|\theta_{k+1})$ به توزیع آمیخته k مؤلفه‌ای اولیه، $f(y_j|\theta_k)$ ، اضافه شده است، آنگاه توزیع آمیخته جدید به صورت

$$f(y_j|\theta_{k+1}) = (1-a)f(y_j|\theta_k) + ag(y_j|\theta_{k+1}) \quad (۴)$$

است، که در آن $0 < a < 1$ و $\theta_{k+1} = (\theta_k, a, \theta_{k+1}^T)^T$ و $\theta_{k+1} = (\mu_{k+1}, \sigma_{k+1}^2, \alpha_{k+1}, v_{k+1})^T$ است. a و $1-a$ به ترتیب نسبت‌های آمیختگی برای توزیع آمیخته اولیه و مؤلفه جدید هستند. سپس پارامترهای جدید a و θ_{k+1} توسط ماکسیم‌سازی لگاریتم تابع درست‌نمایی

$$\begin{aligned} \mathcal{L}_{k+1} &= \sum_{j=1}^n \log f(y_j|\theta_{k+1}) \\ &= \sum_{j=1}^n \log \left[(1-a)f(y_j|\theta_k) + ag(y_j|\theta_{k+1}) \right] \end{aligned} \quad (۵)$$

برآورد می‌شوند. با توجه به تابع درست‌نمایی فوق مشاهده می‌شود که پارامترهای توزیع آمیخته اولیه در طول ماکسیم

^۳ greedy EM

^۴ Expectation and Maximization

^۵ Conditional Expectation and Maximization

می‌کند. به عنوان مثال تابع چگالی توزیع آمیخته با مؤلفه‌های تی-نرمال چوله به صورت

$$f(y|\theta) = \sum_{i=1}^k \pi_i g(y|\mu_i, \sigma_i^2, \alpha_i, v_i)$$

است، به طوری که

$$g(y|\mu, \sigma^2, \alpha, v) = \gamma t(y|\mu, \sigma^2, v) \Phi \left(\alpha \frac{y-\mu}{\sigma} \right)$$

و نمادهای t و Φ به ترتیب تابع چگالی احتمال توزیع تی و تابع توزیع تجمعی نرمال هستند و μ و α نیز به ترتیب پارامترهای مکان و شکل بوده و متعلق به اعداد حقیقی \mathbb{R} هستند همچنین σ و v به ترتیب پارامترهای مقیاس و درجه آزادی هستند که متعلق به بازه $(0, +\infty)$ هستند. بردار پارامترهای توزیع، θ ، به صورت $\theta = (\pi_1, \pi_2, \dots, \pi_k, \theta_1^T, \theta_2^T, \dots, \theta_k^T)^T$ بوده که در آن $0 < \pi_i < 1$ و $\sum_{i=1}^k \pi_i = 1$ ، $\theta_i = (\mu_i, \sigma_i^2, \alpha_i, v_i)^T$ پارامتر k ، تعداد مؤلفه‌ها در یک توزیع آمیخته، از اهمیت زیادی برخوردار است. زیرا اگر k را بسیار بزرگ در نظر بگیریم، ممکن است میانگین برخی از زیرجامعه‌ها بسیار شبیه به یکدیگر باشند و یا π_j (احتمال تعلق به زیرجامعه j -ام) متناظر با برخی از آنها بسیار کوچک باشد و چنانچه k را مقداری کوچک در نظر بگیریم، آنگاه ممکن است μ_j ها نماینده خوبی برای هر گروه نباشند و یا به عبارت دیگر زیرجامعه‌ها دارای واریانس بزرگی باشند. نکته دیگر که اهمیت معلوم بودن تعداد مؤلفه‌های یک توزیع آمیخته را روشن تر می‌کند، تعیین تعداد پارامترهای مدل است. زیرا چنانچه k مقداری معلوم باشد، تعداد پارامترهای توزیع نیز معلوم است. به عنوان مثال یک توزیع آمیخته شامل k مؤلفه نرمال با میانگین و واریانس متفاوت، دارای $2k + (k-1)$ پارامتر است. بنا بر این با توجه به اهمیت تعداد مؤلفه‌ها در یک توزیع آمیخته، روش‌های زیادی برای تعیین تعداد بهینه این مؤلفه‌ها بیان شده است که در این مقاله به بررسی سه مورد از آنها می‌پردازیم. الگوریتم gEM برای حالت خاص توزیع آمیخته با مؤلفه‌های تی-نرمال چوله بیان شده، ولی روش کلی آن را می‌توان برای هر

