

## یک معیار جدید انتخاب مدل مبتنی بر داده تاگی

صدیقه اسحق، حسین باغیشنی، نگار اقبال

گروه آمار، دانشگاه صنعتی شاهرود

تاریخ دریافت: ۱۳۹۶/۳/۲۸ تاریخ آخرین بازنگری: ۱۳۹۷/۱/۲۶

**چکیده:** یک چالش اساسی در استنباط مدل‌های آمیخته، معرفی معیارهای کارا برای انتخاب مدل است. منبع اصلی این چالش نیز برازش و محاسبه ماکسیم تابع درستنمایی مدل می‌باشد. داده تاگی روش جدیدی است که برای برازش کارای مدل‌های آمیخته با روش ماکسیم درستنمایی پیشنهاد شده است. این روش، اخیراً، طرفداران زیادی پیدا کرده است و مشکلات عمده سایر روش‌های استنباط مبتنی بر درستنمایی در مدل‌های آمیخته را ندارد. یکی از معایب این روش، عدم توانایی محاسبه مقدار ماکسیم تابع درستنمایی است. این مقدار یک کمیت کلیدی در معرفی و محاسبه معیارهای انتخاب مدل محسوب می‌شود. بنابراین به نظر می‌رسد با روش داده تاگی نمی‌توان یک معیار اطلاع مناسب، به‌طور مستقیم، برای یافتن بهترین مدل در رده مدل‌های آمیخته، تعریف کرد. این پژوهش تلاشی است در جهت نقض این باور. در این مقاله، یک معیار مبتنی بر روش داده تاگی معرفی می‌شود و عملکرد آن در یک مطالعه شبیه‌سازی مورد ارزیابی قرار می‌گیرد.

واژه‌های کلیدی: الگوریتم MCMC، مدل آمیخته خطی تعمیم‌یافته، معیار انتخاب مدل، داده تاگی.

## ۱ مقدمه

در کاربردهای مختلف ماهیت متغیر پاسخ، ناگاوسی (گسسته) است و نوعی وابستگی بین مشاهدات آن وجود دارد. به عنوان نمونه می‌توان به داده‌های گسسته طولی، خوشه‌ای، و فضایی اشاره کرد. مدل‌های آمیخته

خطی تعمیم‌یافته<sup>۱</sup> (GLMMs) تعمیمی متداول از مدل‌های خطی تعمیم‌یافته<sup>۲</sup> (GLMs) هستند که برای تحلیل داده‌های وابسته گسسته استفاده می‌شوند. در این مدل‌ها، با افزودن اثرات تصادفی (متغیرهای پنهان) به پیش‌گوی خطی، وابستگی مشاهدات وارد می‌شود (بریسلو و کلیتون، ۱۹۹۳). برازش مدل و پیش‌گویی، دو هدف کلیدی در رده GLMM محسوب می‌شوند. دستیابی به این دو هدف با هر دو رهیافت بسامدی و بیزی دارای محدودیت‌ها و مشکلاتی است:

- استنباط‌های مبتنی بر درست‌نمایی، حل عددی انتگرال‌های با بعد بالا را شامل می‌شوند که می‌تواند بسیار پیچیده و هزینه‌بر باشد. به عبارتی، حضور انتگرال‌های رام‌نشده (با بعد بالا) عمده‌ترین مشکل بر سر راه استنباط‌های مبتنی بر درست‌نمایی در رده GLMM است که مساله انتخاب مدل را نیز شامل می‌شود.
- استنباط‌های بیزی، به دلیل پیدایش الگوریتم‌های نمونه‌گیری مونت کارلوی زنجیر مارکوفی<sup>۳</sup> (MCMC)، مشکلات حل انتگرال‌های رام‌نشده و محاسبه مشتق توابع پیچیده را ندارند. در مقابل، این استنباط‌ها همواره با مساله انتخاب توزیع پیشین و انتقاد وابستگی نتایج به این انتخاب همراه هستند. در اغلب موارد، در این مدل‌ها، از توزیع‌های پیشین ناآگاهی‌بخش عینی<sup>۴</sup> استفاده می‌شود که عدم وجود تعریفی واحد برای این نوع از توزیع‌های پیشین، به پیشنهاد‌های مختلف و در نتیجه استنباط‌های ممکن مختلف منتهی خواهد شد (ترابی و همکاران، ۲۰۱۵). البته، اخیراً، برگر و همکاران (۲۰۰۹) با ارایه یک تعریف یکتا و سراسر برای توزیع‌های پیشین مرجع<sup>۵</sup>، که دسته‌ای از پیشین‌های عینی محسوب می‌شوند، تا حدودی این مشکل را مرتفع کرده‌اند. با این وجود نویسندگان مایل هستند دو نکته را یادآوری کنند:

۱. آماردان‌های بیزی ذهنی‌گرا<sup>۶</sup> اعتقادی به استفاده از پیشین‌های عینی ندارند و استنباط بیزی صرفاً به روش‌های بیزی عینی‌گرا محدود نمی‌شود. بنابراین آن‌ها پیشین‌های ذهنی خود را دارند که نحوه انتخاب آن‌ها می‌تواند از طرف سایرین مورد انتقاد قرار گیرد.

<sup>1</sup>Generalized linear mixed models

<sup>2</sup>Generalized linear models

<sup>3</sup>Markov chain Monte Carlo

<sup>4</sup>Objective

<sup>5</sup>Reference prior

<sup>6</sup>Subjective Bayesians

۲. سرعت کند و همگرایی ضعیف الگوریتم‌های MCMC در GLMM به انتخاب توزیع پیشین وابسته است.

تلاش‌های زیادی برای رفع این معایب در هر دو دیدگاه استنباطی مذکور صورت گرفته‌اند. به‌عنوان چند نمونه می‌توان به بریسلو و کلیتون (۱۹۹۳)، پینیرو و بیتس (۱۹۹۵)، و کریستنسن (۲۰۰۴) در رهیافت مبتنی بر درست‌نمایی، و کریستنسن و واگاپترسن (۲۰۰۲)، کریستنسن و همکاران (۲۰۰۶)، فونگ و همکاران (۲۰۱۰)، و حسینی و همکاران (۲۰۱۱) در دیدگاه بیزی اشاره کرد.

