

توزیع‌های آمیخته متناهی با مولفه‌های نرمال چوله

محمد بهرامی^۱

چکیده:

یک توزیع آمیخته ترکیبی از دو یا چند توزیع آماری است. وقتی نمونه‌گیری از یک جامعه غیرهمگن متشکل از دو یا چند زیرجامعه هر یک با توزیعی متفاوت انجام می‌گیرد، با توزیعی آمیخته سر و کار پیدا می‌کنیم. برای مثال، توزیع زمان خرابی آمیخته‌ای از مؤلفه‌های سالم و معیوب، توزیع وزن حیوانات با گروه‌های سنی متفاوت یا مدت زمانی که بیماران قلبی پس از یک عمل جراحی در گروه‌های سنی متفاوت زندگی خواهند کرد. در چنین حالاتی نمودار منحنی مربوط به تابع چگالی اندازه خصوصیت مورد نظر در جامعه آماری می‌تواند دارای چند مد یا نما باشد. همچنین در بسیاری از موارد با خصوصیت‌هایی مواجه می‌شویم که مجموعه مقادیر آنها می‌تواند متعلق به مجموعه اعداد حقیقی باشد و در واقع چون این‌گونه خصوصیت‌ها معمولاً به طور طبیعی پدیدار می‌گردند، به نظر می‌رسد که توزیع نرمال برای آنها مناسب باشد. اما با مطالعه بیشتر در می‌یابیم که وضعیت خصوصیت مورد مطالعه برای افراد جامعه کاملاً متقارن نبوده بلکه به راست یا به چپ دارای چولگی می‌باشد. یعنی تابع چگالی مورد نظر دارای چولگی به راست یا به چپ می‌باشد. بنابراین متغیرهایی با مجموعه مقادیر اعداد حقیقی دارای عدم تقارن مورد نظر خواهد بود. این موضوع می‌تواند منجر به تعریف توزیعی به نام توزیع نرمال چوله گردد. در این مقاله حالتی را بررسی می‌کنیم که مؤلفه‌های یک توزیع آمیخته از توزیع‌های نرمال چوله تشکیل شده‌اند. بنابراین ابتدا توزیع آمیخته و سپس توزیع نرمال چوله را بررسی می‌کنیم و در پایان توزیع آمیخته با مؤلفه‌های نرمال را تعریف و برآورد پارامترهای موجود در این مدل را با استفاده از الگوریتم EM به دست می‌آوریم.

واژه‌های کلیدی: الگوریتم EM ، تابع درست‌نمایی، توزیع آمیخته، توزیع نرمال چوله، چگالی آمیخته، متغیرهای پنهان.

۱ مقدمه

در بسیاری موارد با رسم منحنی فراوانی نسبی داده‌ها و مشاهده دو یا چند نقطهٔ ماکسیمم در آنها، پی می‌بریم که جامعه آماری از لحاظ خصوصیت مورد بررسی یک جامعه ناهمگن^۲ می‌باشد. نمودار فراوانی نسبی زیر با توجه به دونمایی^۳ بودن آن حاکی از ناهمگنی جامعه آماری مورد نظر است. این نمودار در واقع می‌تواند معرف

^۱ استادیار گروه آمار دانشگاه اصفهان
^۲ Non-Homogeneous Population
^۳ Bimodal

یک تابع چگالی آمیخته باشد.

تصادفی در زیر جامعه‌های P_1, P_2, \dots, P_k و f تابع چگالی احتمال X در جامعه P باشد، آنگاه مدل آمیخته را می‌توان به صورت زیر بر حسب تابع چگالی نوشت:

$$\begin{aligned} f(x) &= \pi_1 f_1(x) + \dots + \pi_k f_k(x) \\ &= \sum_{i=1}^k \pi_i f_i(x), x \in S \end{aligned} \quad (2)$$

در حالت کلی لزومی ندارد که مؤلفه‌های f_1, f_2, \dots, f_k متعلق به یک خانواده از چگالی‌ها باشند. توابع چگالی f_i برای $i = 1, 2, \dots, k$ می‌توانند شامل بردار پارامتر $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{ih})$ باشند. بنابراین یک چگالی آمیخته متناهی پارامتری با k مؤلفه به شکل کلی زیر است:

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

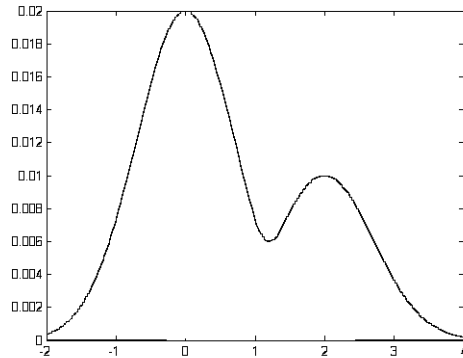
که در آن $\psi = (\pi_1, \pi_2, \dots, \pi_k, \theta'_1, \theta'_2, \dots, \theta'_k)' \in \Psi$ و ψ را بردار پارامتر مدل آمیخته می‌نامیم که در فضای پارامتر Ψ تغییر می‌کند.

۳ نمونه‌گیری از یک جامعه ناهمگن

فرض کنید X_1, X_2, \dots, X_n نمونه‌ای تصادفی به اندازه n از جامعه ناهمگن P باشد. واضح است هر یک از نمونه‌های X_j با احتمال π_j متعلق به زیر جامعه P_i است. بنابراین برای هر مشاهده X_j داریم:

$$f(x_j|\psi) = \sum_{i=1}^k \pi_i f(x_j|\theta_i), \theta_i \in \Theta, x \in S \quad (4)$$

برای مشخص نمودن اینکه نمونه X_j متعلق به زیر جامعه P_i می‌باشد، از یک متغیر تصادفی کمکی استفاده می‌شود که به آن متغیر پنهان^۴ می‌گوییم.