سازی \mathcal{L}_{k+1} ثابت باقی می ماند.

$$\hat{\alpha}_{k+1}^{(m+1)} = \frac{\sum_{j=1}^n \hat{z}_j^{(m)} \hat{\gamma}_{\lambda_j}^{(m)} \hat{\mu}_j^{(m+1)}}{\sum_{j=1}^n \hat{z}_j^{(m)} \hat{\mu}_j^{(m+1)}} \\ \hat{v}_{k+1}^{(m+1)} = \operatorname{argmax} \left\{ \frac{v}{\psi} \log \frac{v}{\psi} - \log \Gamma \left(\frac{v}{\psi} \right) + \frac{v}{\psi} \hat{b}_{\psi}^{(m)} \right\} \quad (V)$$

است که در آن

$$\hat{u}_j^{(m+1)} = (y_j - \hat{\mu}^{(m+1)}) / \hat{\sigma}^{(m+1)} \\ \hat{b}_{\lambda}^{(m+1)} = \sum_{j=1}^n \hat{z}_j^{(m)} \hat{t}_j^{(m)} y_j \\ \hat{b}_{\psi}^{(m+1)} = \sum_{j=1}^n \hat{z}_j^{(m)} y_j \\ \hat{b}_{\psi}^{(m+1)} = \sum_{j=1}^n \hat{z}_j^{(m)} \hat{\gamma}_{\lambda_j}^{(m)} \\ \hat{b}_{\psi}^{(m+1)} = \sum_{j=1}^n \hat{z}_j^{(m)} \left(\hat{k}_j^{(m)} \hat{t}_j^{(m)} \right) / \sum_{j=1}^n \hat{z}_j^{(m)}$$

هستند. گام های الگوریتم EM جزئی را تا رسیدن به همگرایی به کار می بریم. به عنوان مثال تا وقتی که

$$\left| \frac{\hat{\mathcal{L}}_{k+1}^{(m)}}{\hat{\mathcal{L}}_{k+1}^{(m-1)}} - 1 \right| < 10^{-6}$$

گام های الگوریتم را تکرار می کنیم.

اگرچه الگوریتم EM جزئی روشی کارآمد در برآوردیابی پارامترهای توزیع آمیخته حاصل است، اما حساسیت این روش به مقادیر اولیه a و μ_{k+1} غیرقابل اغماض است. بنا بر این برای یافتن مقادیر اولیه مناسب این پارامترها تدبیر جست و جوی کلی را به کار می بریم. اساس این روش به منظور یافتن یک مقدار اولیه مناسب و همچنین تسهیل جست و جوی کامل در فضای پارامتر، جایگزینی تابع درست نمایی (۵) با بسط تیلور در نقطه $a = a$ و سپس یافتن مقدار بهینه پارامترهای Θ_{k+1} است. به ویژه اگر از بسط تیلور مرتبه دو استفاده کنیم، تقریب زیر را مشاهده خواهیم کرد.

$$\hat{\mathcal{L}}_{k+1} = \mathcal{L}_{k+1}(a.) + \frac{\left| \dot{\mathcal{L}}_{k+1}(a.) \right|^2}{2 \ddot{\mathcal{L}}_{k+1}(a.)} \quad (۸)$$

که در آن $\dot{\mathcal{L}}_{k+1}(a.)$ و $\ddot{\mathcal{L}}_{k+1}(a.)$ مشتق های مرتبه اول و دوم $\mathcal{L}_{k+1}(a.)$ حول نقطه $a = a$ است.

نکته حائز اهمیت در این روش تبدیل یک توزیع آمیخته با $k+1$ مؤلفه به توزیعی آمیخته با دو مؤلفه است که منجر به سهولت در محاسبات برآوردیابی می شود. برای برآورد پارامترهای Θ_{k+1} ، پارامترهای توزیع آمیخته جدید با دو مؤلفه، استفاده از الگوریتم EM مناسب است. اما با توجه به این که پارامترهای Θ_k در طول فرایند برآوردیابی به عنوان یک مقدار ثابت عمل می کنند، این الگوریتم بخشی از تابع درست نمایی را ماکسیم می کند و از این رو الگوریتم EM جزئی نامیده می شود. محاسبات متناظر با این برآوردها مشابه با محاسبات متناظر با برآورد پارامترهای توزیع آمیخته اولیه با k مؤلفه است. گام های الگوریتم EM جزئی به شرح زیر است.