یک روش جانشین برای برازش مبتنی بر درست‌نمایی در رده GLMM، روش داده تاگی<sup>۷</sup> (DC) است که اولین بار توسط لهله و همکاران (۲۰۰۷) معرفی و به‌کار گرفته شد. روش DC، روشی ساده برای محاسبه MLE پارامترهای مدل با بهره‌گیری از الگوریتم‌های MCMC است. این روش از یک رهیافت بیزی به‌عنوان ابزار محاسبه MLE پارامترهای مدل استفاده می‌کند و نتایج حاصل از آن نسبت به انتخاب توزیع‌های پیشین، ناورد<sup>۸</sup> هستند. افزون بر این، مراحل برآورد پارامترهای مدل (و دقت آن‌ها) شامل محاسبه ساده مقادیر میانگین و واریانس نمونه‌ای است که با الگوریتم MCMC تولید می‌شود و نیازی به ماکسیم‌سازی عددی و مشتق‌گیری از یک تابع درست‌نمایی پیچیده نیست.

روش DC محدودیت‌هایی را نیز به همراه دارد. یکی از آن‌ها به مقدار ماکسیم تابع درست‌نمایی اختصاص دارد. این روش تنها برآورد ML پارامترهای مدل را فراهم می‌کند و قادر به محاسبه مقدار ماکسیم تابع درست‌نمایی، در آن برآوردها، نیست. ماکسیم تابع درست‌نمایی، کمیتی کلیدی در تعریف معیارهای انتخاب مدل به شمار می‌آید. بنابراین، با استفاده از روش DC نمی‌توان یک معیار انتخاب مدل (مستقیم) در رده‌ای از مدل‌های ممکن تعریف کرد. پیشنهادی برای رفع این مشکل توسط پونچیانو و همکاران (۲۰۰۹) ارائه شد. آن‌ها برای محاسبه نسبت درست‌نمایی، با استفاده از روش DC، یک الگوریتم ساده با نام الگوریتم نسبت درست‌نمایی<sup>۹</sup> DC معرفی کردند و از آن برای محاسبه اختلاف مقادیر معیارهای اطلاع مورد نیاز، به‌منظور مقایسه جفتی بین دو مدل، استفاده کردند. در روش پیشنهادی پونچیانو و همکارانش، برای مقایسه چند مدل نامزد و انتخاب بهترین مدل، از منظر یک معیار انتخاب مدل مثل معیار اطلاع آکاییک<sup>۱۰</sup> (AIC)، باید ابتدا مدل‌ها دو به دو با هم مقایسه و سپس مدل بهتر از ترکیب نتایج همه مقایسه‌های جفتی انتخاب شود. با این فرآیند، به‌راحتی قابل درک است که حتی اگر بعد رده

<sup>7</sup>Data cloning

<sup>8</sup>Invariant

<sup>9</sup>Data cloned likelihood ratio

<sup>10</sup>Akaike information criterion

مدل‌ها متوسط هم باشد، تعداد مقایسه‌های جفتی می‌تواند خیلی بزرگ و انتخاب بهترین، ناممکن شود. در این مقاله، بر اساس روش داده‌تاگی، معیار اطلاعاتی معرفی می‌شود که به‌طور مستقل مدل‌های نامزد حاضر در یک رده را رتبه‌بندی می‌کند. این معیار تقریباً مشابه AIC تعریف می‌شود. در بخش ۲ مدل‌های آمیخته خطی تعمیم‌یافته و در بخش ۳ روش داده‌تاگی معرفی می‌شود. در بخش ۴ معیار اطلاع جدید مبتنی بر DC ارایه می‌شود و در بخش ۵ در یک مطالعه شبیه‌سازی، عملکرد معیار معرفی‌شده مورد ارزیابی قرار می‌گیرد. در پایان، بحث و نتیجه‌گیری ارایه خواهد شد.

## ۲ مدل‌های آمیخته خطی تعمیم‌یافته

در یک GLMM فرض می‌شود متغیرهای پاسخ با شرط معلوم بودن اثرات تصادفی، مستقل از هم بوده و دارای توزیعی از خانواده توزیع‌های نمایی هستند. برای تشریح مدل، داده‌های خوشه‌ای را در نظر بگیرید که اندازه‌های مکرر متغیر پاسخ در یک نمونه تصادفی از  $m$  خوشه رخ داده باشند. فرض کنید بردار  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_m)^T$ ،  $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^T$ ،  $i = 1, \dots, m$ ، مشاهدات پاسخ خوشه  $i$ ام،  $\mathbf{u}_i = (u_{i1}, \dots, u_{iq})^T$  کل مشاهدات و  $n = \sum_{i=1}^m n_i$  حجم نمونه کل باشند. همچنین فرض کنید  $q$  بعدی اثرات تصادفی در خوشه  $i$ ام باشد. با شرط معلوم بودن اثرات تصادفی  $\mathbf{u}_i$ ، توزیع هر یک از متغیرهای پاسخ عضوی از خانواده توزیع‌های نمایی با تابع چگالی

$$f(y_{ij} | \mathbf{u}_i, \boldsymbol{\beta}) = \exp \{ y_{ij} (\mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{v}_{ij}^T \mathbf{u}_i) - a(\mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{v}_{ij}^T \mathbf{u}_i) + c(y_{ij}) \}$$