شکل ۱. نمودار تابع چگالی آمیخته با دو نما.

۲ معرفی توزیع آمیخته

فرض کنید جامعه آماری P شامل k ($k \geq 2$) زیر جامعه P_1, P_2, \dots, P_k باشد. هدف بررسی خصوصیتی در این جامعه است. همچنین فرض کنید متغیر تصادفی X نشان دهنده مقدار خصوصیت مذکور باشد. اگر F_1, F_2, \dots, F_k تابع توزیع متغیر تصادفی X به ترتیب در زیر جامعه‌های P_1, P_2, \dots, P_k باشند، آنگاه تابع توزیع متغیر تصادفی X با تکیه گاه S در جامعه P به صورت زیر خواهد بود:

$$F(x) = \pi_1 F_1(x) + \pi_2 F_2(x) + \dots + \pi_k F_k(x), x \in S \quad (1)$$

که در آن احتمال آن است که X دارای توزیع F_i باشد. $0 < \pi_i < 1$ و برای $i = 1, 2, \dots, k$ و $\sum_{i=1}^k \pi_i = 1$ و مقادیر π_i ها را نسبت‌های آمیختگی می‌نامند. مدل (۱) یک توزیع آمیخته متناهی با k مؤلفه بر حسب تابع توزیع است. چنانچه f_1, f_2, \dots, f_k توابع چگالی متغیر

^۴ Latent Variable

در نتیجه برای نمونه تصادفی X_1, X_2, \dots, X_n تابع درست‌نمایی در داده‌های کامل به صورت زیر خواهد بود:

$$L(\psi|y_1, y_2, \dots, y_n) = L(\psi|\underline{x}, \underline{z}) = \prod_{j=1}^n f(x_j, z_j|\psi) \quad (11)$$

و با به‌کارگیری (۱۰) داریم:

$$L(\psi|y_1, y_2, \dots, y_n) = \prod_{j=1}^n \prod_{i=1}^k [\pi_i f(x_j|\theta_i)]^{z_{ij}} \quad (12)$$

لگاریتم $L(\psi)$ به صورت زیر نوشته می‌شود:

$$\begin{aligned} \log L(\psi|y_1, y_2, \dots, y_n) &= l(\psi|y_1, y_2, \dots, y_n) \\ &= \sum_{j=1}^n \sum_{i=1}^k z_{ij} \log \pi_i + \sum_{j=1}^n \sum_{i=1}^k z_{ij} \log f(x_j|\theta_i) \quad (13) \end{aligned}$$

یکی از مسائل اصلی در توزیع‌های آمیخته، مسئله برآورد پارامترها می‌باشد. معمولاً برآوردهای گشتاوری، برآوردهای دقیقی نبوده و در بسیاری از مواقع این برآوردها بسیار نادرست و غیر واقعی هستند. اگر از (۱۳) نسبت به پارامترهای موجود مشتق گرفته و مساوی صفر قرار دهیم، می‌توانیم برآوردهای ماکسیمم درست‌نمایی این پارامترها را به دست آوریم. از آنجا که حل معادلات بوجود آمده توسط مشتق گیری در توزیعهای آمیخته کاری دشوار است بنابراین باید برای به دست آوردن برآورد پارامترها از روشهای عددی استفاده کنیم. از معروفترین روشهای عددی موجود در به دست آوردن برآوردهای ماکسیمم درست‌نمایی پارامترها می‌توان الگوریتم EM را نام برد که توسط دمپستر [۲] معرفی گردید.

اگر این متغیرها را با Z_j نشان دهیم، آنگاه برای نمونه: $Z_j = (Z_{1j}, Z_{2j}, \dots, Z_{kj})$ بردار X_1, X_2, \dots, X_n به صورت زیر تعریف می‌شود:

$$Z_{ij} = \begin{cases} 1 & \text{اگر } X_j \text{ متعلق به زیر جامعه } P_i \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases} \quad (5)$$

با به‌کارگیری متغیرهای پنهان Z_j ، داده‌های کامل به صورت زیر تعریف می‌شود:

$$\{Y_j; j = 1, 2, \dots, n\} = \{(X_j, Z_j); j = 1, 2, \dots, n\} \quad (6)$$

واضح است که $\sum_{j=1}^n Z_{ij} = 1$ برای $i = 1, 2, \dots, k$ داریم:

$$\sum_{i=1}^k z_{ij} f(x_j|\theta_i) = \prod_{i=1}^k f(x_j|\theta_i)^{z_{ij}} \quad (7)$$

از سوی دیگر $(Z_j|\underline{\theta}, \underline{\pi})$ برای $j = 1, 2, \dots, k$ از هم مستقل بوده و دارای توزیع چند جمله‌ای به صورت زیر است.

$$(Z_j|\underline{\theta}, \underline{\pi}) \sim M(1, \pi_1, \dots, \pi_k)$$

بنابراین داریم:

$$f(z_j|\underline{\theta}, \underline{\pi}) = \frac{1!}{z_{1j}! z_{2j}! \dots z_{kj}!} \prod_{i=1}^k \pi_i^{z_{ij}} \sim \prod_{i=1}^k \pi_i^{z_{ij}} \quad (8)$$

