

تعمین تعداد خوشه‌ها در تحلیل‌ها مخته با استفاده از نتروپی نرمال شده

محمد قربانی^۱ محسن محمدزاده^۲

^۱ گروه مار دانشگاه ترز

^۲ گروه مار دانشگاه تربیت مدرس

چکیده یکی از مسایل مهم در تحلیل خوشه‌ای تعین تعداد خوشه‌هاست علی‌رغم اینکه آن مو و توسعه محققان زیادی مورد بحث واقع شده ولی همچنان به عنوان یک مشکل خوشه بندی مطرح است در این مقاله هدف تعین تعداد خوشه‌ها با استفاده از نتروپی نرمال شده در تحلیل مخته است که این معیار از ارتبای تعین درستنمایی کلی و درستنمایی رده بندی دست می‌دهد همچنین با استفاده از شبیه‌سازی مونت کارلو نشان داده خواهد شد که این معیار بهتر از معیارهای مانند BIC و AIC عمل می‌کند

واژه‌ها کلید واژه تحلیل مخته تعداد خوشه‌ها نتروپی نرمال شده درستنمایی رده بندی

مقدمه

یکی از مسایل مهم در بسیاری از مطالعات ژنتیکی و پزشکی خوشه بندی داده‌ها است که در N مشاهده M ویژگی g گروه همگن افزایش می‌شوند به وری که تشابهات درون گروهها ماکسیمم گردند روشهای مختلفی از جمله روشهای سلسله مراتبی رای خوشه بندی کردن داده‌ها به کار می‌روند که در تعرف فاصله بین دو خوشه با هم تفاوت دارند هارتگان همچنین به دلیل اینکه این روشها بر اساس مدل خاصی نانشده‌اند استنبای ماری بر اساس آنها امکان پذیر نیست و تعداد خوشه‌ها نیز به صورت ابتکاری با تعرف ستانه‌ای دلخواه تعین می‌شود رای رفع این مشکلات لازم است روش خوشه بندی حتی الامکان مبتنی بر سلسله محقق نبوده و بر اساس یک مدل با توزی احتمالی باشد تا توان در مورد ن استنبای ماری انجام داد معمولاً مجموعه مشاهدات تحت بررسی همگی از یک جامعه خاص نیستند و رای تشخیص آن که هر مشاهده از کدام جامعه مده است منتهی است فرض شود که هر مشاهده بر اساس ویژگی‌ها و خصوصیاتش دارای توزی احتمال خاصی است بنابراین جامعه‌ای مرکب از چند زیرجامعه دارای توزی احتمالی مخته صورت

$$f(x|\psi) = \lambda_1 f_1(x|\theta_1) + \dots + \lambda_g f_g(x|\theta_g)$$

است که در n رای $j = 1, \dots, g$ تا چگالی مولفه‌ها $f_j(\cdot)$ و $\lambda_j \leq \lambda_j <$ و معمولاً $\lambda = (\lambda_1, \dots, \lambda_g)$ و $\theta = (\theta_1, \dots, \theta_g)$ و $\psi = (\lambda, \theta)$ است. فرض می‌شود مولفه‌ها دارای توزیع تک مدی $n(\mu_j, \sigma_j^2)$ است و هدف از خوشه بندی تجزیه مولفه‌های چند مدی مبهم و مختلط مولفه‌های ساده تک مدی است. یکی از مسائل مهم خوشه بندی تعیین تعداد خوشه‌ها مولفه‌ها است علی‌رغم آنکه این موضوع توسط محققان زیادی مورد بحث واقع شده ولی همچنان به عنوان یک مشکل خوشه بندی مطرح است. در این مقاله از معاری نام معیار نتروپی نرمال شده رای تعیین تعداد خوشه‌ها استفاده می‌شود و کمک تکنیک شبه سازی مونت کارلو نشان داده خواهد شد این معیار بهتر از معیارهای مانند BIC و AIC عمل می‌کند.