- گام امید ریاضی جزئی: در این مرحله امید ریاضی شرطی متغیرهای پنهان به شرط مشاهدات در m امین تکرار محاسبه می شود. در این صورت

$$\hat{z}_j^{(m)} = \frac{\hat{a}^{(m)} g(y_j | \hat{\Theta}^{(m)})}{(1 - \hat{a}^{(m)}) f(y_j | \hat{\Theta}_k^{(m)}) + \hat{a}^{(m)} g(y_j | \hat{\Theta}^{(m)})} \\ \hat{t}_j^{(m)} = \frac{\hat{v}^{(m)} + 1}{\hat{v}^{(m)} + \hat{u}_j^{(m)}} \\ \hat{\gamma}_{\lambda_j}^{(m)} = \hat{\alpha}^{(m)} \hat{u}_j^{(m)} + \frac{\varphi(\hat{\alpha}^{(m)} \hat{u}_j^{(m)})}{\Phi(\hat{\alpha}^{(m)} \hat{u}_j^{(m)})} \\ \hat{k}_j^{(m)} = DG \left(\frac{\hat{v}^{(m)} + 1}{\psi} \right) - \log \left(\frac{\hat{v}^{(m)} + \hat{u}_j^{(m)}}{\psi} \right) \quad (۶)$$

است، به طوری که $\hat{u}_j^{(m)} = (y_j - \hat{\mu}^{(m)}) / \hat{\sigma}^{(m)}$ است.

- گام ماکسیم سازی جزئی: پارامترهای جدید موجود در (a, Θ_{k+1}) ، توسط ماکسیم سازی امید ریاضی شرطی تابع درست نمایی (۵) به شرط مشاهدات، برآورد می شوند. برآوردهای حاصل به صورت

$$\hat{a}^{(m+1)} = \frac{\sum_{j=1}^n \hat{z}_j^{(m)}}{n} \\ \hat{\mu}_{k+1}^{(m+1)} = \frac{\hat{b}_{\lambda}^{(m)} + \hat{\alpha}^{(m)} \hat{b}_{\psi}^{(m)} - \hat{\sigma}^{(m)} \hat{\alpha}^{(m)} \hat{b}_{\psi}^{(m)}}{\sum_{j=1}^n \hat{z}_j^{(m)} \hat{t}_j^{(m)} + \hat{\alpha}^{(m)} \sum_{j=1}^n \hat{z}_j^{(m)}} \\ \hat{\sigma}_{k+1}^{(m+1)} = \frac{\sum_{j=1}^n \hat{z}_j^{(m)} \hat{t}_j^{(m)} (y_j - \hat{\mu}^{(m+1)})^2}{\sum_{j=1}^n \hat{z}_j^{(m)}}$$

۳. مقادیر اولیه مناسب برای a و Θ_{k+1} را توسط روش جست و جوی کلی انتخاب می‌کنیم.

۴. پارامترهای a و Θ_{k+1} را با استفاده از الگوریتم EM جزئی برآورد می‌کنیم.

۵. اگر $\hat{\mathcal{L}}_{k+1} \leq \hat{\mathcal{L}}_{k+p}$ باشد (که در آن $p > 0$ یک عبارت جریمه است)، الگوریتم gEM پایان می‌یابد. در غیر اینصورت یک مؤلفه جدید در نظر می‌گیریم، به مرحله (۲) بازگشته و k را برابر با $k+1$ قرار می‌دهیم.