است، که در آن  $\mathbf{x}_{ij}$  و  $\mathbf{v}_{ij}$  برای  $i = 1, \dots, m$  و  $j = 1, \dots, n_i$  بردارهای متغیرهای تبیینی به ترتیب  $p$  و  $q$  بعدی متناظر با اثرات ثابت  $\boldsymbol{\beta}$  و تصادفی  $\mathbf{u}$  و  $a(\cdot)$  و  $c(\cdot)$  توابعی معلوم هستند. در این مدل، پیش‌گوی خطی  $\mu_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{v}_{ij}^T \mathbf{u}_i$  پارامتر کانونی خانواده توزیع‌های نمایی است. اگر میانگین شرطی مدل باشد، آن‌گاه  $g(\mu_{ij}) = \eta_{ij}$  که در آن  $g(\cdot)$  تابع پیوند مدل است. علاوه بر این، توزیع اثرات تصادفی  $\mathbf{u}_i$  معمولاً گاوسی  $\mathbf{u}_i | \boldsymbol{\Sigma}_\theta \sim N(\mathbf{0}, \boldsymbol{\Sigma}_\theta)$  در نظر گرفته می‌شود که در آن  $\boldsymbol{\theta}$  بردار  $r$  بعدی پارامترهای وابستگی مدل (مولفه‌های واریانس) است و ماتریس کوواریانس  $\boldsymbol{\Sigma}_\theta$  معمولاً به صورت پارامتری در نظر گرفته می‌شود. به عبارت دیگر با معلوم بودن پارامترهای  $\boldsymbol{\theta}$  ماتریس کوواریانس معلوم خواهد بود.

برای انجام استنباط مبتنی بر درستنمایی باید تابع درستنمایی کناری<sup>۱۱</sup> مدل به صورت

$$L(\beta, \theta) \propto \prod_{i=1}^m \int \prod_{j=1}^{n_i} f(y_{ij} | \mathbf{u}_i, \beta) \phi_q(\mathbf{u}_i; \circ, \Sigma_\theta) d\mathbf{u}_i \quad (1)$$

محاسبه شود، که در آن  $\phi_q(\cdot; \circ, \Sigma_\theta)$  تابع چگالی نرمال  $q$  متغیره با بردار میانگین صفر و ماتریس کوواریانس  $\Sigma_\theta$  است. محاسبه تابع درستنمایی (۱) مستلزم حل انتگرالی با بعد  $q$  است و برای داده‌های حجیم، محاسبه آن کاری دشوار خواهد بود. ماکسیم‌سازی انتگرال محاسبه‌شده، مرحله بعدی برای محاسبه MLE پارامترهای مدل است. داده تاگی راهی جذاب برای گریز از محاسبه چنین انتگرال‌هایی و ماکسیم‌سازی نتیجه آن است.

### ۳ داده تاگی

داده تاگی، ترفندی محاسباتی است که به‌عنوان روشی کارا برای استنباط مبتنی بر درستنمایی در مدل‌های آمیخته به‌کار می‌رود. این روش از الگوریتم‌های نمونه‌گیری MCMC برای تولید نمونه از یک توزیع پسین به‌طور مصنوعی ساخته‌شده<sup>۱۲</sup>، با نام توزیع پسین مبتنی بر داده<sup>۱۳</sup> تاگی، برای محاسبه MLE پارامترها و برآورد واریانس آن‌ها استفاده می‌کند. مدل مبتنی بر داده تاگی، یک مدل بیزی کامل ساخته‌شده بر اساس  $k$  نسخه تکرار شده داده‌های اصلی است: ابتدا بردار  $n_k$  بعدی  $\mathbf{y}^{(k)} = (\mathbf{y}, \dots, \mathbf{y})$  از کپی کردن  $k$  بار بردار  $n$  بعدی مشاهدات  $\mathbf{y}$  تشکیل می‌شود. به همین طریق مقادیر متغیرهای تبیینی نیز  $k$  بار تکرار می‌شوند. همچنین با استفاده از توزیع احتمالی اثرات تصادفی،  $k$  تحقق از آن تولید و بردار همسانه‌سازی اثرات تصادفی  $\mathbf{u}^{(k)} = (\mathbf{u}_1, \dots, \mathbf{u}_k)$  نتیجه می‌شود. برای پارامترهای نامعلوم مدل نیز توزیع پیشین  $\pi(\beta, \theta)$  در نظر گرفته می‌شود. با فرض مستقل بودن نسخه‌های کپی‌شده داده‌ها از هم، تابع درستنمایی جدید که همان تابع درستنمایی اصلی به توان  $k$  است، به دست می‌آید و توزیع پسین مبتنی بر داده تاگی با یک رهیافت MCMC به‌طور تقریبی محاسبه می‌شود.

از آن‌جا که ایده پشتیبان روش DC بر ویژگی‌های جانبی توزیع‌های پسین بنا نهاده شده است (باغیشنی و محمدزاده، ۲۰۱۱)، تعداد نسخه‌ها، یعنی  $k$ ، باید به اندازه کافی بزرگ اختیار شود تا میانگین

<sup>11</sup>Marginal likelihood function

<sup>12</sup>Artificially constructed distribution

<sup>13</sup>DC-based posterior distribution

و ماتریس واریانس مقیاس بندی شده توزیع پسین نتیجه شده پارامترها به MLE و برآورد ماتریس واریانس متناظرشان همگرا شوند. برای مقیاس بندی ماتریس واریانس نیز از ضریب  $k^{-1}$  استفاده می شود. لازم به ذکر است که اگرچه در واقعیت  $k$  تکرار مستقل یک آزمایش مشابه، مجموعه های داده یکسانی را نتیجه نمی دهند و به طور نظری احتمال چنین رخدادی صفر است، اما این کار توسط برنامه های رایانه ای تولید مقادیر تصادفی امکان پذیر است.