تعیین معیارها خوشه بندی بر اساس مدل

فرض کنید X_1, \dots, X_n یک نمونه تصادفی از جامعه g از زیر جامعه‌های $1, \dots, g$ باشند و $\phi(x|\mu_j, \Sigma_j)$ تا چگالی متغیر تصادفی X_i در زیر جامعه j ام باشد. در این صورت اگر λ_j احتمال تعلق X_i به جامعه j ام باشد توزیع X_i عبارت است از

$$f(x; \psi) = \sum_{j=1}^g \lambda_j \phi(x|\mu_j, \Sigma_j) \quad (1)$$

که در n $\theta_j = (\mu_j, \Sigma_j)$ $\psi = (\lambda_1, \dots, \lambda_g, \theta_1, \dots, \theta_g)$ و $\sum_{j=1}^g \lambda_j = 1$ و $\lambda_j > 0$ تا درستی نمونه تصادفی به صورت

$$L(\psi|x) = \prod_{i=1}^n \sum_{j=1}^g \lambda_j f_j(x|\theta_j) \quad (2)$$

خواهد بود. رای $C_j = \{i, X_i \in j\}$ معیار خوشه بندی حاصل از ماکسیم کردن تا درستی

$$L(x|C) = \prod_{j=1}^g \lambda_j^{n_j} \prod_{X_i \in C_j} f(x_i|\theta_j) \quad (3)$$

معادل معیار حاصل از ماکسیم کردن خواهد بود. مکمل و فرالی و رافتی و واحدی و همکاران معمولاً مولفه اصلی X_i ها معلوم نیستند و رای مشخص کردن مولفه اصلی X_i متغیرهای گروه بندی Z_{ij} به صورت

$$Z_{ij} = \begin{cases} X_i \in j \\ X_i \notin j \end{cases}$$

تعرف می‌شوند. اساس مشخصه‌های گروه بندی لگاریتم تا درستنمای را می‌توان به صورت

$$\log L(\psi|x, Z) = \sum_{j=1}^g \sum_{i=1}^n z_{ij} \log \{\lambda_j f(x_i/\theta_j)\} \quad ()$$

نوشت اگر z_{ij} ها مشخص شدند عناصر خوشه j ام به صورت $C_j = \{j; z_{ij} > z_{ij'} \quad j \neq j'\}$ خواهد بود ولی اگر z_{ij} ها معلوم نباشند رای انجام تحلیل خوشه‌ای از ورود مقدار مورد انتظار آنها استفاده می‌شود پس تحلیل خوشه‌ای را می‌توان به عنوان ورود مقدار مورد انتظار z_{ij} معنی

$$E(Z_{ij}) = P(X_i \in j) = \frac{\lambda_j f(x_i|\theta_j)}{\sum_{j=1}^g \lambda_j f(x_i|\theta_j)}$$

تلقی نمود که لازم است رای هر j پارامترهای نامعلوم λ_j و θ_j ورود شوند الگوریتم امید رای و ماکسیم سازی (EM) روشی کار ردی در محاسبات تکراری رای به دست وردن ورود ماکسیم درستنمای پارامترها در توزی های مخته است فرض کند $\psi^{(0)}$ مقدار اول ψ باشد در ان صورت مراحل الگوریتم EM به صورت زیر خواهد بود مکین و کرشنان مرحله E محاسبه امید رای لگاریتم تا درستنمای در نقه $\psi^{(0)}$ به شر مشاهده داده‌های کامل

$$Q(\psi, \psi^{(0)}) = E_{\psi^{(0)}} \{\log L_c(\psi|x)\}$$

مرحله M دست وردن مقداری مانند ψ^* رای ψ به وری که

$$Q(\psi^*, \psi^{(0)}) = \max_{\psi \in \Omega} (Q(\psi, \psi^{(0)}))$$

مراحل E و M تا زمانی تکرار می‌شوند که شر همگرایی $|L(\psi^{(k+1)}) - L(\psi^{(k)})| < \epsilon$: $\forall \epsilon >$ رقرار شود جی اف وو