۳ روش ادغام آمیختگی

بادری و همکاران [۱] روش دیگری جهت تعیین تعداد بینه مؤلفه‌ها در توزیع‌های آمیخته، تحت عنوان ادغام آمیختگی^۶ معرفی کردند. چنانچه داده‌ها به علت داشتن توزیعی با شکل غیر معمول و زیرجامعه‌های غیر قابل تفکیک^۷ و یا دارای ویژگی‌های خاصی باشند که منجر به روی هم افتادن مؤلفه‌ها و در نتیجه بیش‌برآورد تعداد آنها شود، روش ادغام آمیختگی یک برآورد مقاوم^۸ برای تعداد مؤلفه‌های توزیع ارائه می‌کند. در حالی که الگوریتم gEM یک روش مناسب برای انتخاب تعداد مؤلفه‌های یک توزیع آمیخته تحت شرایط متعارف (مانند محدب بودن مؤلفه‌ها) فراهم می‌کند. نخست روش ادغام آمیختگی برای توزیع‌های آمیخته با مؤلفه‌های متقارن، طراحی شد. اما به‌منظور رفع نیاز به مؤلفه‌هایی که در حالت متقارن توضیح دهنده چولگی توزیع واقعی داده‌ها هستند، تکنیک ادغام آمیختگی برای توزیع‌های آمیخته با مؤلفه نامتقارن، گسترش یافت. به عنوان مثال چنانچه k مؤلفه‌ی یک توزیع آمیخته دارای توزیع نرمال چوله باشد و بخواهیم این توزیع آمیخته را توسط یک توزیع آمیخته دیگر با مؤلفه‌های نرمال بازنویسی کنیم، در این صورت توزیع آمیخته جدید دارای $2k$ مؤلفه نرمال است. در حقیقت هر توزیع نرمال چوله را می‌توان توسط یک توزیع آمیخته با دو مؤلفه نرمال تقریب زد که یکی از این مؤلفه‌ها در واقع توضیح دهنده چولگی

با بررسی رابطه فوق می‌توان نشان داد که ماکسیم \mathcal{L}_{k+1} حول نقطه $a_0 = 0.5$ به صورت

$$\hat{\mathcal{L}}_{k+1} = \sum_{j=1}^n \log \left(\frac{f(y_j|\hat{\Theta}_k) + g(y_j|\Theta_{k+1})}{2} \right) + \frac{\left[\sum_{j=1}^n \delta_j(\Theta_{k+1}) \right]^2}{2 \sum_{j=1}^n \delta_j^2(\Theta_{k+1})} \quad (9)$$

است. به طوری که

$$\delta_j(\Theta_{k+1}) = \frac{f(y_j|\hat{\Theta}_k) - g(y_j|\Theta_{k+1})}{f(y_j|\hat{\Theta}_k) + g(y_j|\Theta_{k+1})}$$

بنا بر این مقدار مناسب نسبت آمیختگی a عبارت‌اند از

$$\hat{a} = \frac{1}{2} \left(1 - \frac{\sum_{j=1}^n \delta_j(\Theta_{k+1})}{2 \sum_{j=1}^n \delta_j^2(\Theta_{k+1})} \right). \quad (10)$$

چنانچه مقدار \hat{a} خارج از بازه $(0, 1)$ قرار گیرد، می‌توان برای $k=1$ مقدار \hat{a} را برابر با 0.5 و برای $k \geq 2$ را برابر با $\frac{2}{k+1}$ در نظر گرفت.

هو و همکاران [۳] ریشه پنجم نیمی از واریانس نمونه $(\sqrt[5]{s_y^2/2})$ را یک انتخاب مناسب برای $\hat{\sigma}_{k+1}^{(.)}$ و همچنین مقادیر ثابت صفر و ده را برای $\hat{\nu}_{k+1}^{(.)}$ و $\hat{\alpha}_{k+1}^{(.)}$ در نظر گرفتند. همچنین چندک‌های ۵-ام، ۱۰-ام و... و ۹۵-ام مشاهدات y را یافته و سپس مقداری از آنها که رابطه (؟؟) را ماکسیم می‌کند را به عنوان $\hat{\mu}_{k+1}^{(.)}$ پیشنهاد کردند.

در حالت کلی، گام‌های الگوریتم gEM به شرح زیر است.

۱. با مقدار $k=1$ شروع کرده و پارامترهای تک مؤلفه توزیع آمیخته اولیه را برآورد می‌کنیم. چنانچه الگوریتم در گام اول با توزیع آمیخته بیشتر از یک مؤلفه ($k > 1$) آغاز شود، پارامترهای Θ_k با استفاده از الگوریتم EM قابل برآورد هستند.
۲. توزیع آمیخته جدید با دو مؤلفه را به صورت رابطه (؟؟) می‌سازیم.