توزیع پسین مبتنی بر DC به صورت

$$\pi^{(k)}(\beta, \theta, \mathbf{u} | \mathbf{y}) \propto \pi^{(k)}(\mathbf{u} | \mathbf{y}, \beta, \theta) \pi^{(k)}(\beta, \theta | \mathbf{y})$$

ساخته می شود. بر اساس ویژگی های مجانبی توزیع پسین مبتنی بر DC کناری پارامترها،  $\pi^{(k)}(\beta, \theta | \mathbf{y})$ ، باغیشتی و محمدزاده (۲۰۱۱) نشان دادند اگر  $k \rightarrow \infty$ ، آن گاه

$$\begin{aligned} E^{(k)}(\beta, \theta | \mathbf{y}) &\rightarrow (\hat{\beta}, \hat{\theta}) \\ \text{Var}^{(k)}(\beta, \theta | \mathbf{y}) &\rightarrow k^{-1} \times \text{Var}(\hat{\beta}, \hat{\theta}) \end{aligned}$$

که در آن  $E^{(k)}(\beta, \theta | \mathbf{y})$  و  $\text{Var}^{(k)}(\beta, \theta | \mathbf{y})$  به ترتیب امید ریاضی و واریانس توزیع پسین مبتنی بر DC کناری پارامترهای مدل و  $(\hat{\beta}, \hat{\theta})$  برآوردهای ML پارامترهای  $(\beta, \theta)$  هستند. دقت شود که برای محاسبه توزیع کناری  $\pi^{(k)}(\beta, \theta | \mathbf{y})$  نیازی به انتگرال گیری از توزیع پسین مبتنی بر DC نسبت به  $\mathbf{u}$  نیست و به آسانی با کنار گذاشتن نمونه های تولید شده  $\mathbf{u}$  در الگوریتم MCMC، نتیجه نمونه ای از توزیع کناری مذکور خواهد بود. یکی از مزایای روش DC ناوردایی نتایج به دست آمده نسبت به انتخاب توزیع های پیشین است (لهله و همکاران، ۲۰۰۷).

### ۱.۳ تشخیص تعداد نسخه ها

همان طور که اشاره شد برای اطمینان از همگرایی برآوردهای نتیجه شده با روش DC به MLE، باید تعداد نسخه ها به اندازه کافی بزرگ اختیار شود. اما برای هر موقعیت عملی، سوال چند نسخه؟ مطرح خواهد بود. لهله و همکاران (۲۰۱۰) معیارهایی برای تشخیص همگرایی الگوریتم DC ارائه کردند که بر اساس توزیع پسین توام مبتنی بر DC تعریف می شوند. این معیارها عبارتند از ماکسیمم مقدار ویژه ماتریس کوواریانس توزیع پسین مبتنی بر DC ( $\lambda_{\max}$ )، میانگین توان دوم خطا (MSE) و یک آماره برازش

شبه معیار همبستگی ( $r^2$ ). معیار ماکسیم مقدار ویژه، تباهدگی توزیع پسین مبتنی بر DC در نقاط MLE را نشان می‌دهد؛ در حالی که دو معیار دیگر بیانگر مناسب بودن تقریب نرمال هستند. هر سه معیار با افزایش تعداد نسخه‌ها، یعنی  $k$ ، باید به صفر همگرا شوند. بنا به لاله و همکاران (۲۰۱۰)، مقدار  $k$  که به ازای آن معیارهای معرفی شده کمتر از ۰/۰۵ باشند، مناسب و کافی است.

#### ۴ یک معیار اطلاع مبتنی بر روش DC

همان‌طور که اشاره کردیم، مقدار ماکسیم تابع درستنمایی یک کمیت کلیدی در تعریف معیارهای انتخاب مدل است و روش DC قادر به محاسبه مقدار ماکسیم تابع درستنمایی نیست. بنابراین، با استفاده از این روش نمی‌توان یک معیار انتخاب مدل (مستقیم) در رده‌ای از مدل‌های ممکن تعریف کرد. پونچیانو و همکاران (۲۰۰۹) توانستند الگوریتمی برای محاسبه اختلاف مقادیر معیارهای اطلاع معرفی کنند و مدل‌های نامزد را به صورت جفتی مقایسه کنند. این رهیافت، با توجه به بعد رده مدل‌ها، می‌تواند هزینه‌بر و طاقت‌فرسا باشد. در این بخش، معیاری معرفی می‌شود که بر مبنای روش داده‌تاگی به دست می‌آید و اعضای رده مدل‌ها را به طور مستقیم و تکی رتبه‌بندی می‌کند.

فرض کنید  $g(y)$  تابع چگالی مدل واقعی باشد که داده‌ها از آن تولید شده‌اند و  $f(y, \psi)$  تابع چگالی مدل پارامتری (نامزد) باشد که برای داده‌ها در نظر گرفته شده است. چندین روش برای اندازه‌گیری فاصله تقریب پارامتری  $f(y, \psi)$  تا چگالی واقعی  $g(y)$  وجود دارند. فاصله‌ای که مرتبط با روش ماکسیم درستنمایی است، فاصله کولبک-لیبلر<sup>۱۴</sup> است که به صورت

$$KL(g(\cdot), f(\cdot, \psi)) = \int g(y) \log \frac{g(y)}{f(y, \psi)} dy$$

تعریف می‌شود. لگاریتم تابع درستنمایی با  $\ell_n(\psi) = \log L_n(\psi)$  نشان داده می‌شود، که در آن اندیس  $n$  برای تاکید وابستگی به حجم نمونه است. تابع لگاریتم درستنمایی داده‌های کپی شده به صورت

$$\ell_n^k(\psi) = \log[L_n(\psi)]^k = \log\left[\prod_{i=1}^n f(y_i, \psi)\right]^k$$

<sup>14</sup>Kullback-Leibler distance

تعریف می‌شود. با توجه به استقلال نسخه‌های کپی‌شده،  $g^k(y) = [g(y)]^k$  تابع چگالی مدل واقعی برای داده‌های کپی‌شده است. بنابراین باید به دنبال مدلی بود که فاصله

$$\text{KL}(g^k(\cdot), f^k(\cdot, \psi)) = \int g^k(y) \log \frac{g^k(y)}{f^k(y, \psi)} dy \quad (۲)$$

را می‌نیم کند. برای هر مدل دلخواه  $f(\cdot, \psi)$ ، مولفه  $\int g^k(y) \log g^k(y) dy$  مقداری ثابت است. از طرفی، با استفاده از قانون قوی اعداد بزرگ

$$\frac{1}{n} \ell_n^k(\psi) \xrightarrow{a.s.} \int g^k(y) \log f^k(y, \psi) dy = E_{g^k} \left( \log f^k(y, \psi) \right).$$