فرض کند $X_1 \dots X_n$ داده‌های ناکامل و $Y_i = (X_i, Z_i)$ داده‌های کامل در الگوریتم EM باشند در ان صورت ورود پارامترهای استفاده از ان الگوریتم به صورت زیر خواهند بود

$$\lambda_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n z_{ji}^{(k)}$$

1) Expectation-Maximization Algorithm

$$\mu_i^{(k+1)} = \frac{\sum_{j=1}^n z_{ji}^{(k)} X_j}{\sum_{j=1}^n z_{ij}^{(k)}}$$

$$(\Sigma)^{(k+1)} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^n z_{ij}^{(k)} (X_j - \mu_i^{(k+1)})(X_j - \mu_i^{(k+1)})'$$

$$\hat{z}_{ij}^{(k)} = \frac{\pi_i^{(k)} f(x_j | \theta_i^{(k)})}{\sum_{i=1}^k \pi_i^{(k)} f(x_j | \theta_i^{(k)})}$$

با ورود z_{ij} به وسيله الگوریتم EM می توان تحلیل خوشه‌ای را متناهی‌تر مقدار $\hat{z}_{ij}^{(k)}$ انجام داد. واحدی اصل محمدزاده و قرانی و مک‌لن و کرشنان ولی کی از مسائل اساسی خوشه‌بندی تعداد خوشه‌هاست. این روش استفاده از معادله‌های مختلف توسط اغلب محققان واحدی و محمدزاده و قرانی فرآیند و رافتری بررسی شده ولی همچنان به عنوان یک مشکل خوشه‌بندی مطرح است.

تعیین تعداد خوشه‌ها

رایج‌ترین معیار نتروپی نرمال شده در تعیین تعداد خوشه‌ها نخست به تعریف نتروپی می‌پردازیم. متغیر تصادفی X با تابع جرم احتمالی $P(x) = P(X = x)$ و مقادیر $\{x_i, i = 1, \dots, n\}$ در نظر بگیرد. عدم قطعیت مرتبه n را می‌توان مشاهده از متغیر X را نتروپی متغیر تصادفی X نامیده و با نماد $H(X)$ نشان می‌دهند که عبارت است از

$$H(X) = - \sum P(x) \log P(x)$$

حال لگاریتم تا درست‌نمایی را در نظر بگیرد

$$L(g) = \sum_{i=1}^n \ln \sum_{j=1}^g \hat{\lambda}_j f_j(x_i | \hat{\theta}_j)$$

که در آن $\hat{\lambda}_j$ و $\hat{\theta}_j$ ورودهای ماکزیمم درست‌نمایی λ_j و θ_j می‌باشند. استفاده از محاسبات مستقیم می‌تواند نشان داد

$$L(g) = C(g) + E(g)$$

که در آن

$$C(g) = \sum_{j=1}^g \sum_{i=1}^n \hat{z}_{ij} \log \{ \hat{\lambda}_j f(x_i | \hat{\theta}_j) \}$$

$$E(g) = - \sum_{j=1}^g \sum_{i=1}^n \hat{z}_{ij} \ln \hat{z}_{ij} \geq$$

راساس روا فوق لگارتم تا درستنمای $L(g)$ لگارتم تا درستنمای رده بندی $C(g)$ و معار نترویی $E(g)$ تفکک شده است اگر چه نترویی $E(g)$ معار را لا است ولی نمی توان مستقما از ن رای تع ن تعداد خوشه ها استفاده کرد ز را $L(g)$ ک تا افزایش از g است و امد نرمالزه شود دارم

$$= \frac{C(g) - C(\cdot)}{L(g) - L(\cdot)} + \frac{E(g) - E(\cdot)}{L(g) - L(\cdot)} \quad g > \quad ()$$