همچنین طبق (۷) داریم:

$$f(x_j|z_j, \underline{\theta}, \underline{\pi}) = \sum_{i=1}^k z_{ij} f(x_j|\theta_i) = \prod_{i=1}^k f(x_j|\theta_i)^{z_{ij}} \quad (9)$$

بنابراین می‌توان نوشت:

$$\begin{aligned} f(x_j, z_j|\underline{\theta}, \underline{\pi}) &= f(x_j|z_j, \underline{\theta}, \underline{\pi}) f(z_j|\underline{\theta}, \underline{\pi}) \\ &= \prod_{i=1}^k [\pi_i f(x_j|\theta_i)]^{z_{ij}} \quad (10) \end{aligned}$$

۴ توزیع نرمال چوله

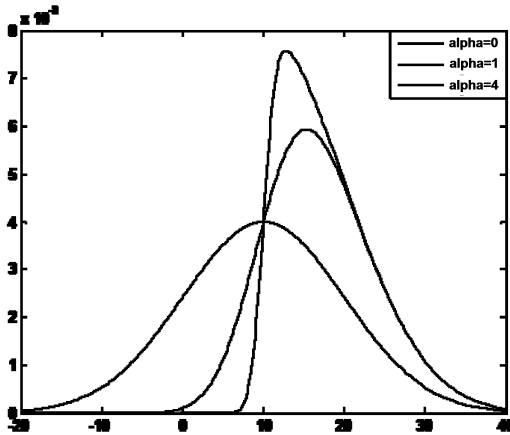
فرض کنید X یک متغیر تصادفی با مجموعه مقادیر \mathbb{R} باشد، اگر تابع چگالی این متغیر تصادفی به صورت:

$$f(x|\mu, \sigma^2, \alpha) = \frac{2}{\sigma} \phi\left(\frac{x-\mu}{\sigma}\right) \Phi\left(\alpha \frac{x-\mu}{\sigma}\right) \quad (14)$$

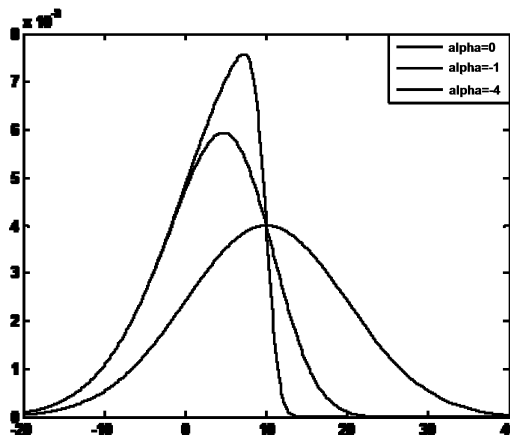
باشد، در این صورت می‌گوییم X دارای توزیع نرمال چوله^۶ با پارامترهای μ ، σ^2 و α بوده و آن را با نماد: $X \sim SN(\mu, \sigma^2, \alpha)$ نشان می‌دهیم. این توزیع اولین بار توسط آزالینی [۱] معرفی گردید. در (۱۴)، Φ و ϕ به ترتیب تابع چگالی و توزیع نرمال استاندارد می‌باشند. α را پارامتر شکل (Shape Parameter) می‌نامیم. به عبارت دیگر تابع چگالی نرمال چوله را می‌توان به صورت زیر نشان داد.

$$f(x|\mu, \sigma^2, \alpha) = \frac{2}{\sigma} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} \times \int_{-\infty}^{\alpha\left(\frac{x-\mu}{\sigma}\right)} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz$$

باید دقت داشته باشیم هرگاه: $\alpha = 0$ باشد در این صورت: $f(x|\mu, \sigma^2, \alpha = 0) = \phi\left(\frac{x-\mu}{\sigma}\right)$ که نشان دهنده توزیع نرمال با پارامترهای μ و σ^2 است. همچنین اگر α مثبت باشد چولگی به سمت راست و در صورتی که α منفی باشد چولگی نمودار تابع چگالی به سمت چپ خواهد بود. چنانچه $\alpha \rightarrow \infty$ آن‌گاه تابع چگالی به سمت تابع چگالی نیم-نرمال (Half-Normal) میل خواهد کرد. نمودارهای زیر وضعیت چولگی تابع چگالی نرمال چوله را به ازای مقادیر مختلف نشان می‌دهد.



شکل ۲. نمودار تابع چگالی نرمال چوله به ازای مقادیر مثبت $(\mu = 10, \sigma = 10)$



شکل ۳. نمودار تابع چگالی نرمال چوله به ازای مقادیر منفی $(\mu = 10, \sigma = 10)$

چنانچه در (۱۴) $\mu = 0$ و $\sigma^2 = 1$ باشد، در این صورت تابع چگالی نرمال چوله استاندارد را خواهیم داشت. و در این حالت تنها پارامتر موجود، پارامتر چولگی بوده و تابع چگالی آن را به صورت: $f(z|\alpha) = 2\phi(z)\Phi(\alpha z)$ نشان می‌دهیم. پیسی [۴] در مورد توزیع نرمال معرفی شده توسط آزالینی [۱] مسائلی را حل نمود. رینالدو [۳] نیز

قضیه ۴-۱ فرض کنید متغیر تصادفی U دارای توزیع نرمال با میانگین صفر و واریانس σ^2 و متغیر تصادفی T دارای توزیع بریده شده نرمال به صورت: $\{t > 0\} I \{t > 0\} \sim TN(0, \sigma^2)$ باشد. اگر T و U مستقل باشند، آنگاه متغیر تصادفی: $X = \mu + \frac{\alpha}{\sqrt{1+\alpha^2}}T + \frac{1}{\sqrt{1+\alpha^2}}U$ دارای توزیع نرمال چوله با پارامترهای μ ، σ^2 و α خواهد بود. اثبات: با