بنابراین چون برآوردگر  $\hat{\psi}$  ماکسیم‌کننده  $\ell_n^k(\psi)$  است، تحت شرایط مناسب به  $\psi_0$ ، یعنی می‌نیم‌کننده فاصله (۲)، میل می‌کند، پس

$$\hat{\psi} \xrightarrow{a.s.} \psi_0 = \arg \min \text{KL}(g^k(\cdot), f^k(\cdot, \psi)).$$

در ادامه، چند کمیت تعریف می‌شوند. کمیت‌های

$$u(y, \psi) = \frac{\partial}{\partial \psi} \log f^k(y, \psi), \quad I(y, \psi) = \frac{\partial^2}{\partial \psi \partial \psi^T} \log f^k(y, \psi)$$

به ترتیب بردار امتیاز و ماتریس اطلاع متناظر با داده‌های کپی‌شده نامیده می‌شوند. چون  $\psi_0$  فاصله (۲) را می‌نیم می‌کند، پس

$$E_{g^k} (u(y, \psi_0)) = \int g^k(y) u(y, \psi_0) dy = 0.$$

با قرار دادن

$$J = -E_{g^k} I(y, \psi_0), \quad K = \text{Var}_{g^k} u(y, \psi_0).$$

وقتی  $f^k(y, \psi_0)$  با  $g^k(y)$  یکی باشد،  $J$  و  $K$  با بعد  $\kappa$ ، ماتریس‌های همانی خواهند شد که در آن  $\kappa$



بعد بردار پارامتر  $\psi$  است. با شرایط و تعاریف بالا ثابت می‌شود (هیورت و پلارد، ۱۹۹۳)

$$\hat{\psi} = \psi_0 + J^{-1}\bar{U}_n + o_p(n^{-\frac{1}{2}}) \quad (۳)$$

که در آن  $\bar{U}_n = n^{-1} \sum_{i=1}^n u(y_i, \psi_0)$  از قضیه حد مرکزی

$$\sqrt{n}\bar{U}_n \xrightarrow{d} U' \sim N_{\kappa}(\circ, K)$$

که در ترکیب با (۳) نتیجه می‌شود

$$V_n = \sqrt{n}(\hat{\psi} - \psi_0) \xrightarrow{d} J^{-1}U' = N_{\kappa}(\circ, J^{-1}KJ^{-1}). \quad (۴)$$

با قرار دادن  $\hat{\psi}$  در (۲) و ثابت بودن مولفه  $\int g^k(y) \log g^k(y) dy$ ، کافی است

$$R_n = \int g^k(y) \log f^k(y, \hat{\psi}) dy$$

به دست آورده شود، که یک متغیر تصادفی است و از طریق برآوردگر  $\hat{\psi}$  به داده‌ها وابسته است. به همین دلیل از امید ریاضی آن، تحت تابع چگالی مدل واقعی، استفاده می‌شود. بنابراین

$$Q_n = E_{g^k}(R_n) = E_{g^k}\left(\int g^k(y) \log f^k(y, \hat{\psi}) dy\right).$$

به طور کلی استراتژی معیارهای انتخاب مدل مبتنی بر ماکسیم تابع درستنمایی، برآورد  $Q_n$  برای هر مدل نامزد و انتخاب مدل با بالاترین  $Q_n$  برآورد شده است. این فرآیند معادل با جستجو برای مدلی با کوچک‌ترین فاصله کولبک-لیبلر برآورد شده است. یک حالت ممکن برای برآورد  $Q_n$ ، استفاده از توزیع تجربی داده‌های کپی شده است. در نتیجه

$$\hat{Q}_n = \frac{1}{n} \sum_{i=1}^n \log f^k(y_i, \hat{\psi}) = \frac{1}{n} \ell_n^k(\hat{\psi}).$$

در داده‌های کپی شده، هر مشاهده  $k$  بار تکرار شده است. بنابراین مخرج  $\hat{Q}_n$  برابر  $n$  است نه  $nk$ .

۳۰ ..... یک معیار جدید انتخاب مدل

اکنون  $Z_i = \log f^k(y_i, \psi_0) - Q_0$  قرار داده می‌شود، که در آن

$$Q_0 = \int g^k(y) \log f^k(y, \psi_0) dy.$$

متغیرهای  $Z_i$ ، متغیرهایی با میانگین صفر هستند.  $\bar{Z}_n$  نیز میانگین  $Z_i$ ها تعریف می‌شود. بنابراین می‌توان نشان داد

$$\hat{Q}_n - R_n = \bar{Z}_n + n^{-1} V_n^T J V_n + o_p(n^{-1}). \quad (5)$$

با توجه به (۴) می‌توان نوشت

$$V_n^T J V_n \xrightarrow{d} W = (U')^T J^{-1} J J^{-1} U' = (U')^T J^{-1} U'.$$

بنابراین

$$\begin{aligned} E(\hat{Q}_n - R_n) = E(\hat{Q}_n - Q_n) &= E(\bar{Z}_n + n^{-1} V_n^T J V_n + o_p(n^{-1})) \quad (6) \\ &= E(\bar{Z}_n) + n^{-1} E(W) \approx \frac{\kappa^*}{n} \end{aligned}$$

که در آن  $\kappa^* = E(W)$  اثر<sup>۱۵</sup> ماتریس  $J^{-1} K$  و جمله جریمه در تعریف معیارهای انتخاب مدل است.

تذکره ۱: اگر  $f^k(y, \psi_0) = g^k(y)$ ، آن‌گاه  $K = J$  و بنابراین  $\kappa^* = \kappa$ .

با نتیجه به دست آمده، یک نسخه با اریبی تصحیح شده از برآوردگر خام  $\hat{Q}_n$  به صورت

$$\hat{Q}_n - \frac{\kappa^*}{n} = n^{-1} \{ \ell_n^k(\hat{\psi}) - \kappa^* \}$$

حاصل می‌شود.

