



Shrinkage Estimators in Semi-Parametric Heteroscedastic Hierarchical Models with Restricted Joint Empirical likelihood

Shantia, V.¹ , Ghoreishi, S. K.² 

¹Department of Statistics, Science and Research Branch, Islamic Azad University, Tehran, Iran.

²Department of Statistics, University of Qom, Qom, Iran.

Corresponding author: S. K. Ghoreishi, atty_ghoreishi@yahoo.com

Received: 16/6/2023 **Revised:** 4/5/2024 **Accepted and Published Online:** 6/5/2024.

Introduction

Applying shrinkage estimators in hierarchical models has garnered significant interest in statistical practice, with seminal contributions dating back to James and Stein (1961). Various methodologies, including the moment method, maximum likelihood method, and unbiased risk method, have been proposed for parameter estimation. Along with homoscedastic, heteroscedastic hierarchical models have been explored, with the latter gaining attention among statisticians, Xie et al. (2012, 2016). Despite advancements, a key challenge remains the normality assumption validity for the second level of the hierarchical models. This paper introduces semi-parametric hierarchical models and employs the Restricted Joint Empirical Likelihood method to examine the impact of distribution dispersion at the second level. We demonstrate the superiority of shrinkage estimates derived from the restricted joint empirical likelihood method through simulation studies, particularly in scenarios with outlier data and heavy-tailed distributions.

Material and Methods

The paper outlines the definition of semi-parametric heteroscedastic hierarchical models $X_i | \theta_i \sim N(\theta_i, A_i)$, $\theta_i \sim \pi_{\mu, \lambda}(\theta_i)$, $i = 1, \dots, n$. The structure of shrinkage estimations within these models $\hat{\theta}_i^S = \frac{\hat{\lambda}^S}{\hat{\lambda}^S + A_i} X_i + \frac{A_i}{\hat{\lambda}^S + A_i} \hat{\mu}^S$, and introduces the Restricted Joint Empirical Likelihood method for estimating hyper-parameters. A simulation study is conducted to evaluate the

performance of the Restricted Joint Empirical Likelihood $L_R(\mu, \lambda)$ in hyper-parametric estimation compared to existing methods for the following four scenarios in the presence of outliers

1. $X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim N(\mu, \lambda).$
2. $X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim Laplace(\mu, 2\lambda).$
3. $X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim U(\mu - \sqrt{3}\lambda, \mu + \sqrt{3}\lambda).$
4. $X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim N(20, \lambda), i = 1, \dots, [\frac{n}{10}],$ and $\theta_i \sim N(\mu, \lambda),$
 $i = [\frac{n}{10}] + 1, \dots, n,$

Results and Discussion

Simulation results indicate that the shrinkage estimation obtained from the restricted joint empirical likelihood method outperforms existing methods, such as the moment method and Stein’s unbiased risk estimator, particularly in the presence of outlier data. The method’s non-parametric nature enhances its applicability to diverse real-world datasets.

Conclusion

This study introduces semi-parametric heteroscedastic hierarchical models and proposes the Restricted Joint Empirical Likelihood method, for obtaining shrinkage estimators in such models. Simulation studies demonstrate this method’s superior performance, highlighting its effectiveness in analyzing semi-parametric heterogeneity hierarchical models. Moreover, applying this approach to real data analysis, exemplified by the Iran Airlines Flight Delay dataset, underscores its practical relevance and usefulness.

Keywords: Moment Estimator, Stein’s Unbiased Risk Estimator, Estimating Equations, Empirical Maximum Likelihood Estimator.

Mathematics Subject Classification (2010): 62F15, 62J07.





