

## آمیخته‌ای از مدل‌های تحلیل عاملی با مؤلفه‌های توزیع $t$ -چندمتغیره

مرضیه مهریار - مجتبی خزائی

دانشکده علوم ریاضی، دانشگاه شهید بهشتی

**چکیده:** استفاده از مدل تحلیل‌گر عاملی آمیخته در طی سال‌های اخیر در تحلیل داده‌ها رواج یافته است. این مدل که در ابتدا مبتنی بر فرض نرمال بودن مؤلفه‌ها معرفی گردید، برای حصول نتایج مطلوب در حضور مشاهدات دور افتاده به حالتی که مؤلفه‌ها دارای توزیع  $t$ -چندمتغیره هستند نیز تعمیم داده شده. در این مقاله به معرفی مدل تحلیل‌گر عاملی آمیخته با مؤلفه‌های  $t$ -چندمتغیره به تاسی از کار مک‌لاکلان و همکاران (۲۰۰۷b) پرداخته و کاربرد آن را در تحلیل داده‌های ریزآرایه مربوط به سرطان روده بررسی خواهیم کرد.

**واژه‌های کلیدی:** مدل تحلیل‌گر عاملی، مدل تحلیل‌گر عاملی آمیخته با مؤلفه‌های  $t$ -چندمتغیره، توزیع  $t$ -چندمتغیره، الگوریتم AECM، ریزآرایه

### ۱ مقدمه

تحلیل عاملی یکی از روش‌های آماری است، که چنانچه هدف تحقیق خلاصه کردن مقدار زیاد اطلاعات و تقلیل آن‌ها به حداقل تعداد عامل‌های مشترک باشد، کاربرد دارد. در مدل تحلیل عاملی معمول فرض بر این است که مشاهده‌ها به یک جامعه همگن تعلق دارند. حال آن‌که این فرض در بسیاری از کاربردها برقرار نبوده و مجموعه مشاهده‌های تحت بررسی را می‌توان تلفیقی از مشاهده‌های چند زیرجامعه یا چند گروه دانست. قهرمانی و هیئتون (۱۹۹۷) مدلی تحت عنوان مدل تحلیل‌گر عاملی آمیخته ارائه کردند که در شرایط فوق به کار می‌آید، این مدل توسط مک‌لاکلان و پیل (۲۰۰۰) گسترش یافت. در این مدل‌ها فرض می‌شود توزیع زیرجامعه‌ها، نرمال چندمتغیره است. این فرض باعث گردیده در مسائل با داده‌های دور افتاده در مجموعه مشاهده‌ها، نتایج به دست آمده چندان قابل اعتماد نباشد، به طوری که به نظر می‌رسد مدل‌هایی که در آن‌ها از توزیع‌های دم‌کلفت استفاده می‌کنند بهتر قادرند رفتار مشاهده‌های تحت بررسی را تبیین کنند. برای این منظور مدل تحلیل‌گر عاملی آمیخته  $t$  توسط مک‌لاکلان و همکاران (۲۰۰۷b) معرفی شد. آن‌ها این مدل را در تحلیل داده‌های ریزآرایه سرطان

سینه به کار بردند. در تکمیل کار این افراد در این مقاله، کار با مجموعه کل ژن‌های موجود در ریزآرایه شروع می‌شود و با انجام آزمون‌های آماری، کاهش ژن انجام می‌شود و نهایتاً از مدل تحلیل‌گر عاملی آمیخته  $t$  جهت شناسایی خوشه‌ای از ژن‌ها که بیشترین قدرت تفکیک‌پذیری را میان بافت‌های سالم و سرطانی ایجاد می‌نماید، استفاده می‌شود. برای این منظور، الگوریتم  $EMMIX - GENE$  (الگوریتمی که از مدل تحلیل‌گر عاملی آمیخته نرمال برای تحلیل‌ها استفاده کرده) مورد استفاده قرار گرفته است.

## ۲ توزیع $t$ ی چند متغیره

گوییم بردار تصادفی  $p$  بعدی  $\mathbf{Y}$  دارای توزیع  $t$ ی  $p$  متغیره با بردار مکانی  $\boldsymbol{\mu}$ ، ماتریس مقیاس  $\boldsymbol{\Sigma}$  و درجه آزادی  $v$  است و می‌نویسیم  $\mathbf{Y} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, v)$ ، هرگاه چگالی آن به صورت زیر باشد:

$$f_p(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, v) = \frac{\Gamma(\frac{v+p}{2})|\boldsymbol{\Sigma}|^{-\frac{1}{2}}}{(\pi v)^{\frac{1}{2}p}\Gamma(\frac{v}{2})\{1 + \frac{1}{v}\delta(\mathbf{y}, \boldsymbol{\mu}; \boldsymbol{\Sigma})\}^{\frac{1}{2}(v+p)}}, \quad \mathbf{y} \in BR^p$$