استفاده از تبدیل زیر می‌توان قضیه را اثبات نمود:

$$\begin{cases} X = \mu + \frac{\alpha}{\sqrt{1+\alpha^2}}T + \frac{1}{\sqrt{1+\alpha^2}}U \\ W = T \end{cases}$$

نتیجه ۴-۱ با توجه به روابط فوق داریم:

$$T|X = x \sim TN(\delta(\alpha)(x - \mu), \sigma^2(1 - \delta^2(\alpha))) I \{t > 0\} \quad (15)$$

فرع ۴-۱ فرض کنید متغیر تصادفی Y دارای توزیع بریده شده با پارامترهای μ و σ^2 در فاصله (a_1, a_2) باشد. در اینصورت می‌توان نشان داد:

$$E(Y) = \mu - \sigma \frac{\phi(\frac{a_2 - \mu}{\sigma}) - \phi(\frac{a_1 - \mu}{\sigma})}{\Phi(\frac{a_2 - \mu}{\sigma}) - \Phi(\frac{a_1 - \mu}{\sigma})}$$

و

$$E(Y^2) = \mu^2 + \sigma^2 - 2\mu\sigma \frac{\phi(\frac{a_2 - \mu}{\sigma}) - \phi(\frac{a_1 - \mu}{\sigma})}{\Phi(\frac{a_2 - \mu}{\sigma}) - \Phi(\frac{a_1 - \mu}{\sigma})} - \sigma^2 \frac{(\frac{a_2 - \mu}{\sigma})\phi(\frac{a_2 - \mu}{\sigma}) - (\frac{a_1 - \mu}{\sigma})\phi(\frac{a_1 - \mu}{\sigma})}{\Phi(\frac{a_2 - \mu}{\sigma}) - \Phi(\frac{a_1 - \mu}{\sigma})}$$

بنابراین با استفاده از (۱۵) و فرع (۴-۱) داریم:

$$E(T|X = x) = \delta(\alpha)(x - \mu) + \frac{\phi(\frac{\delta(\alpha)(x - \mu)}{\sigma\sqrt{1 - \delta^2(\alpha)}})}{\Phi(\frac{\delta(\alpha)(x - \mu)}{\sigma\sqrt{1 - \delta^2(\alpha)}})} \times \sigma\sqrt{1 - \delta^2(\alpha)} \quad (16)$$

یک حالت کلی تر از توزیع نرمال چوله را که به آن توزیع نرمال چوله تعمیم یافته^۷ گفته می‌شود معرفی کرد که تابع چگالی آن به صورت زیر تعریف می‌شود.

$$f(x|\alpha_1, \alpha_2) = 2\phi(x)\Phi\left(\frac{\alpha_1 x}{\sqrt{1 + \alpha_2 x^2}}\right); x \in R, \alpha_1 \in R, \alpha_2 \geq 0$$

لازم به ذکر است که تابع چگالی فوق مربوط به توزیع نرمال استاندارد چوله تعمیم یافته با پارامترهای α_1 و α_2 بوده و آن را با نماد: $X \sim SGN(0, 1, \alpha_1, \alpha_2)$ نشان می‌دهیم. با تبدیل: $Y = \mu + \sigma X$ شکل کلی تابع چگالی نرمال چوله تعمیم یافته به صورت زیر به دست می‌آید.

$$f(y|\mu, \sigma, \alpha_1, \alpha_2) = \frac{2}{\sigma}\phi\left(\frac{y - \mu}{\sigma}\right)\Phi\left\{\frac{\alpha_1(y - \mu)}{\sqrt{\sigma^2 + \alpha_2(y - \mu)^2}}\right\}$$

برای به دست آوردن تابع چگالی نرمال چوله می‌توانیم از قضیه (۱) استفاده کنیم. اما ابتدا تعریف زیر را خواهیم داشت.

تعریف ۴-۱ فرض کنید Y یک متغیر تصادفی باشد. می‌گوییم این متغیر تصادفی دارای توزیع نرمال بریده شده^۸ در فاصله (a_1, a_2) می‌باشد و آن را با نماد $Y \sim TN(\mu, \sigma^2) I \{a_1 < y < a_2\}$ نشان می‌دهیم، هرگاه تابع چگالی آن به صورت زیر باشد.

$$f(y|\mu, \sigma^2) = \left\{ \Phi\left(\frac{a_2 - \mu}{\sigma}\right) - \Phi\left(\frac{a_1 - \mu}{\sigma}\right) \right\}^{-1} \times \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}; a_1 < y < a_2$$

$$; j = 1, 2, \dots, n \quad (19)$$

$$T_j | z_{ij} = 1 \sim TN(0, \sigma_i^2) I\{t_j > 0\} \quad (20)$$

$$Z_j \sim M(\lambda, \pi_1, \pi_2, \dots, \pi_k); j = 1, 2, \dots, n \quad (21)$$

که $\sigma_{T_i} = \sigma_i \sqrt{1 - \delta^2(\alpha_i)}$ و $\mu_{T_{ij}} = \delta(\alpha_i)(x_j - \mu_i)$ اگر داده‌های کامل را $(\underline{X}, \underline{T}, \underline{Z})$ در نظر بگیریم، در این صورت مدل سلسله مراتبی^۹ برای داده‌های کامل به صورت زیر است:

$$f(x_j, t_j, z_j | \psi) = f(x_j, t_j | z_j, \psi) f(z_j | \underline{\pi}) \quad (22)$$