<sup>15</sup>Trace

#### ۱.۴ معیار انتخاب مدل پیشنهادی

اگر داده‌ها مستقل باشند، آنگاه درجه آزادی مدل برابر حجم نمونه منهای تعداد پارامترهای برآوردشده مدل است. این تعریف روشن و عام هست و در معرفی جمله جریمه در معیارهای انتخاب مدل معروف مثل AIC و BIC<sup>۱۶</sup> از آن استفاده می‌شود. اما زمانی که داده‌ها وابسته هستند، تعریف روشنی از درجه آزادی مدل وجود ندارد. یعنی مشخص نیست که چه تعداد از داده‌ها برای برآورد پارامترها کافی است. بنابراین، مساله اصلی در معرفی معیارهای انتخاب مدل برای مدل‌های آمیخته، تعریف یک جمله جریمه مناسب است.

کارهای متفاوتی به منظور معرفی یک جریمه مناسب برای AIC در مدل‌های آمیخته انجام شده‌اند. به‌عنوان نمونه می‌توان به وایدا و بلانچارد (۲۰۰۵)، رفتری و همکاران (۲۰۰۷)، گراون و نیب (۲۰۱۰)، و لیانگ و همکاران (۲۰۰۸) اشاره کرد. رفتری و همکاران (۲۰۰۷) با استفاده از میانگین هارمونیک مقادیر تابع درستنمایی به ازای نمونه‌های حاصل از توزیع پسین، برآوردی برای ماکسیمم تابع درستنمایی ارایه دادند. این برآوردگر در عین سادگی، به‌شدت ناپایدار است. آن‌ها برای رفع این مشکل، از توزیع (تقریبی)  $\ell_{\max} - \ell_t$  استفاده کردند و نشان دادند

$$\ell_{\max} - \ell_t \sim \text{Gamma}(\alpha, 1)$$

که در آن  $\{\ell_t; t = 1, \dots, B\}$  مقادیر تابع درستنمایی به ازای  $B$  نمونه تولیدشده از توزیع پسین،  $\ell_{\max}$  مقدار ماکسیمم تابع درستنمایی و  $\alpha = \frac{\kappa}{\bar{\ell}}$  هستند. با توجه به آن که  $E(\ell_{\max} - \ell_t) = \alpha$  و  $\text{Var}(\ell_t) = \alpha$ ، برآوردهای گشتاوری به صورت

$$\hat{\ell}_{\max} = \bar{\ell} + s_{\ell}^2, \quad \hat{\alpha} = s_{\ell}^2$$

نتیجه می‌شوند، که در آن‌ها  $\bar{\ell}$  و  $s_{\ell}^2$  میانگین نمونه و واریانس  $\ell_t$  ها هستند. با توجه به معیار

$$\text{AIC} = 2\ell_{\max} - 2\kappa$$

<sup>16</sup>Bayesian information criterion

رفتاری و همکاران (۲۰۰۷) یک برآورد مونت کارلویی از آن را به صورت

$$AICM = \psi \hat{\ell}_{\max} - \psi \hat{\kappa} = \psi(\bar{\ell} + s_{\ell}^{\psi}) - \psi s_{\ell}^{\psi} = \psi \bar{\ell} - \psi s_{\ell}^{\psi} = \psi(\bar{\ell} - s_{\ell}^{\psi})$$

ارایه دادند، که با الهام از آن و داشتن یک نمونه از توزیع پسین مبتنی بر DC، معیار

$$DCAIC = \psi(\bar{\ell}^k(\hat{\psi}) - s_{\ell^k}^{\psi})$$

یا به طور معادل

$$DCAIC = -\psi \bar{\ell}^k(\hat{\psi}) + \psi s_{\ell^k}^{\psi} \quad (۷)$$

پیشنهاد می‌شود. این معیار نسخه تعمیم‌یافته‌ای از AIC با استفاده از روش داده‌تاگی است. بنابراین کاملاً منطقی است که تمام ویژگی‌های AIC، از جمله سازگاری را به ارث ببرد. یعنی با احتمال متمایل به یک قادر است مدل واقعی را در بین مدل‌های نامزد، در صورت وجود آن، انتخاب کند. همچنین AIC تمایل به انتخاب مدل‌های پیچیده‌تر دارد، یعنی مدل انتخابی توسط این معیار می‌تواند بیش‌برازش<sup>۱۷</sup> باشد. این موضوع برای DCAIC نیز برقرار است. مدل با کمترین مقدار DCAIC، برحسب رابطه (۷)، به‌عنوان مدل برتر انتخاب می‌شود.

## ۵ ارزیابی عملکرد معیار پیشنهادی

برای ارزیابی معیار ارائه‌شده، بردار پاسخ با حجم  $n$  از یک مدل آمیخته خطی تعمیم‌یافته پواسون به صورت

$$Y_i | \lambda_i \sim \text{Poisson}(\lambda_i), \quad i = 1, \dots, n \quad (۸)$$

$$\lambda_i = \exp(\mathbf{x}_i^T \beta + \alpha_i) \quad (۹)$$

$$\alpha_i \sim N(0, \sigma^2)$$

<sup>17</sup>Overfitted

تولید شده است، که در آن  $x_i = (1, x_{1i}, x_{2i}, x_{3i})^T$  و مقادیر واقعی پارامترها

$$(\beta_0, \beta_1, \beta_2, \beta_3, \sigma) = (1, -0.5, 1.5, 0.7, 0.5)$$

انتخاب شدند. متغیرهای تبیینی  $x_1, x_2$  و  $x_3$  به ترتیب از توزیع‌های یکنواخت استاندارد، نرمال استاندارد، و برنولی با احتمال موفقیت 0.5 تولید شده‌اند. سه حجم نمونه مختلف 20، 50، 100 نیز در نظر گرفته شده‌اند. علاوه بر این، تعداد نسخه‌ها برای داده‌های کپی‌شده برابر  $k = 1, 5$  انتخاب شدند. دلیل انتخاب  $k = 5$  برای روش DC، مقادیر نتیجه‌شده معیارهای تشخیص همگرایی است که در جدول 1 برای حجم‌های نمونه مختلف گزارش شده‌اند. تنها برای حجم نمونه  $n = 20$  این انتخاب می‌تواند کمی مورد تردید قرار گیرد. اما مقدار معیار  $r^2$  برای  $k = 5$  در این حالت نزدیک به 0.5 است. همه محاسبات در نرم‌افزار R