و معار نترویی نرمال شده که امد رای تع ن تعداد خوشه ها می نیم شود عبارت است از

$$NEC(g) = \frac{E(g)}{L(g) - L(\cdot)}, \quad E(\cdot) =$$

اما $NEC(\cdot) = \infty$ خوش تعریف نیست زیرا $NEC(\cdot) = \infty$ ارنایکی سلوکس و گوورت لذا ه ور مستقم قادر ه مقاسه حالت $g =$ در رار $g >$ در استفاده از $NEC(g)$ نخواهم ود لذا می است $NEC(g)$ را رای رف ان مشکل خاص سه ده م رای تصم مگری ن $g =$ و $g >$ می توان از ان نقه نرکه $C(g)$ معاری رای اندازه گیری دقت افزاز داده ها در خوشه و $L(\cdot) = C(\cdot)$ معاری رای اندازه گیری دقت رازش تک خوشه ه داده هاست استفاده کرد هنگام مقاسه دو تعداد از خوشه ها g و g' ممکن است ک مدل مخته ا تعداد پارامترهای زاد مناسب ه نر رسد ولی از دیدگاه ساده منقی است اگر رای $g > g'$ $C(g) < C(g')$ اشد g ر ترجیح داده شود بنا بران رای انتخا $g >$ در رار $g =$ لازم است $C(g) > C(\cdot) = L(\cdot)$ در صورتی که $C(g) > L(\cdot)$ تمام جزهای معادله نامنفی خواهند ود که در ان صورت $NEC(g) \leq$ که تنها حالت $g >$ ان شر را فراهم می کند پس اگر $NEC(g) \leq$ نباشد دللی رانتخا ش از ک خوشه وجود ندارد

مقاسه معارها

در ان بخش معارهای NEC و BIC و AIC راساس تکنک شبه سازی مونت کارلو مورد مقاسه قرار می گردند رای ان من ور نمونه های ه حجم n از توز نرمال تک متغره ا و ارنس های رار ک و همچن نرمال مخته دو متغره ا ماترس های و ارنس رار I تولد

می‌شود

رای توزی نرمال تک متغیره چهار نو توزی پارامترهای مختلف زرد در: ر گرفته شده است
الف توزی نرمال استاندارد

توزی نرمال مخته دو مولفه‌ای مانگن‌های $\mu_1 =$ و $\mu_2 =$ و نسبت‌های مخته‌گی
رار

ج توزی مخته نرمال دو مولفه‌ای مانگن‌های $\mu_1 =$ و $\mu_2 =$ و نسبت‌های مخته‌گی
 $\lambda_1 =$ و $\lambda_2 =$

د توزی مخته نرمال سه مولفه‌ای مانگن‌های $\mu_1 =$ $\mu_2 =$ و $\mu_3 =$ و نسبت‌های
مخته‌گی مساوی

رای حالت دو متغیره سه نو توزی زرد در: ر گرفته شده است
الف توزی نرمال دو متغیره استاندارد

توزی مخته نرمال دو مولفه‌ای ردار مانگن $(,) = \mu'_1$ و $(,) = \mu'_2$ و نسبت‌های
مساوی

ج توزی مخته نرمال سه مولفه‌ای ردار مانگن $(,) = \mu'_1$ $(,) = \mu'_2$ و $(,) = \mu'_3$ و
 $(, -) = \mu'_4$ و نسبت‌های مساوی

از هر یک از توزی‌ها دو نمونه با حجم‌های $n =$ و $n =$ هر کدام با تعداد رار
شبه‌سازی شده است. ر اساس نمونه‌های تولد شده پارامترهای مدل‌های مخته ر اساس
الگوریتم EM رورد شده‌اند. در جدول نتایج مرو h مانگن و انحراف معیار نترویی
 NEC ان شده است که در d نشان دهنده عدد فای نمونه‌ای و g نشان دهنده تعداد
خوشه‌هاست. ر اساس نتایج با دست مده NEC رای $n =$ هتر از $n =$
عمل می‌کند. با عنوان مثال در جدول رای حالت الف توزی مخته نرمال تک متغیره
مانگن نترویی رای $g =$ کمتر از $g =$ حالت‌هاست لذا فرض تک خوشه ودن پذیرفته می‌شود
همچون رای حالت توزی دو متغیره وقتی که $n =$ است مانگن معیار نترویی رای
 $g =$ رار است که کمتر از مقدار n رای تعداد خوشه‌های $g =$ و $g =$
است و فرض دو مولفه‌ای ودن را تا د می‌کند. جدول درصد فراوانی انتخابی توزی مخته
 g مولفه‌ای ر اساس استفاده از معارهای AIC و BIC و NEC رای حالت‌های مختلف مذکور
و رای تعداد خوشه‌های مختلف ان می‌کند.