که در (۲۲) داریم:

$$\begin{aligned} f(x_j, t_j | z_{ij}, \psi) &= \prod_{i=1}^k z_{ij} f(x_j, t_j | \psi) \\ &= \prod_{i=1}^k [f(x_j, t_j | \psi)]^{z_{ij}} \end{aligned}$$

و یا

$$\begin{aligned} f(x_j, t_j | z_{ij}, \psi) &= \prod_{i=1}^k \left[\frac{1}{\pi \sigma_i^2 \sqrt{1 - \delta^2(\alpha_i)}} \right. \\ &\times \exp\left\{ -\frac{1}{2\sigma_i^2(1 - \delta^2(\alpha_i))} [(x_j - \mu_i) - 2\delta(\alpha_i) \right. \\ &\times (x_j - \mu_i)t_j + t_j^2] \left. \right\} \left. \right]^{z_{ij}} \end{aligned}$$

بنابراین می‌توانیم تابع درستنمایی را برای داده‌های کامل به فرم زیر بنویسیم:

$$\begin{aligned} L(\psi | \underline{x}, \underline{t}, \underline{z}) &= \prod_{j=1}^n \prod_{i=1}^k \left[\frac{1}{\pi \sigma_i^2 \sqrt{1 - \delta^2(\alpha_i)}} \right. \\ &\times \exp\left\{ -\frac{1}{2\sigma_i^2(1 - \delta^2(\alpha_i))} [(x_j - \mu_i) \right. \\ &\left. - 2\delta(\alpha_i)(x_j - \mu_i)t_j + t_j^2] \right\} \left. \right]^{z_{ij}} \prod_{i=1}^n \prod_{i=1}^k \pi_i^{z_{ij}} \end{aligned}$$

$$\begin{aligned} E(T^2 | X = x) &= \delta^2(\alpha)(x - \mu)^2 + \sigma^2(1 - \delta^2(\alpha)) \\ &+ \frac{\phi\left(\frac{\delta(\alpha)(x - \mu)}{\sigma\sqrt{1 - \delta^2(\alpha)}}\right)}{\Phi\left(\frac{\delta(\alpha)(x - \mu)}{\sigma\sqrt{1 - \delta^2(\alpha)}}\right)} \delta(\alpha)(x - \mu)\sigma^2(1 - \delta^2(\alpha)) \quad (17) \end{aligned}$$

در اینجا نیز باید دقت داشته باشیم که برآورد ماکسیمم درستنمایی برای پارامترها در توزیع نرمال چوله به صورت دقیق میسر نبوده و باید از روشهای عددی مانند الگوریتم EM استفاده کنیم.

۵ تابع چگالی آمیخته با مولفه‌های نرمال چوله

فرض کنید یک متغیر تصادفی با تابع چگالی آمیخته زیر باشد:

$$f(x | \psi) = \sum_{i=1}^k \pi_i g(x | \mu_i, \sigma_i^2, \alpha_i); x \in R \quad (18)$$

به طوری که $g(x | \mu_i, \sigma_i^2, \alpha_i)$ نشان دهنده تابع چگالی نرمال چوله با پارامترهای μ_i, σ_i^2 و پارامتر شکل α_i است. بردار پارامتر ψ به صورت: $\psi = (\theta_1, \theta_2, \dots, \theta_k, \pi_1, \pi_2, \dots, \pi_k)$ در آن: $\theta_i = (\mu_i, \sigma_i^2, \alpha_i)$ می‌باشد. فرض کنید نمونه تصادفی X_1, X_2, \dots, X_n از x انتخاب و بردار متغیرهای پنهان باشد. $Z_j = (Z_{1j}, Z_{2j}, \dots, Z_{kj})$ همچنین فرض کنید T_1, T_2, \dots, T_n نیز یک نمونه تصادفی از توزیع نرمال بریده شده در فاصله $(0, +\infty)$ باشد. بنابراین طبق (۱۵) داریم:

$$T_j | x_j, z_{ij} = 1 \sim TN(\mu_{T_{ij}}, \sigma_{T_i}^2) I\{t_j > 0\}$$

لگاریتم تابع درست‌نمایی فوق را به دست می‌آوریم، می‌توان نوشت:

لگاریتم تابع درست‌نمایی فوق را به دست می‌آوریم، می‌توان نوشت:

$$l(\psi|\underline{x}, \underline{t}, \underline{z}) = \sum_{j=1}^n \sum_{i=1}^k Z_{ij} \{ \log(\pi_i) - \log(\pi_i \sigma_i^2) \} \\ - \frac{1}{2} \log(1 - \delta^2(\alpha_i)) \\ - \frac{t_j^2 - 2\delta(\alpha_i)(x_j - \mu_i) + (x_j - \mu_i)^2}{2\sigma_i^2(1 - \delta^2(\alpha_i))} \quad (۲۳)$$

$$\hat{Z}_{ij} = E_{\hat{\psi}}(Z_{ij}|\underline{x}) \\ = \frac{\hat{\pi}_i^{(k)} g(x_j|\hat{\mu}_i^{(k)}, \hat{\sigma}_i^{2(k)}, \hat{\alpha}_i^{(k)})}{\sum_{i=1}^k \hat{\pi}_i^{(k)} g(x_j|\hat{\mu}_i^{(k)}, \hat{\sigma}_i^{2(k)}, \hat{\alpha}_i^{(k)})} \quad (۲۴)$$