که  $\delta(\mathbf{y}, \boldsymbol{\mu}; \boldsymbol{\Sigma}) = (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})$  مربع فاصله مایلانویس میان  $\mathbf{y}$  و  $\boldsymbol{\mu}$  است. توزیع  $t$ ی چند متغیره را می‌توان به صورت آمیخته‌ای (نامتناهی) از توزیع نرمال و گاما نوشت. اثبات می‌شود اگر  $\mathbf{Y}|W = w \sim N_p(\boldsymbol{\mu}, \frac{1}{w}\boldsymbol{\Sigma})$  و  $W \sim \gamma(\frac{v}{2}, \frac{v}{2})$ ، آن‌گاه  $\mathbf{Y} \sim t_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, v)$ ، که در آن،  $Y$  مشاهده می‌شود ولی  $W$  قابل مشاهده نیست. از این رو، در برآورد ماکسیمم درست‌نمایی پارامترهای توزیع  $t$ ی چند متغیره در حضور مقادیر گمشده  $w$ ها از الگوریتم EM استفاده می‌شود. در این حالت به سادگی می‌توان نشان داد که

$$w_j^{(k)} = E_{\Psi^{(k)}}(W_j | \mathbf{y}_j) = \frac{v+p}{v + \delta(\mathbf{y}_j, \boldsymbol{\mu}^{(k)}; \boldsymbol{\Sigma}^{(k)})}, \quad (1)$$

رابطه (۱) نشان می‌دهد  $w_j^{(k)}$ ها یا وزن‌های کوچک متناظر با  $\delta(\mathbf{y}_j, \boldsymbol{\mu}^{(k)}; \boldsymbol{\Sigma}^{(k)})$  بزرگ بوده و بیان‌گر این است که مشاهده  $\mathbf{y}_j$  از  $\boldsymbol{\mu}^{(k)}$  دور است. به عبارت دیگر نشان از دور افتاده بودن مشاهده  $\mathbf{y}_j$  دارد.

## ۳ مدل تحلیل‌گر عاملی آمیخته $t$

فرض کنید که  $y_1, \dots, y_n$  بردارهای تصادفی  $p$  بعدی از جامعه  $\Pi$  باشند که  $\Pi$  نیز ترکیبی از  $g \geq 2$  زیر جامعه  $\Pi_1, \dots, \Pi_g$  است و  $f_1, \dots, f_g$  تابع چگالی توزیع  $t$ -چند

متغیره بردار تصادفی  $y$ ، در زیرجامعه‌های  $\Pi_1, \dots, \Pi_g$  است. مدل تحلیل گره‌های عاملی آمیخته با مؤلفه‌های  $t$ -چندمتغیره به شکل زیر بیان می‌شود:

$$f(y_j; \psi) = \sum_{i=1}^g \pi_i f(y_j; \mu_i, \Sigma_i, \nu_i), \quad (2)$$

که در آن

$$\Sigma_i = B_i B_i' + D_i,$$

و بردار پارامتری  $\Psi$  متشکل از درجه‌های آزادی مؤلفه‌ها،  $\nu_i$ ها، نسبت‌های آمیختگی،  $\pi_i$ ها، عناصر  $\mu_i$ ،  $B_i$  و  $D_i$  می‌باشد.

مدل تحلیل گره عاملی آمیخته، آمیخته‌ی متناهی از مدل‌های خطی است و در کل مدلی غیرخطی در نظر گرفته می‌شود. این مدل علاوه بر این که یک شیوه کاهش بعد موضعی خطی، و در کل ناخطی است، در خوشه‌بندی داده‌ها شامل ساختارهای گروه‌بندی نیز مفید است.

مشابه با مدل تحلیل عاملی معمول، مدل آمیخته (۲) را به شکل (۳) می‌توان تعریف کرد. برای این منظور فرض می‌کنیم  $Z_j$  یک بردار تصادفی  $g$  بعدی باشد، که  $\pi_i$ ها را مؤلفه آن  $Z_{ij}$ ، به صورت زیر تعریف می‌شود:

$$Z_{ij} = \begin{cases} 1 & y_j \in \Pi_i \\ 0 & y_j \notin \Pi_i \end{cases}$$

و احتمال اینکه  $Z_{ij}$  مقدار ۱ را اختیار کند برابر است با  $\pi_i$ . بر این اساس توزیع  $Z_j$ ، توزیع چندجمله‌ای زیر است:

$$Z_j \sim MB_g(1, (\pi_1, \dots, \pi_g)').$$

تحت شرایط فوق، به طور معادل می‌توان گفت که مشاهده  $y_j$  با احتمال  $\pi_i$  به صورت زیر است:

$$y_j = \mu_i + B_i U_{ij} + e_{ij} \quad i = 1, \dots, g \quad (3)$$

مشروط بر این که

$$U_{ij} | z_{ij} = 1 \sim t_q(0, I_q, \nu_i),$$

$$e_{ij} | z_{ij} = 1 \sim t_q(0, D_q, \nu_i),$$

$$\text{cov}[(e_{ij}, U_{ij}) | z_{ij} = 1] = 0.$$

به عبارت دیگر، توزیع خطاها و عامل‌ها، توزیع  $t$ -چندمتغیره‌ای با درجه‌های آزادی یکسان و ناهمبسته‌اند.

## ۴ برآورد پارامترهای مدل تحلیل گر عاملی

در این بخش با چگونگی برآورد ماکسیمم درستنمایی پارامترهای مدل تحلیل گر عاملی آمیخته  $t$  آشنا می شویم، با توجه به نامعلوم بودن  $Z_{ij}$  ها و  $U_{ij}$  ها (و البته  $W_j$  ها) می توان از الگوریتم EM استفاده کرد، که الگوریتم EM یکی از ابزارهای متداول برای یافتن برآورد ماکسیمم درستنمایی، در حضور مقادیر گمشده است. در استفاده از الگوریتم EM، ایده آل این است که نقطه ماکسیمم در گام  $M$ ، به سادگی و توسط روابط بسته ای قابل حصول باشند. اما متأسفانه در بعضی از موقعیت ها این گونه نیست و این نقاط را باید به روش های عددی محاسبه کرد که در جای خود می توانند پیچیده بوده و به خصوص سرعت همگرایی را پایین آورند. در مواردی که اطلاع از مقدار بعضی از پارامترها کار ماکسیمم سازی پارامترهای دیگر را بسیار ساده ترمی کند، استفاده از الگوریتم ECM توسط منگ و روبین (۱۹۹۳) پیشنهاد شد. در مواردی مثل این جا که علاوه بر این، مقادیر گمشده از انواع مختلفی است، روش  $AECM$ <sup>۱</sup> که توسط منگ و ون دیک (۱۹۹۷) ارائه شد، برای محاسبه برآورد پارامترهای مدل تحلیل گر عاملی  $t$  استفاده می شود.

در ادامه با جزئیات الگوریتم  $AECM$  برای برآورد بردار پارامترهای نامعلوم،  $\Psi$ ، آشنا می شویم. در ابتدا بردار  $\Psi$  به دو بخش  $(\Psi'_1, \Psi'_2)$  افزایش می شود که در آن  $\Psi_1$  شامل عناصر  $\mu_i$  ها، نسبت های آمیختگی  $\pi_i$  ها و درجه های آزادی  $\nu_i$  ها و  $\Psi_2$  شامل عناصر  $B_i$  ها و  $D_i$  ها است. هر تکرار الگوریتم  $AECM$  از دو چرخه تشکیل شده است، که هر چرخه شامل دو مرحله  $E$  و  $M$  است:

### چرخه اول:

مرحله  $E$ : امید ریاضی لگاریتم درستنمایی داده های کامل، به شرط داده های مشاهده شده، به ازای برآورد پارامترها،  $\Psi^{(k)}$ ، در تکرار  $(k+1)$  ام به صورت زیر تعریف می شود:

$$Q_1(\theta; \theta^{(k)}) = \sum_{j=1}^n \sum_{i=1}^g \ln[\pi_i f_p(\mathbf{y}_j; \mu_i, \Sigma_i, \nu_i)]_{\Psi^{(k)}}(Z_{ij} | \mathbf{y}_j) \quad (4)$$

برای محاسبه امید شرطی (۴) کافیهست  $E_{\Psi^{(k)}}(Z_{ij} | \mathbf{y}_j)$  محاسبه شود.

$$E_{\Psi^{(k)}}(Z_{ij} | \mathbf{y}) = \tau_{ij}^{(k)} = \frac{\pi_i^{(k)} f(\mathbf{y}_j; \mu_i^{(k)}, \Sigma_i^{(k)}, \nu_i^{(k)})}{f(\mathbf{y}_j; \Psi^{(k)})}$$

<sup>۱</sup> Alternating Expectation Conditional Maximization

مرحله  $M$ : شامل ماکسیمم کردن تابع  $Q_1(\Psi; \Psi^{(k)})$  نسبت به پارامترهای  $\mu_i, \pi_i$  و  $\nu_i$  است. برآورد  $\mu_i$  و  $\pi_i$  به صورت زیر برآورد خواهد شد.