^۶ Merging mixture

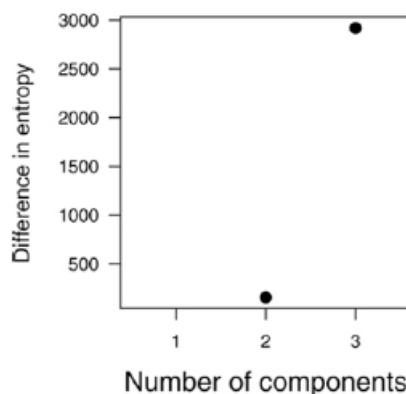
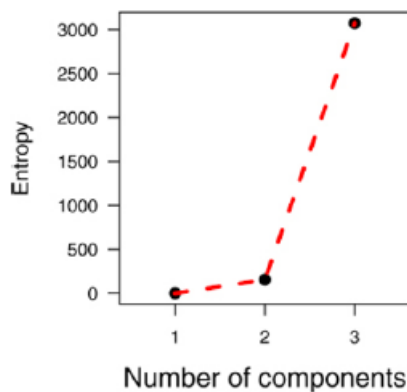
^۷ low separation

^۸ Robust

نمودار آنتروپی زیر کلاسها در مقابل تعداد آنها دارد و یا نقطه ماکسیم نمودار تعداد زیر کلاسها در مقابل تفاضل آنتروپی آنها پس از ادغام است. به عنوان مثال نمودارهای موجود در شکل ۱ (الف) بیان کننده تعداد دو مؤلفه و نمودارهای شکل ۱ (ب) بیان کننده تعداد ده مؤلفه برای داده‌های متناظرشان هستند.

۴ روش ناپارامتری تعیین تعداد مؤلفه‌های توزیع آمیخته

در این بخش برای محاسبه برآورد تعداد مؤلفه‌های مدل آمیخته روشی متفاوت از دو روش مطرح شده در بخش‌های قبل ارائه می‌کنیم. در حقیقت این روش، تکنیکی ناپارامتری است که توسط تعریف متغیرهای نشانگر تعداد بهینه مؤلفه‌های مدل را به شرح زیر است.



(الف)

توزیع هستند. ماکسیم آنتروپی ادغام در تکرار زیر کلاس‌های روی هم افتاده تا زمانی که در نتیجه این ادغام‌ها، توزیع آمیخته مورد بررسی دارای یک مؤلفه شود، نقش اساسی را در این روش ایفا می‌کند. روش ادغام آمیختگی طی گام‌های زیر اجرا می‌شود.

۱. ابتدا بیشترین تعداد مؤلفه برای توزیع آمیخته در نظر گرفته می‌شود و میانگین آنتروپی آن توسط فرمول

$$Ent(k) = - \sum_{j=1}^n \sum_{i=1}^k \hat{z}_{ij} \log \hat{z}_{ij} \geq 0 \quad (11)$$

محاسبه می‌شود. که در آن \hat{z}_{ij} نشان دهنده احتمال پسین پیشامد $(\Theta_k = \hat{\Theta}_k)$ است.

۲. برای تمام زوج زیر کلاسهای ممکن (ℓ, ℓ') دو زیر کلاس ℓ و ℓ' را با یکدیگر ترکیب کرده تا ضابطه زیر ماکسیم شود.

$$\left(- \sum_{i=1}^n \{ \hat{z}_{i\ell} \log \hat{z}_{i\ell} + \hat{z}_{i\ell'} \log \hat{z}_{i\ell'} \} + \sum_{i=1}^n (\hat{z}_{i\ell} + \hat{z}_{i\ell'}) \log (\hat{z}_{i\ell} + \hat{z}_{i\ell'}) \right)$$

۳. سپس دو زیر کلاس مذکور با یکدیگر ادغام می‌شود.