با عنوان مثال رای حالت اول توزی نرمال تک متغیره در درصد اوقات NEC جوا
صحیح داده است در حالیکه در درصد اوقات AIC پاسخ صحیح می‌دهد. در حالت
کلی ر اساس جدول نتیجه می‌شود که کفایت تصمیم‌گیری NEC ن AIC و BIC قرار
دارد. معیار AIC تعداد مولفه‌های توزی مخته را اندکی ش تخمین می‌کند و ر عکس BIC
تعداد مولفه‌ها را اندکی کمتر تخمین می‌زند.

جدول مانگن وانحراف معار رای معار NEC ر اساس توز های مخته نرمال تک متغره و دو متغره

d	n	پارامترهای توز	تعداد خوشهها			
			g = 1	g = 2	g = 3	g = 4
		$p_1 = 1$ $\mu_1 = 0$	۱۶/۴۱ (۶/۱۵)	۱۷/۹۳ (۷/۹۷)	۲۲/۸۳ (۸/۶۱)	۲۹/۳۲ (۱۰/۳۵)
۱	۲۰۰	$p_1 = 0.5$ $\mu_1 = 0, \mu_2 = 2$	۱۴/۸۳ (۵/۸۴)	۱۰/۹۸ (۴/۳۶)	۱۵/۷۳۵ (۵/۴۹)	۲۲/۱۰ (۶/۸۱)
		$p_1 = 0.7$ $\mu_1 = 0, \mu_2 = 2$	۱۵/۹۸ (۶/۱۲)	۱۵/۷۸ (۵/۹۵)	۲۲/۶۷ (۷/۳۱)	۳۰/۵۲ (۱۰/۱۱)
		$p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$	۲۰/۴۵ (۱۲/۸۸)	۷/۴۲ (۲/۳۵)	۵/۸۹ (۰/۹۷)	۹/۱۶ (۳/۴۱)
		$p_1 = 1$ $\mu_1 = 0$	۱۷/۹۲ (۷/۲۱)	۱۸/۲۲ (۸/۹۱)	۲۳/۲۵ (۱۰/۱۲)	۳۰/۲۹ (۱۲/۷۴)
۱	۵۰	$p_1 = 0.5$ $\mu_1 = 0, \mu_2 = 2$	۱۴/۲۷ (۷/۴۲)	۱۲/۲۳ (۵/۸۲)	۱۷/۱۸ (۷/۷۹)	۲۴/۹۴ (۶/۹۱)
		$p_1 = 0.7$ $\mu_1 = 0, \mu_2 = 2$	۱۸/۵۴ (۱۰/۱۷)	۱۶/۷۶ (۷/۶۴)	۲۴/۵۲ (۹/۸۶)	۳۲/۸۶ (۱۳/۲۲)
		$p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$	۲۸/۳۸ (۱۱/۸۲)	۵/۸۲ (۲/۱۰)	۵/۱۵ (۰/۶۲)	۷/۵۵ (۲/۸۴)
		$p_1 = 1$ $\mu_1 = [0, 0]$	۳/۴۱ (۳/۰۲)	۱۶/۸۴ (۶/۱۲)	۸/۶۹ (۶/۶۴)	۲۴/۴۷ (۷/۷۶)
۲	۲۰۰	$p_1 = 0.5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$	۸/۸۶ (۷/۴۲)	۶/۲۲ (۵/۸۲)	۸/۱۲ (۷/۷۹)	۱۰/۲۲ (۶/۹۱)
		$p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$	۱۱/۱۹ (۱۱/۸۲)	۳/۵۹ (۲/۱۰)	۲/۰۲ (۰/۶۲)	۲/۹۰ (۲/۸۴)
		$p_1 = 1$ $\mu_1 = [0, 0]$	۹/۶۷ (۴/۱۵)	۱۱/۵۸ (۶/۹۹)	۲۳/۹۶ (۷/۴۴)	۳۰/۹۵ (۸/۹۴)
۲	۵۰	$p_1 = 0.5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$	۱۰/۸۱ (۶/۶۸)	۸/۹۷ (۴/۱۲)	۱۰/۲۲ (۵/۸۳)	۱۲/۷۰ (۷/۵۶)
		$p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$	۲/۹۲ (۰/۸۸)	۲/۳۵ (۲/۱۰)	۴/۴۳ (۰/۶۲)	۱۱/۹۷ (۲/۸۴)