۶ برآوردهای ماکسیمم درست‌نمایی پارامترها

$$\hat{S}_{T_{ij}}^{(k)} = E_{\hat{\psi}}(Z_{ij} T_j^2 | \underline{x}) = \hat{Z}_{ij}^{(K)T} [\hat{\mu}_{T_{ij}}^{(K)T} + \hat{\sigma}_{T_i}^{2(K)}] \\ \times \frac{\phi\{\hat{\alpha}_i^{(k)}(\frac{x_j - \hat{\mu}_i^{(k)}}{\hat{\sigma}_i^{(k)}})\}}{\Phi\{\hat{\alpha}_i^{(k)}(\frac{x_j - \hat{\mu}_i^{(k)}}{\hat{\sigma}_i^{(k)}})\}} \hat{\mu}_{T_{ij}}^{(K)} \hat{\sigma}_{T_i}^{(K)} \quad (۲۵)$$

بنابراین با پیش بینی مقادیر Z_{ij} و T_j ها، اکنون می‌توانیم مقادیر برآورد پارامترها را در مرحله $(k+1)$ -ام به دست آورده و الگوریتم EM را ادامه داده تا به همگرایی برآورد پارامترها برسیم.

اگر مقدار پیش بینی \hat{Z}_{ij} و مقدار پیش بینی \hat{T}_j باشد، در این صورت برآورد ماکسیمم درست‌نمایی برای σ_i^2 ، μ_i ، π_i و α_i به صورت زیر خواهد بود.

$$\hat{\pi}_i = \frac{\sum_{j=1}^n \hat{Z}_{kj}}{n}; i = 1, 2, \dots, k \quad (۲۶)$$

$$\hat{\mu}_i = \frac{\sum_{j=1}^n \hat{Z}_{ij} X_j - \delta^{-1}(\hat{\alpha}_i) \sum_{j=1}^n \hat{Z}_{ij} \hat{T}_j}{\sum_{j=1}^n \hat{Z}_{ij}} \quad (۲۷)$$

$$\hat{\sigma}_i^2 = \frac{1}{2(1 - \delta^2(\hat{\alpha}_i))} \left\{ \sum_{j=1}^n \hat{Z}_{ij} \hat{T}_j^2 - 2\delta^{-1}(\hat{\alpha}_i) \sum_{j=1}^n \hat{Z}_{ij} \hat{T}_j (X_j - \hat{\mu}_i) + \sum_{j=1}^n \hat{Z}_{ij} (X_j - \hat{\mu}_i)^2 \right\} \quad (۲۸)$$

با مشتق‌گیری از لگاریتم تابع درست‌نمایی (۲۳) نسبت به α_i ، σ_i^2 ، μ_i ، π_i ها و مساوی صفر قرار دادن معادلات به دست آمده، داریم:

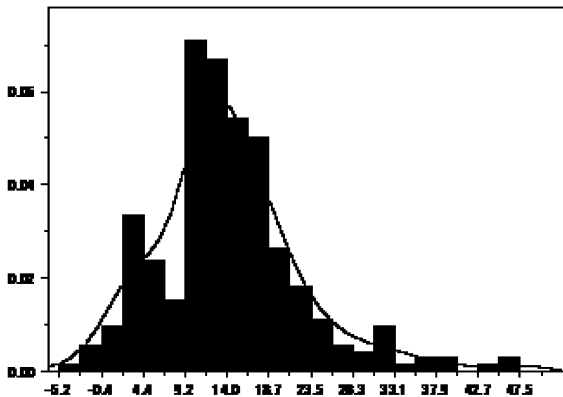
$$\mu_i = \frac{\sum_{j=1}^n Z_{ij} X_j - \delta^{-1}(\alpha_i) \sum_{j=1}^n Z_{ij} T_j}{\sum_{j=1}^n Z_{ij}}, \pi_i = \frac{\sum_{j=1}^n Z_{kj}}{n} \\ \sigma_i^2 = \frac{\sum_{j=1}^n Z_{ij} T_j^2 - 2\delta^{-1}(\alpha_i) \sum_{j=1}^n Z_{ij} T_j (X_j - \mu_i) + \sum_{j=1}^n Z_{ij} (X_j - \mu_i)^2}{2(1 - \delta^2(\alpha_i)) \sum_{j=1}^n Z_{ij}}$$

و

حل معادله آخر می‌تواند مقدار $\delta(\alpha_i)$ و در نتیجه مقدار α_i را در اختیار ما قرار دهد. همان‌گونه که مشاهده می‌شود در تمام معادلات فوق Z_{ij} و T_j ها ظاهر شده‌اند که عملاً مقادیر نامعلوم بوده و باید پیش بینی شوند. لذا برای پیش بینی آنها از الگوریتم EM استفاده می‌کنیم. با در نظر گرفتن مقادیر

همان‌گونه که مشاهده می‌شود در تمام معادلات فوق Z_{ij} و T_j ها ظاهر شده‌اند که عملاً مقادیر نامعلوم بوده و باید پیش بینی شوند. لذا برای پیش بینی آنها از الگوریتم EM استفاده می‌کنیم. با در نظر گرفتن مقادیر

آوریم. ابتدا هیستوگرام داده‌ها را رسم می‌کنیم. چولگی و آمیختگی داده‌ها کاملاً مشخص است.



شکل ۴. هیستوگرام داده‌های مربوط به چگالی آمیخته از دو توزیع نرمال چوله.

با نوشتن برنامه مربوط به الگوریتم EM در نرم افزار $MATLAB$ و اجرای این برنامه، برآورد پارامترها در جدول زیر به دست آمده است.