$$\pi_i^{(k+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(k)}}{n}, \quad \mu_i^{(k+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(k)} w_{ij}^{(k)} y_j}{\sum_{j=1}^n \tau_{ij}^{(k)} w_{ij}^{(k)}}$$

برای تعیین برآورد  $\nu_i$  نیاز به حل معادله زیر است.

$$\left\{ \psi\left(\frac{\nu_i}{\gamma}\right) + \log \frac{\nu_i}{\gamma} + 1 + \frac{1}{n_i^{(k)}} \sum_{j=1}^n \tau_{ij}^{(k)} (\log w_{ij}^{(k)} - w_{ij}^{(k)}) + \psi\left(\frac{\nu_i^{(k)} + p}{\gamma}\right) + \log\left(\frac{\nu_i^{(k)} + p}{\gamma}\right) \right\} = 0 \quad (5)$$

که در آن  $n_i^{(k)} = \sum_{j=1}^n \tau_{ij}^{(k)}$  معادله (5) جواب صریحی نداشته و باید با به کارگیری روش های عددی مانند روش نیوتن-رافسون برآورد  $\nu_i$  را به دست آورد.

### چرخه دوم:

در چرخه دوم الگوریتم  $AEEM$  داده های کامل،  $y_c$  به صورت زیر تعریف می شود.

$$y_c = (y', z_1', \dots, z_n', u_1', \dots, u_n', w_1, \dots, w_n)'$$

مرحله  $E$ : امید ریاضی لگاریتم درستنمایی داده های کامل به شرط مشاهده ها به ازای برآورد جاری پارامترها،  $Q_2(\Psi; \Psi^{(k+\frac{1}{p})})$  محاسبه می گردد.

$$Q_2(\Psi; \Psi^{(k+\frac{1}{p})}) = \sum_{j=1}^n \sum_{i=1}^g \ln[\pi_i f_p(y_j; \mu_i, \Sigma_i, \nu_i)] E_{\Psi^{(k)}}(Z_{ij} | y_j)$$

مرحله  $M$ : شامل ماکسیمم کردن تابع  $Q_2(\Psi; \Psi^{(k+\frac{1}{p})})$  به ازای  $B_i$  و  $D_i$  است. برآورد  $B_i$  و  $D_i$  در تکرار  $(K+1)$  به ترتیب به صورت زیر برآورد خواهد شد.

$$B_i^{(k+1)} = V_i^{(k+\frac{1}{p})} \gamma_i^{(k)} \left( \gamma_i^{(k)'} V_i^{(k+\frac{1}{p})} \gamma_i^{(k)} + \Omega_i^{(k)} \right)^{-1},$$

$$D_i^{(k+1)} = \text{diag} \left\{ V_i^{(k+\frac{1}{p})} - V_i^{(k+\frac{1}{p})} \gamma_i^{(k)} B_i^{(k+1)'} \right\}$$

که

$$V_i^{(k+\frac{1}{p})} = \frac{\sum_{j=1}^n w_{ij}^{(k+\frac{1}{p})} \tau_{ij}^{(k+\frac{1}{p})} (y_j - \mu_i^{(k+1)}) (y_j - \mu_i^{(k+1)})'}{\sum_{j=1}^n \tau_{ij}^{(k+\frac{1}{p})}},$$

$$\gamma_i = (B_i B_i' + D_i)^{-1} B_i, \quad \Omega_i = I_q - \gamma_i' B_i$$

چرخه‌ی اول و دوم تا همگرایی الگوریتم تکرار می‌شوند. بر اساس آنچه به آن پرداخته شد در مدل تحلیل‌گر عاملی  $t$  فرض کردیم که تعداد مؤلفه‌ها و تعداد عامل‌ها معلوم هستند، اما متأسفانه در عمل تعداد مؤلفه‌ها و تعداد عامل‌ها مشخص نیستند. لذا در ادامه، در انتخاب یک مدل مناسب از میان یک مجموعه متناهی از مدل‌ها، از ملاک  $BIC$  استفاده می‌شود. ملاک اطلاع‌بیزی  $BIC$  که توسط شوارتز (۱۹۷۸) ارائه شد، عبارت است از:

$$BIC = -2 \log l(y, \hat{\theta}) + k \log n$$

که در آن  $\log l(y, \hat{\theta})$  مقدار لگاریتم تابع درست‌نمایی در نقطه ماکسیمم و  $k$  تعداد پارامترهای مدل است. مطابق تعریف فوق، مدلی که  $BIC$  کمتری داشته باشد مدل بهتری خواهد بود.