۴. آنتروپی حاصل پس از ادغام به صورت

$$Ent(k-1) = - \sum_{j=1}^n \left\{ \sum_{i \neq \ell, \ell'} \hat{z}_{ij} \log \hat{z}_{ij} + \hat{z}_{i, \ell \cup \ell'} \log \hat{z}_{i, \ell \cup \ell'} \right\} \quad (12)$$

محاسبه می‌شود، که در آن $\hat{z}_{i, \ell \cup \ell'} = \hat{z}_{i, \ell} + \hat{z}_{i, \ell'}$ احتمال پسین متناظر با زیر کلاس حاصل از ادغام یعنی زیر کلاس $\ell \cup \ell'$ است.

۵. مقدار \hat{z}_j که شامل احتمال پسین ادغام شده و ادغام نشده است، به هنگام می‌شود.

۶. مقدار k برابر با $k-1$ در نظر گرفته می‌شود، به مرحله ۲ بازگشته و این کار تا $k=1$ تکرار می‌شود.

۷. در نهایت تعداد بهینه مؤلفه‌ها در توزیع آمیخته مورد بررسی برابر با نقطه‌ای است که یک پرش ناگهانی در

جدول ۱. انتخاب تعداد مؤلفه‌های مدل

→	۱	۲	...	M
↓				
y_1				
y_2				
...				
y_N				

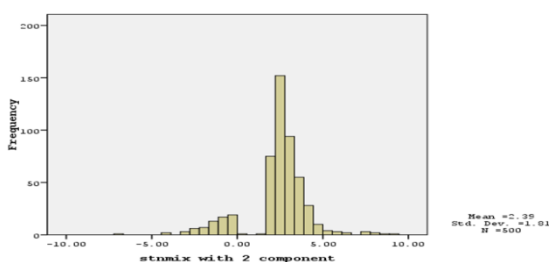
با به کارگیری برنامه زیر در نرم‌افزار Open-Bugs می‌توان تعداد مؤلفه‌های یک توزیع آمیخته را توسط روش ناپارامتری فوق تعیین کرد.

blueModel

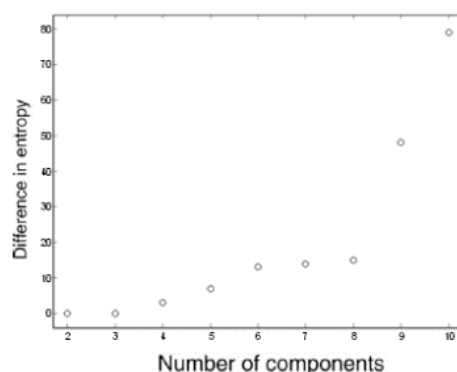
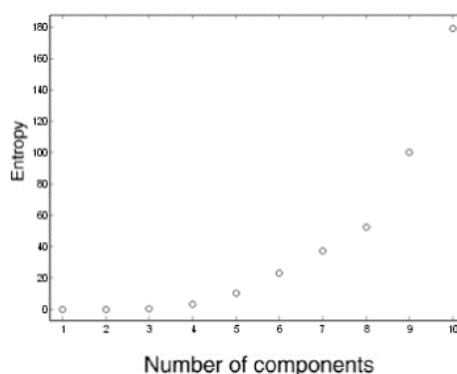
```
{for (i in 1:n) {
for ( j in 1:M) { Memb[i,j] <-
equals(y[i], j)}}
Mstar<- sum(Fmemb[]) ; for (j in 1:M){ Fmemb[j] <-
step(sum(Memb[,j])-1)}}}
```

۵ بررسی تعداد مؤلفه‌های داده‌های شبیه‌سازی شده از توزیع آمیخته

در این بخش ابتدا یک نمونه ۵۰۰ تایی از توزیع آمیخته تی-نرمال چوله به گونه‌ای تولید می‌کنیم که دارای دو زیرجامعه StN با پارامترهای (۳، ۱، -۳، ۴) و (۲، ۱، ۵، ۴) باشند. سپس با استفاده از روش ناپارامتری ارائه شده در این مقاله، تعداد مؤلفه‌های این مجموعه داده‌ها را برآورد می‌کنیم. نمودار فراوانی داده‌های تولیدی به صورت زیر است.