جدول ۱. مقادیر برآورد پارامترهای چگالی آمیخته از دو

چگالی نرمال چوله.

Parameter	True value	Estimation	Std.
μ_1	۵	۵/۳۶۶۲	۰/۱۵۸۶
μ_2	۱۰	۹/۹۶۴۱	۰/۰۳۰۵
σ_1	۵	۴/۸۲۳۰	۰/۰۸۴۹
σ_2	۱۰	۱۱/۱۷۶۹	۰/۹۹۳۱
α_1	۰/۵	۰/۵۲۸۴	۰/۰۷۹۷
α_2	۲	۱/۹۶۸۴	۰/۰۲۷۴
π	۰/۳	۰/۲۷۴۷	۰/۰۰۶۶

برای اطمینان از درست بودن برآوردهای به دست آمده نمودارهای همگرایی برآورد پارامترها را رسم می‌کنیم.

با استفاده از الگوریتم EM و جایگزینی $\hat{S}_{ij}^{(k)}$ و $\hat{S}_{\nu ij}^{(k)}$ به جای \hat{T}_j و \hat{T}_j^{ν} می‌توانیم برآورد پارامترها را در مرحله $(k+1)$ - ام به صورت زیر به دست آوریم.

$$\hat{\pi}_i^{(k+1)} = \frac{\sum_{j=1}^n \hat{Z}_{kj}^{(k)}}{n}; i = 1, 2, \dots, k \quad (29)$$

$$\hat{\mu}_i^{(k+1)} = \frac{\sum_{j=1}^n \hat{Z}_{ij}^{(k)} X_j - \delta^{-1}(\hat{\alpha}_i^{(k)}) \sum_{j=1}^n \hat{S}_{ij}^{(k)}}{\sum_{j=1}^n \hat{Z}_{ij}^{(k)}} \quad (30)$$

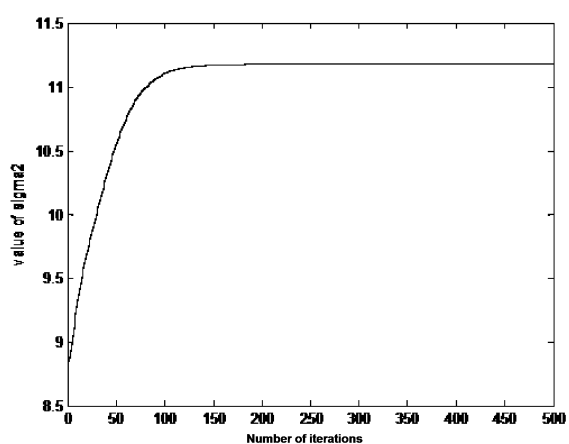
$$\hat{\sigma}_i^{\nu(k+1)} = \frac{1}{\nu(1 - \delta^{\nu}(\hat{\alpha}_i^{(k)}))} \left\{ \sum_{j=1}^n \hat{S}_{\nu ij}^{(k)} - \nu \delta^{-1}(\hat{\alpha}_i^{(k)}) \sum_{j=1}^n \hat{S}_{ij}^{(k)} (X_j - \hat{\mu}_i^{(k)}) + \sum_{j=1}^n \hat{Z}_{ij}^{(k)} (X_j - \hat{\mu}_i^{(k)})^{\nu} \right\} \quad (31)$$

با ثابت نگه داشتن $\mu_i = \hat{\mu}_i^{(k+1)}$ و $\sigma_i^{\nu} = \hat{\sigma}_i^{\nu(k+1)}$ در معادله زیر، $\hat{\alpha}_i^{(k+1)}$ را به عنوان برآورد پارامتر α_i ($i = 1, 2, \dots, k$) به دست می‌آوریم.

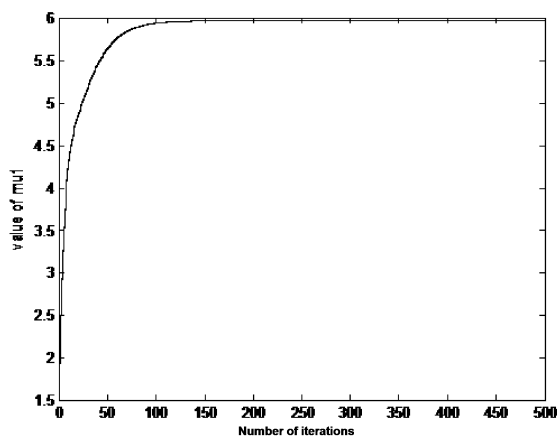
$$\sigma_i^{\nu} \delta(\alpha_i) (1 - \delta^{\nu}(\alpha_i)) \sum_{j=1}^n Z_{ij} + (1 + \delta^{\nu}(\alpha_i)) \times \sum_{j=1}^n Z_{ij} T_j (X_j - \mu_i) - \delta(\alpha_i) \sum_{j=1}^n Z_{ij} T_j^{\nu} - \delta(\alpha_i) \sum_{j=1}^n Z_{ij} (X_j - \mu_i)^{\nu} = 0$$

در اینجا مثالی را در مورد آمیخته دو چگالی نرمال چوله بررسی می‌کنیم.