روش خوشه‌بندی، روشی برای شناسایی گروه‌های مشابه از اشیا در بین داده‌هاست که با استفاده از آن، حجم وسیعی از داده‌ها، به خوبی سازماندهی و خلاصه می‌شوند، به طوری که هر کدام از خوشه‌ها شامل مشاهده‌های مشابه بوده و همچنین هر خوشه تا جای ممکن متمایز از دیگر خوشه‌ها باشد. برای خوشه‌بندی داده‌ها روش‌های زیادی وجود دارد، از این رو جهت مقایسه‌ی روش‌های مختلف خوشه‌بندی دو شاخص رند و رند تعدیل‌شده در این مقاله معرفی خواهند شد.

یک روش عمومی برای ارزیابی نتایج خوشه‌بندی یک روش معین، مقایسه خوشه‌بندی به دست آمده از این روش با یک خوشه‌بندی درست و از پیش تعریف شده است. یک معیار معمول برای چنین مقایسه‌ای، شاخص رند است که اولین بار توسط رند (۱۹۷۱) معرفی شد. فرض کنید  $C = \{C_1, \dots, C_s\}$  ساختار خوشه‌بندی حاصل از یک روش خوشه‌بندی معین و  $P = \{P_1, \dots, P_t\}$  خوشه‌بندی درست و از پیش تعریف شده باشد. شاخص رند تطابق بین دو ساختار خوشه‌بندی  $P$  و  $C$  را به صورت زیر اندازه می‌گیرد.

$$RI(P, C) = \frac{a + d}{a + b + c + d} = \frac{a + d}{\binom{n}{2}}$$

که در آن  $a$  تعداد جفت مشاهده‌هایی است که متعلق به خوشه‌های مشابهی در  $P$  و همچنین به خوشه‌های مشابهی در  $C$  تخصیص داده شده‌اند،  $d$  جفت مشاهده‌هایی که متعلق به خوشه‌های متفاوتی در  $P$  و همچنین خوشه‌های متفاوتی در  $C$  هستند،  $b$  جفت مشاهده‌هایی که متعلق به خوشه‌های متفاوتی در  $P$  اما در خوشه‌های مشابهی در  $C$  قرار گرفته‌اند،  $c$  تعداد جفت مشاهده‌هایی که متعلق به خوشه‌های مشابهی در  $P$  هستند ولی در خوشه‌های متفاوتی در  $C$  قرار گرفته‌اند و  $n$  تعداد کل مشاهدات است.

این شاخص مقادیر بین صفر و یک می‌پذیرد و هرچه مقدار آن به یک نزدیک‌تر باشد، تطابق بین دو ساختار خوشه‌بندی بیشتر است.

موردی که به عنوان نقطه ضعف شاخص رند تلقی می‌شود این است که، مقدار متوسط این شاخص نامشخص است و با تغییر ساختار خوشه‌بندی، مقدار متوسط این شاخص تغییر می‌کند. برای اجتناب از این مشکل هوبرت و عربی (۱۹۸۵) شاخص رند تعدیل شده،  $ARI^2$ ، را ارائه دادند و آن را به شکل زیر تعریف کردند.

$$ARI = \frac{\frac{a+d}{\binom{n}{2}} - \frac{(a+b)(a+c)+(c+d)(b+d)}{\binom{n}{2}^2}}{1 - \frac{(a+b)(a+c)+(c+d)(b+d)}{\binom{n}{2}^2}}$$

این شاخص مقادیر کمتر از یک را با مقدار متوسط ثابت صفر می‌پذیرد و رفتاری مشابه رفتار شاخص رند دارد به این معنا که مقادیر بزرگ شاخص رند تعدیل شده بیان‌گر تطابق بیشتر بین دو ساختار خوشه‌بندی است. برای جزئیات بیشتر می‌توانید به یونگ و روزو (۲۰۰۰) مراجعه نمایید.

## ۵ تحلیل داده‌های ریزآرایه

فناوری استفاده از ریزآرایه درنگرش بشر، در مسائل زیستی، تحول عظیمی ایجاد نموده است. این فناوری پژوهشگران را قادر می‌سازد به‌جای بررسی تک به تک ژن‌ها، به طور موازی چندین هزار ژن را مورد بررسی قرار دهند تا عملکرد و نوع تعامل آن‌ها با یکدیگر بهتر مشخص شود. در اکثر داده‌های ریزآرایه، تعداد ژن‌ها از مرتبه چند هزار و چند ده هزار است، اما تعداد نمونه‌ی بافت‌ها به ندرت از ۱۰۰ مورد تجاوز می‌کند. به همین دلیل استفاده از تکنیک‌ها و روش‌های انتخاب ژن و کاهش بعد امری اجتناب‌ناپذیر می‌باشد. هدف خوشه‌بندی، گروه‌بندی ژن‌ها و نمونه‌ها با ویژگی‌های شبیه هم است. این مسأله می‌تواند به عنوان کاهش ابعاد سیستم نیز در نظر گرفته شود.