شکل ۲. نمودار فراوانی داده‌های شبیه‌سازی شده از توزیع آمیخته با دو مؤلفه StN



(ب)

شکل ۱. نمودار تعداد مؤلفه‌های توزیع آمیخته در مقابل آنتروپی توزیع و تفاضل آنتروپی‌های متناظر

ابتدا بیشترین تعداد مؤلفه (M) را برای توزیع آمیخته مورد نظر، فرض می‌کنیم. این تعداد را می‌توان با توجه به (نمودار فراوانی) بافت‌نگاشت داده‌ها تخمین زد. سپس جامعه را به M زیرجامعه‌ها چه تفکیک کرده و بررسی می‌کنیم که هر یک از زیرجامعه‌ها چه تعداد داده را تحت پوشش خود قرار می‌دهند و در نهایت تعداد زیرجامعه‌هایی که شامل هیچ داده‌ای نیستند (r) را شناسایی کرده و از زیر جامعه‌ها حذف می‌کنیم. بنا بر این تعداد مؤلفه‌های برآورد شده توسط این روش برابر با $M - r$ است. جدول ۱ چگونگی این فرایند را بیان می‌کند. زمانی که مقدار M مشخص شد، اگر مشاهده i -ام متعلق به زیرجامعه

z -ام باشد، در خانه سطر i -ام و ستون z -ام عدد یک و در بقیه خانه‌های آن سطر عدد صفر را قرار می‌دهیم. این عمل را برای $i = 1, 2, \dots, N$ و به ازای $M, \dots, 2, 1 = z$ تکرار می‌کنیم. در نهایت اعداد هر ستون را جمع کرده و سپس مشاهده می‌شود که تعداد ستون‌هایی که مجموع متناظر با آنها بزرگ‌تر از یک باشد، تعداد مؤلفه‌های مدل را نشان می‌دهند.

۶ نتیجه گیری

با توجه به اهمیت پارامتر تعداد مؤلفه‌های یک توزیع آمیخته، در این مقاله سه روش متفاوت برای تعیین تعداد پارامترهای توزیع آمیخته بیان شده است. چنانچه تنها در مورد آمیخته بودن توزیع داده‌های مورد بررسی اطلاع داشته باشیم و از توزیع مؤلفه‌های آن اطلاعی نداشته باشیم، روش ناپارامتری بیان شده در بخش ۴، استفاده می‌شود و در صورتی که الگوریتم EM از نظر اجرایی و تحلیلی در توزیع آمیخته مورد بررسی ساده باشد، الگوریتم gEM و همچنین اگر توزیع داده‌ها دارای شکلی نامعمول با زیرجامعه‌های غیر قابل تفکیک باشد روش ادغام آمیختگی به عنوان روشی مناسب توصیه می‌شود.

با توجه به این نمودار، واضح است که داده‌ها دارای دو زیر جامعه هستند، ولی با این وجود عدد M ، یعنی بیشترین تعداد مؤلفه در این توزیع آمیخته را با توجه به نمودار عدد ۴ در نظر گرفته و با اجرای این برنامه تعداد بهینه مؤلفه‌ها تعیین می‌شود. جدول ۲ حاصل از اجرای برنامه مذکور برای داده‌های شبیه‌سازی شده است. با توجه به نتایج جدول ۲، تعداد مؤلفه‌ها، ۲ برآورد شده است.

جدول ۲. جدول حاصل

از خروجی نرم‌افزار این باگز، شامل برآورد تعداد مؤلفه‌های توزیع آمیخته

	mean	sd	MC_error	val2.5pc	median	val97.5pc	start	sample
Mstar	2.0	0.0	2.236E-12	2.0	2.0	2.0	2001	2000

مراجع

- [1] Baudry, J. P., Raftery, A. F., Celeux, G., Lo, K. and Gottardo, R. (2010). Combining Mixture Components for Clustering. *Journal of Computational and Graphical Statistics* **9**, 332-353.
- [2] Congdon, P. (2006). Bayesian Statistical Modeling. (Second Edition). *Wiley Series in probability and statistics*.
- [3] Ho, H. J., Lin, T.I., Chang, H. H., Haase, S. B., Huang, S. and Pyne, S. (2011). Parametric Modeling of Cellular State Transitions as Measured with Flow Cytometry, *from first IEEE international Conference on Computational Advances in Bio and medical Sciences*.
- [4] Vlassis, N. and Likas, A. (2002). A Greedy EM Algorithm for Gaussian Mixture Learning. *Neural Processing Letters* **15**, 77-87.