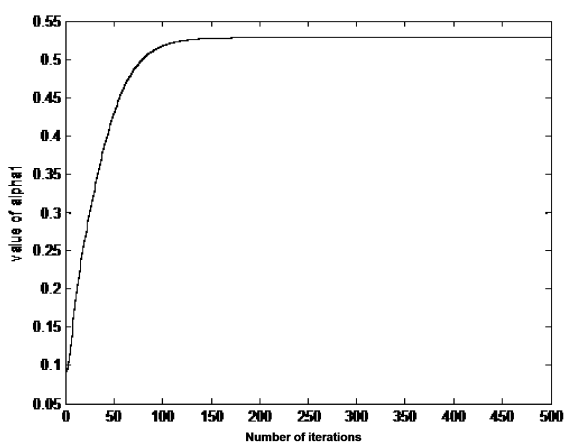
مثال ۶-۱ با به کارگیری قضیه (۱) تعداد ۳۰۰ داده را از دو توزیع نرمال چوله آمیخته با پارامترهای $\mu_1 = 5, \mu_2 = 10, \sigma_1 = 5, \sigma_2 = 10$ و پارامترهای شکل $\alpha_1 = 0.5$ و $\alpha_2 = 2$ با نسبتهای 0.3 و 0.7 تولید کرده‌ایم. می‌خواهیم برآورد پارامترها را به دست



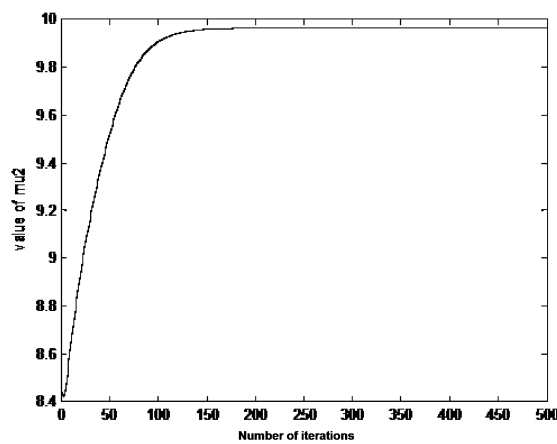
شکل ۸. نمودار همگرایی برآورد پارامتر σ_2 .



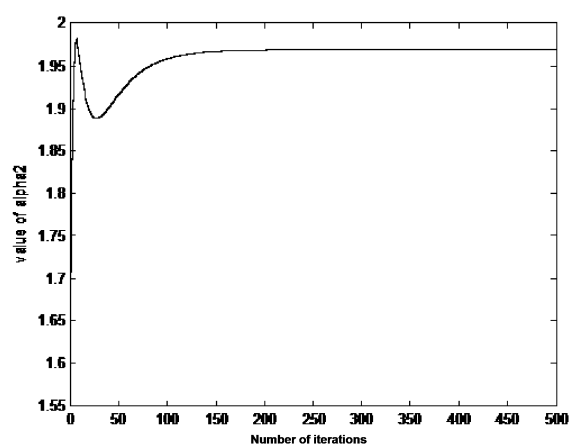
شکل ۵. نمودار همگرایی برآورد پارامتر μ_1 .



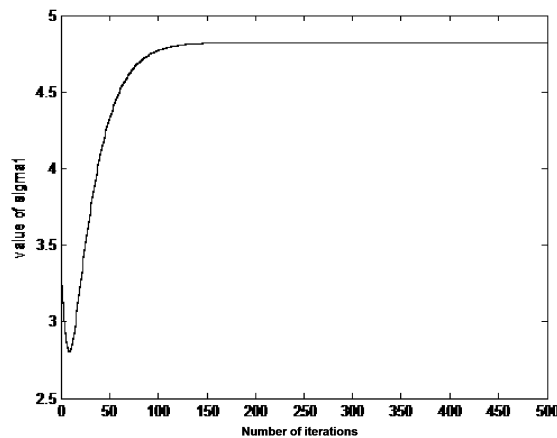
شکل ۹. نمودار همگرایی برآورد پارامتر α_1 .



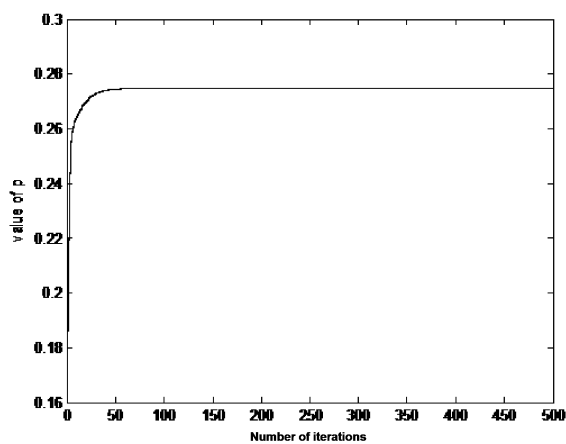
شکل ۶. نمودار همگرایی برآورد پارامتر μ_2 .



شکل ۱۰. نمودار همگرایی برآورد پارامتر α_2 .



شکل ۷. نمودار همگرایی برآورد پارامتر σ_1 .



شکل ۱۱. نمودار همگرایی برآورد پارامتر π .

مراجع

- [1] Azzalini, A. (1986). Further results on a class of distributions which includes the normal ones, *Statistica*, 46,199-208.
- [2] Dempster, A. P., Laird, N. M and Rubin D. (1977). Maximum Likelihood from incomplete Data via the EM algorithm, *Jornal of Royal Statistical Society*, B, 39, 1-38.
- [3] Reinaldo, B. (2003). *A New Class of Skew-Normal Distributions*, Technical Report, North Carolina State University, Institue of Statistics, 1-15.
- [4] Pewsey, A. (2001). Problems of inference for Azzalini's skew-normal distribution, *Journal of Applied Statistics*, 27, 859-870.