در این مقاله به معرفی روش  $EMMIX-GENE$  بر روی مجموعه داده‌های ریزآرایه سرطان روده بزرگ پرداخته خواهد شد. برای کسب اطلاعات به مک‌لاکلان و همکاران (۲۰۰۲) و مک‌لاکلان و همکاران (۲۰۰۷a) مراجعه نمایید. روش  $EMMIX-GENE$  شامل سه مرحله زیر است:

### مرحله ۱: تشخیص و حذف ژن‌های غیر مرتبط

برای کشف ژن‌های غیر مرتبط، هر ژن را به صورت جداگانه بررسی، و این که آیا ژن مربوطه قادر به تشخیص دو مؤلفه متناظر با نمونه‌های بافتی سالم و سرطانی هست یا

<sup>۲</sup> Adjusted Rand Index

نه را با آزمون  $G = 1 : H_0$  در مقابل  $G = 2 : H_1$  بررسی می‌کنیم. در این روش جهت کاستن حساسیت آزمون نسبت به نقاط دورافتاده، توزیع مؤلفه‌های آمیخته  $t$ ، در نظر گرفته می‌شود. در انتخاب یک ژن، شرط‌های زیر باید برقرار باشد:

$$-2Ln\lambda > b_1 \quad (۶)$$

و

$$\min\{N_1, N_2\} \geq b_2, \quad (۷)$$

که  $N_1$  و  $N_2$  تعداد بافت‌هایی است که به مؤلفه‌ها تخصیص داده شده‌اند و نقاط آستانه  $b_1$  و  $b_2$  از قبل توسط محقق تعیین می‌شود. اگر شرط (۷) برقرار ولی شرط (۸) برقرار نباشد، یکی از خوشه‌ها شامل تعداد کمی مشاهده و در واقع دورافتاده‌ها است. در این حالت ممکن است مؤلفه با تعداد مشاهده‌های بیشتر خود شامل دو زیرمؤلفه باشد. برای این‌که شانس مجددی برای شناسایی چنین ژن‌هایی وجود داشته باشد، آزمون  $G = 2 : H_0$  در مقابل  $G = 3 : H_1$  انجام می‌شود. شرط (۷) به همراه این شرط که حداقل حجم دو مؤلفه از سه مؤلفه، حداقل  $b_2$  است، در نظر گرفته می‌شود. اگر هر دو شرط برقرار باشد، ژن انتخاب و در غیر این صورت کنار گذاشته می‌شود.

### مرحله ۲: خوشه‌بندی ژن‌ها

منظور از یک خوشه‌ی ژنی، گروهی از ژن‌ها هستند که در میان همه‌ی ژن‌ها بردار بیان شبیه و نزدیک به هم دارند. ژن‌هایی که دارای بردار بیان شبیه و نزدیک به هم هستند، با یکدیگر همبستگی داشته و این همبستگی باعث می‌شود که قدرت تمییز یکسانی در دسته‌بندی نمونه‌ها داشته باشند. با این انگیزه، ژن‌های باقی‌مانده از مرحله اول با استفاده از الگوریتم  $K$ -میانگین خوشه‌بندی می‌شوند. تعداد خوشه‌های مدنظر توسط محقق تعیین می‌گردد. این تعداد باید طوری انتخاب شود که علاوه بر کم بودن مجموع تغییرپذیری‌های درون خوشه‌ها، متوسط تعداد ژن‌ها در هر خوشه خیلی بزرگ نباشد (از مشاهده‌ها کمتر باشد).

### مرحله ۳: خوشه‌بندی بافت‌ها

در این مقاله، بخش آخر الگوریتم را به مدل تحلیل‌گر عاملی آمیخته  $t$ ، تعمیم داده و در هر خوشه‌ی ژنی مرحله قبل سعی می‌کنیم بعد فضای ویژگی را کوچک‌تر و افراد را خوشه‌بندی نماییم. به منظور تعیین تعداد عامل‌ها در هر خوشه ژنی، معیار  $BIC$  مورد استفاده قرار می‌گیرد. در نهایت از میان خوشه‌های ژنی، خوشه‌ای که رند و رند تعدیل شده بزرگ‌تری داشته باشد انتخاب می‌شود.