جدول 1. مقادیر معیارهای تشخیص همگرایی الگوریتم DC

$r^2$	MSE	$\lambda_{\max}$	$k$	$n$
0.305	35.26	0.813	1	20
0.079	2.642	0.113	5	
0.056	1.122	0.007	10	
0.007	0.228	0.066	1	50
0.002	0.027	0.013	5	
0.001	0.045	0.007	10	
0.002	0.019	0.127	1	100
0.0004	0.004	0.021	5	
0.001	0.029	0.011	10	

(تیم هسته R، 2016) و به کمک بسته dclone (سلیموس، 2010) انجام شده‌اند. برای بررسی توانایی DCAIC در انتخاب مدل واقعی و کنکاش ویژگی بیش‌برازشی آن، دو رده مدل (نادرست) کم‌برازش<sup>18</sup> و بیش‌برازش در نظر گرفته شده‌اند. رده مدل‌های کم‌برازش، شامل سه مدل با پیش‌گوه‌های خطی به صورت

$$M_1 : \eta_i = \beta_0 + \beta_1 x_{1i} + \alpha_i$$

<sup>18</sup>Underfitted

$$M_2 : \eta_i = \beta_0 + \beta_2 x_{2i} + \alpha_i$$

$$M_3 : \eta_i = \beta_0 + \beta_2 x_{2i} + \beta_3 x_{3i} + \alpha_i$$

هستند. برای رده مدل‌های بیش‌برازش، تنها یک مدل با پیش‌گوی خطی

$$M_4 : \eta_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 x_{4i} + \alpha_i$$

در نظر گرفته شد، که در آن  $x_4$  از توزیع یکنواخت استاندارد شبیه‌سازی شده است. در این مطالعه، ۱۰۰ مجموعه داده از مدل واقعی تولید شدند و برای هر مجموعه، هر ۵ مدل درست ( $M_0$ ) و نادرست ( $M_1$ ) تا  $M_4$ ) برازش و معیار DCAIC محاسبه شد. نتایج برای تعداد نسخه‌های ۱ و ۵ در جدول ۲ گزارش شده‌اند که اعداد آن بیانگر تعداد دفعات انتخاب مدل واقعی توسط معیار DCAIC در مقابل مدل‌های رقیب در دو گروه مدل‌های کم‌برازش و بیش‌برازش، در ۱۰۰ مجموعه داده، هستند. برای  $k = 1$  که در واقع فقط از داده‌های اصلی استفاده می‌شود، معیار DCAIC همان معیار پیشنهادی رفتاری و همکاران (۲۰۰۷) خواهد بود که تقریبی از AIC در مدل آمیخته شبیه‌سازی شده است. مشاهده می‌کنید که وقتی مدل رقیب یک مدل کم‌برازش است، با افزایش حجم نمونه، معیار پیشنهادی به انتخاب مدل واقعی تمایل دارد، در حالی‌که برای مدل رقیب بیش‌برازش شانس انتخاب مدل واقعی تقریباً برابر شانس عدم انتخاب آن است. از طرفی، در حالتی که حجم نمونه کوچک است، مدل واقعی معمولاً در مقابل رقبای خود بازنده است. این نتیجه ویژگی بیش‌برازشی معیار DCAIC را که از معیار AIC به ارث می‌برد، به وضوح، نشان می‌دهد. در مقابل زمانی که تعداد نسخه‌های داده به  $k = 5$  افزایش می‌یابد، معیار DCAIC در اغلب موارد برای حالتی که مدل رقیب یک کم‌برازش است، مدل واقعی را انتخاب می‌کند. این ویژگی (سازگاری) با افزایش حجم نمونه نیز قوی‌تر می‌شود. در مقابل، زمانی که مدل رقیب یک بیش‌برازش است، در اغلب موارد، مدل واقعی شکست می‌خورد که نشان از نمود بارزتر ویژگی بیش‌برازشی در این حالت است.

با توجه به آن‌که شالوده یک توزیع پسین مبتنی بر DC بر پایه توزیع‌های پسین بیزی شکل می‌گیرد، بنا بر پیشنهاد یکی از داوران مقاله، مقایسه معیار پیشنهادی با یک معیار انتخاب مدل بیزی مثل معیار اطلاع‌کبیش<sup>۱۹</sup> (DIC) می‌تواند مفید باشد. برای این منظور، در مطالعه شبیه‌سازی تعداد دفعات انتخاب مدل واقعی توسط DIC نیز محاسبه شدند که در جدول ۲ گزارش شده‌اند. همانطور که ملاحظه می‌شود روند عملکرد دو معیار برای زمانی که از داده‌های واقعی استفاده می‌شود، تقریباً مشابه است، اما وقتی

<sup>19</sup>Deviance information criterion

جدول ۰۲. تعداد دفعات انتخاب مدل واقعی با دو معیار

k	مدل	DCAIC			DIC		
		n			n		
		۱۰۰	۵۰	۲۰	۱۰۰	۵۰	۲۰
۱	کم‌برازش	۳۹	۵۱	۷۲	۴۲	۶۳	
	بیش‌برازش	۴۶	۴۷	۵۹	۴۶	۷۳	
۵	کم‌برازش	۷۲	۷۲	۸۶	۶۷	۶۵	
	بیش‌برازش	۱۱	۱۹	۱۶	۴۲	۶۳	

از داده‌های کپی‌شده استفاده می‌شود، دو معیار نتایج متفاوتی را ارائه می‌دهند. در واقع معیار DCAIC در انتخاب مدل‌های کم‌برازش، به‌ویژه در حجم نمونه کوچک، توفیق خیلی بیشتری نسبت به DIC دارد و این نتیجه برای مدل‌های رقیب بیش‌برازش برعکس است. البته نویسندگان از معتبر بودن معیار DIC برای انتخاب مدل بر اساس داده‌های کپی‌شده، مطمئن نیستند. بنابراین مقایسه بین دو معیار برای حالت  $k \neq 1$ ، مورد سوال است.