جدول درصد فراوانی انتخاب g خوشه راساس توزی های مخته نرمال تک متغره و دو متغره

d	n	پارامترها توز	تعداد خوشه ها	معیارها		
				AIC	BIC	NEC
		$p_{\lambda} =$ $\mu_{\lambda} =$				
		$p_{\lambda} = /$ $\mu_{\lambda} = , \mu_{\tau} =$				
		$p_{\lambda} = /$ $\mu_{\lambda} = , \mu_{\tau} =$				
		$p_{\lambda} = p_{\tau} = p_{\sigma} = \frac{1}{\sigma}$ $\mu_{\lambda} = , \mu_{\tau} = , \mu_{\sigma} =$	2			
		$p_{\lambda} =$ $\mu_{\lambda} =$				
		$p_{\lambda} = /$ $\mu_{\lambda} = , \mu_{\tau} =$				
		$p_{\lambda} = /$ $\mu_{\lambda} = , \mu_{\tau} =$				
		$p_{\lambda} = p_{\tau} = p_{\sigma} = \frac{1}{\sigma}$ $\mu_{\lambda} = , \mu_{\tau} = , \mu_{\sigma} =$				
		$p_{\lambda} =$ $\mu_{\lambda} = [,]$				
		$p_{\lambda} = /$ $\mu_{\lambda} = [,] , \mu_{\tau} = [,]$				
		$p_{\lambda} = p_{\tau} = p_{\sigma} = \frac{1}{\sigma}$ $\mu_{\lambda} = [,] , \mu_{\tau} = [,] , \mu_{\sigma} = [, -]$				
		$p_{\lambda} =$ $\mu_{\lambda} = [,]$				
		$p_{\lambda} = /$ $\mu_{\lambda} = [,] , \mu_{\tau} = [,]$				
		$p_{\lambda} = p_{\tau} = p_{\sigma} = \frac{1}{\sigma}$ $\mu_{\lambda} = [,] , \mu_{\tau} = [,] , \mu_{\sigma} = [, -]$				

مراج

- [1] Biernacki, C., Celeux, G. and Govert, G. (1999), An Improvement of the NEC Criterion for Assessing the Number of Clusters in a Mixture model, Pattern Recognition Letters, 20, 267-272.
- [2] Celeux, G. and Soromenho, G. (1996), An Entropy Criterion for Assessing the Number of Clusters in a Mixture Model. Journal of Classification , 13, 195-212.
- [3] Fray, C. and Raftery, A. E. (1998), How Many Clusters? Which Clustering Method? Answer via Model-Based Cluster Analysis. Technical Report, No. 329. Seattle: Department of Statistics, University of Washington.
- [4] Fraley, C. and Raftery, A. E. (1999), MCLUST: Software for Model-Based Cluster Analysis, J. Classification, 16, 297-306.
- [5] Hartigan, J. A. (1975), Clustering Algorithms. Wiley, New York.
- [6] Jeff C. F., (1983), On the Convergence of the EM Algorithm, Annals of Statistics, 11, 95-103.
- [7] McLachlan, G. J. (1982), The Classification and Mixture Maximum Likelihood Approaches to Cluster Analysis. In Krishnan, P. R. and Kanal, L. N. (eds), Handbook of Statistics, 2, 199-208. North-Holland, Amsterdam.
- [8] McLachlan, G. J. and Krishnan, T. (1997), The EM Algorithm and Extensions. Wiley, New York.

محمد قاسم وحدی اصل محسن محمدزاده و محمد قرانی
 خوشه بندی احتمالاتی ر اساس معار الا زی مجموعه مقالات ششمین کنفرانس
 نالمللی مارا ران ص

Archiving at SID