## ۶ *EMMIX - GENE* در تحلیل ریزآرایه سرطان روده بزرگ

در این بخش به تحلیل یک مجموعه داده بیان ژنی ریزآرایه مربوط به سرطان روده بزرگ حاوی ۲۰۰۰ ژن از ۶۲ بافت نمونه خواهیم پرداخت. که ۴۰ نمونه از بافت‌های سرطانی بیماران و ۲۲ نمونه از بافت‌های سالم بیماران هستند. این داده‌ها، از داده‌های موجود در سایت «<http://www.maths.uq.edu.au/gim/emmix-gene>» است.

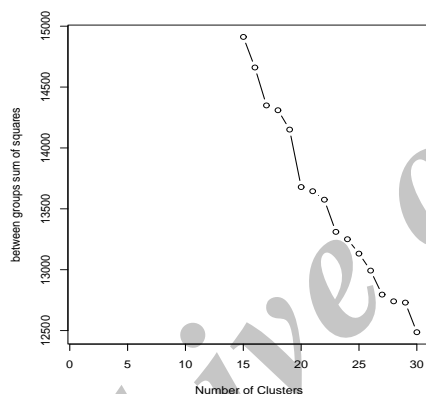
جدول ۱: مقدار معیار *RI* و *ARI* در هر خوشه‌ی ژنی

شماره خوشه	<i>q</i>	<i>ARI</i>	<i>RI</i>	شماره خوشه	<i>q</i>	<i>ARI</i>	<i>RI</i>
۱	۵	۰/۰۷	۰/۵۶	۱۱	۶	۰/۰۷	۰/۵۴
۲	۳	۰/۰۶	۰/۵۶	۱۲	۷	۰/۰۳	۰/۵۲
۳	۳	۰/۱۶	۰/۵۸	۱۳	۵	۰/۰۷	۰/۵۶
۴	۷	۰/۰۳	۰/۵۵	۱۴	۴	۰/۰۹	۰/۵۷
۵	۳	-۰/۰۳	۰/۴۹	۱۵	۳	۰/۰۹	۰/۵۷
۶	۶	۰/۱۱	۰/۵۶	۱۶	۴	۰/۰۸	۰/۵۵
۷	۴	۰/۰۸	۰/۵۵	۱۷	۵	۰/۱۱	۰/۵۶
۸	۶	۰/۰۲	۰/۵۲	۱۸	۹	۰/۰۱	۰/۵۰
۹	۸	۰/۰۴	۰/۵۴	۱۹	۴	-۰/۰۲	۰/۴۹
۱۰	۸	۰/۰۴	۰/۵۴	۲۰	۹	۰/۰۸	۰/۵۵

در اجرای روش *EMMIX - GENE* در مرحله اول،  $b_1 = b_2 = 8$  در نظر می‌گیریم، ۱۵۴۷ ژن به عنوان ژن‌های نامرتب حذف می‌شود و ۴۵۳ ژن باقی می‌ماند. پس از این مرحله، ژن‌های ریزآرایه را با استفاده از الگوریتم  $k$ -میانگین در  $k$  خوشه قرار خواهیم داد. جهت تعیین تعداد خوشه‌ها، مجموع تغییرپذیری‌های درون خوشه‌ای را به ازای  $k = 3, \dots, 15$  محاسبه کرده و در شکل ۱ رسم کرده‌ایم. همان‌طور که ملاحظه می‌شود با افزایش  $k$  مجموعه تغییرپذیری‌های درون خوشه‌ای همواره کاهش می‌یابد، اما چند کاهش زیاد به‌خصوص از ۱۹ به ۲۰ خوشه اشاره‌ای به این است که تعداد ۲۰ خوشه می‌تواند انتخاب مناسبی باشد. در ادامه با اجرای روش تحلیل‌گر عاملی آمیخته  $t$  با  $G = 2$  و  $q = 3, \dots, 9$  در هر خوشه، بافت‌ها را به دو گروه تقسیم می‌کنیم. برای ارزیابی نتایج خوشه‌بندی در هر خوشه از دو ملاک رند و رند تعدیل شده استفاده می‌کنیم که نتایج آن در جدول ۲ آمده است. مطابق جدول ۲، خوشه سوم بهترین نتیجه را در میان خوشه‌ها داراست. این خوشه در میان بافت‌های سرطانی ۱۴ نمونه و در میان بافت‌های سالم ۴ نمونه به اشتباه تخصیص داده شده‌اند. هرچند هنوز نرخ خطا نسبتاً بالا است، ولی مقایسه روش‌های فوق با روش‌های معمول خوشه‌بندی سلسله مراتبی،  $k$ -میانگین و مدل-پایه که در مکنیکلاس و مورفی (۲۰۱۰) ارائه شده است، نشان از برتری این روش دارد (جدول ۲).