## بحث و نتیجه‌گیری

در این مقاله، یک معیار اطلاع جدید مبتنی بر روش داده‌تاگی ارائه شد که برای انتخاب مدل در رده مدل‌های آمیخته می‌تواند یک گزینه ممکن دلخواه باشد. عملکرد این معیار، بر اساس یک مطالعه شبیه‌سازی، نشان داد که هر دو ویژگی سازگاری و تمایل به انتخاب مدل‌های پیچیده معیار AIC در ذات این معیار جدید نیز وجود دارند. اما، این معیار نسخه قابل اعتمادتری از نسخه ارائه‌شده توسط رفتی و همکاران (۲۰۰۷) است.

با توجه به ویژگی بیش‌برازشی DCAIC، انتظار می‌رود بتوان نسخه سازواری از آن را که مشابه BIC عمل می‌کند، به‌دست آورد که مشکل بیش‌برازشی را با توجه به جریمه سنگین‌تر موجود در BIC نسبت به AIC، مرتفع کند. امکان‌سنجی بسط DCAIC به مواردی که ویژگی‌های سایر معیارهای اطلاع پرمصرف را دارا باشد، توسط نویسندگان در حال بررسی است و به‌عنوان یک موضوع پژوهشی آینده مطرح است.

## تقدیر و تشکر

نویسندگان از داوران محترم برای بیان نظرات ارزشمندشان که در ارتقای کیفیت مقاله موثر بودند، قدردانی می‌کنند.

## مراجع

- Baghishani, H. and Mohammadzadeh, M. (2011), A Data Cloning Algorithm for Computing Maximum Likelihood Estimates in Spatial Generalized Linear Mixed Models, *Computational Statistics and Data Analysis*, **55**, 1748-1759.
- Berger, J. O., Bernardo, J. M. and Sun, D. (2009), The Formal Definition of Reference Priors, *The Annals fo Statistics*, **37**, 905-938.
- Breslow, N. and Clayton, D. G. (1993), Approximate Inference in Generalized Linear Mixed Models, *Jornal of the American Statistical Association*, **88**, 9-25.
- Christensen, O. F. (2004), Monte Carlo Maximum Likelihood in Model-Based Geostatistics, *Journal of Computational and Graphical Statistics*, **13**, 702-718.
- Christensen, O. F. and Waagepetersen, R. P. (2002), Bayesian Prediction of Spatial Count Data Using Generalized Linear Mixed Models, *Biometrics*, **58**, 280-286.
- Christensen, O. F., Roberts, G. O. and Skold, M. (2006), Robust Markov Chain Monte Carlo Methods for Spatial Generalized Linear Mixed Models, *Journal of Computational and Graphical Statistics*, **15**, 1-17.
- Fong, Y., Rue, H. and Wakefield, J. (2010), Bayesian Inference for Generalized Linear Mixed Models, *Biostatistics*, **11**, 397-412.
- Greven, S. and Kneib, T. (2010), On the Behavior of Marginal and Conditional Akaike Information Criteria in Linear Mixed Models, *Biometrika*, **97**, 773-789.
- Hjort, N. L. and Pollard, D. B. (1993), Asymptotics for Minimisers of Convex Processes, *Statistical Research Report*, Department of Mathematics, University of Oslo.
- Hosseini, F., Eidsvik, J. and Mohammadzadeh, M. (2011), Approximate Bayesian Inference for Generalized Linear Mixed Models with Skew Normal Latent Variables, *Computational Statistics and Data Analysis*, **55**, 1791-1806.



- Lele, S. R., Dennis, B. and Lutscher, F. (2007), Data Cloning: Easy Maximum Likelihood Estimation for Complex Ecological Models Using Bayesian Markov Chain Monte Carlo Methods, *Ecology Letters*, **10**, 551–563.
- Lele, S. R., Nadeem, K. and Schmuland, B. (2010), Estimability and Likelihood Inference for Generalized Linear Mixed Models Using Data Cloning, *Journal of the American Statistical Association*, **105**, 1617-1625.
- Liang, H., Wu, H. and Zou, G. (2008), A Note on Conditional AIC for Linear Mixed-effects Models, *Biometrika*, **95**, 773-778.
- Pinheiro, J. C. and Bates, D. M. (1995), Approximations to the Log-Likelihood Function in the Nonlinear Mixed-Effects Model, *Journal of Computational and Graphical Statistics*, **4**, 12-35.
- Ponciano, J. M., Taper, M. L., Dennis, B. and Lele, S. R. (2009), Hierarchical Models in Ecology: Confidence Intervals, Hypothesis Testing, and Model Selection Using Data Cloning, *Ecology*, **90**, 356–362.
- Raftery, A. E., Newton, M. A., Satagopan, J. M. and Krivitsky, P. N. (2007), Estimating the Integrated Likelihood via Posterior Simulation Using the Harmonic Mean Identity, *In. Bernardo et al. (eds) Bayesian Statistics*, Oxford University Press.
- R Core Team (2016), *R: a language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Solymos P. (2010), dclone: Data Cloning in R. *The R Journal*, **2**, 29-37.
- Torabi, M., Lele, S. R. and Parsad, N. G. N. (2015), Likelihood Inference for Small Area Estimation Using Data Cloning, *Computational Statistics and Data Analysis*, **89**, 158-171.
- Vaida, F. and Blanchard, S. (2005), Conditional Akaike Information for Mixed-Effects Models, *Biometrika*, **92**, 351-370.