## ۷ بحث و نتیجه گیری

مسئله‌ای که محققین در کاربردهای علمی با آن مواجه می‌شوند، بعد بالای مشاهده‌ها است. از این رو، مدل‌های تحلیل گر عاملی آمیخته دارای کاربردهای وسیعی در زمینه‌های مختلف علوم کامپیوتر، پزشکی و ژنتیک می‌باشد. یکی از حوزه‌هایی که در آن کمتر از این تکنیک استفاده شده، تحلیل داده‌های ریزآرایه است، که در این مقاله این روش را در تحلیل ریزآرایه‌های سرطان روده بزرگ استفاده کردیم. روش‌های آماری تحلیل داده‌های ریزآرایه، بدلیل تعداد متغیرهای بسیار زیاد و مشاهده‌های کم، هنوز به عملکرد خیلی رضایت‌بخشی منجر نشده‌اند. مک‌نیکلاس و مورفی (۲۰۱۰) تکنیک‌های دیگری را در تحلیل این داده‌ها به کار گرفتند و روش‌های مختلف خوشه‌بندی را بر اساس شاخص‌های رند و رند تعدیل شده مقایسه نمودند. نتایج به دست آمده توسط آن‌ها نشان می‌دهد، روش‌های معمول خوشه‌بندی سلسله‌مراتبی،  $k$ -میانگین و مدل-پایه نتوانسته‌اند نتایج بهتری از مدل تحلیل گر عاملی آمیخته  $t$  با  $ARI = 0/16$  و  $RI = 0/58$  کسب کنند.



شکل ۱: مجموع تغییرپذیری‌های درون خوشه‌ای به ازای تعداد خوشه‌های متفاوت

جدول ۲: شاخص  $ARI$  و  $RI$  حاصل از روش‌های معمول خوشه‌بندی داده‌های ریزآرایه سرطان روده بزرگ

روش	$ARI$	$RI$
سلسله‌مراتبی کامل	-۰/۰۲	۰/۵۰
سلسله‌مراتبی میانگین	-۰/۰۱	۰/۵۳
سلسله‌مراتبی منفرد	-۰/۰۱	۰/۵۳
$k$ -میانگین	-۰/۰۲	۰/۵۰
مدل-پایه با ماتریس وارینانس-کووارینانس $VII$	-۰/۰۱	۰/۵۰

## مراجع

- Ghahramani, Z. and Hinton, G. E. (1997), The EM algorithm for factor analyzers, *Technical Report NO. CRG-TR-96-1. Toronto: The University of Toronto*.
- Hubert, L. and Arabie, P. (1985), Comparing partitions, *Journal of classification*, **2**, 193-218.
- McLachlan, G. J. and Bean, R. and Wen, S. (2007a), Application of gene shaving and mixture models to cluster microarray gene expression data, *Cancer Inform*, **2**, 25-43.
- McLachlan, G. J. and Bean, R. W. and Ben-Tovim Jones, L. (2007b), Extension of the mixture of factor analyzers model to incorporate the multivariate t-distribution, *Computational Statistics and Data Analysis*, **51**, 5327-5338.
- McLachlan, G. J. and peel, D. (2000), Mixtures of factor analyzers, *Proceedings of the Seventeenth International Conference on Machine Learning*.
- McLachlan, G. J. and Peel, D. and Bean, R. W. (2002), A mixture model-based approach to the clustering of microarray expression data, *Bioinformatics*, **18**, 413-422.
- McNicholas, P. D. and Murphy, T. B. (2010), model-based clustering of microarray expression data via latent Gaussian mixture models, *Bioinformatics*, **21**, 2705-2712.
- Meng, X. L. and Rubin, D. B. (1993), Maximum likelihood estimation via the ECM algorithm: a general framework, *Biometrika*, **80**, 267-278.
- Meng, X. L. and van Dyk, D. (1997), The EM algorithm-an old folk song sung to a fast new tune (with discussion), *J. Roy. Soc. B*, **59**, 511-567.
- Rand, W. M. (1971), Objective criteria for the evaluation of clustering methods, *J. Amer. Statist. Assoc*, **66**, 846-850.
- Schwartz, G. (1978), Estimating the dimension of a model, *Annals of Statistics*, **6**, 31-38.
- Yeung, K. Y. and Ruzzo, W. L. (2000), An empirical study on principal component analysis for clustering gene expression data, *Tech. Report, University of Washington*, **17**, 763-774.