مجله علوم آماری، بهار و تابستان ۱۴۰۳

جلد ۱۸، شماره ۱، ص ۹۱ -- ۱۰۲

DOI: 10.52547/jss.18.1.06

مقاله پژوهشی

برآوردهای انقباضی در مدل‌های سلسله مراتبی نیم-پارامتری خطی با تابع درستنمایی تجربی توأم مقید

ویدا شنتیاء^۱، سید کامران قریشی^۲

گروه آمار، دانشگاه آزاد اسلامی، واحد علوم و تحقیقات تهران

^۲ گروه آمار، دانشگاه قم

نویسنده مسئول: سید کامران قریشی، atty_ghoreishi@yahoo.com

تاریخ دریافت: ۱۴۰۲/۳/۲۶ تاریخ بازنگری: ۱۴۰۳/۲/۱۵ تاریخ پذیرش و انتشار: ۱۴۰۳/۲/۱۷

چکیده: در این مقاله ابتدا مدل‌های سلسله مراتبی نیم-پارامتری معرفی می‌شود. سپس با ارائه نسخه جدیدی از تابع درستنمایی تجربی (تابع درستنمایی تجربی توأم مقید)، از آن برای برآورد پارامترهای انقباضی در مدل‌های سلسله مراتبی نیم-پارامتری استفاده خواهد شد. تحت فرض‌های مختلف کارآمدی استفاده از تابع درستنمایی تجربی توأم مقید در تحلیل مدل‌های سلسله مراتبی نیم-پارامتری با یک مطالعه شبیه‌سازی بررسی می‌گردد. همچنین از روش معرفی شده در این مقاله برای تحلیل داده‌های تعداد تأخیر پروازهای شرکت‌های مختلف هواپیمایی استفاده خواهد شد. **واژه‌های کلیدی:** برآوردگر گشتاوری، برآوردگر مخاطره ناریب اشتاین، معادلات برآورد ساز، برآوردگر ماکسیم درستنمایی تجربی
کد موضوع بندی ریاضی (۲۰۱۰): 62F15، 62J07.

۱ مقدمه

استفاده از برآوردهای انقباضی^۱ برای برآورد پارامترها در مدل‌های سلسله مراتبی همواره مورد توجه آمارشناسان بوده است. **جیمز و اشتاین (۱۹۶۱)** و **نیز اشتاین (۱۹۶۲)** پیش از دیگران مدل‌های سلسله مراتبی همگن را توسعه دادند.



©نویسندگان). ناشر انجمن آمار ایران است.

این مقاله با دسترسی آزاد تحت شرایط و ضوابط (CC BY-NC 4.0) توزیع شده است.

¹Shrinkage

تعیین ابرپارامترها در برآورد انقباضی اهمیت ویژه‌ای دارد، لذا روش‌های آماری مختلفی برای برآورد این کمیت‌ها توسط جیمز و اشتاین پیشنهاد شده است که از آن جمله می‌توان به روش گشتاوری، روش ماکسیمم درست‌نمایی، روش برآورد بر اساس برآورد نااریب مخاطره اشاره نمود. به دلیل محدودیت استفاده از مدل‌های سلسله مراتبی همگن در تحلیل داده‌های واقعی، مدل‌های سلسله مراتبی ناهمگن نیز توسعه یافتند. در این راستا خواص مجانبی (با مینیمم مخاطره) برآوردهای انقباضی در مدل‌های سلسله مراتبی ناهمگن، مشابه با مدل‌های سلسله مراتبی همگن، مورد توجه آمارشناس قرار گرفت که در این بین می‌توان به تحقیقات شی و همکاران (۲۰۱۶)، قریشی (۲۰۱۷) و کرمی و آرشی (۱۳۹۳) اشاره نمود.

اگرچه برآوردهای انقباضی در مدل‌های سلسله مراتبی برای همه توزیع‌های متعلق به خانواده نمایی قابل استفاده است اما بیشتر تحقیقات صورت گرفته بر اساس مدل‌های سلسله مراتبی نرمال است که برای نمونه‌های X_1, \dots, X_n به صورت

$$X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim N(\mu, \lambda), \quad i = 1, \dots, n, \quad (1)$$

داده می‌شود، با این فرض که A_1, \dots, A_n کمیت‌های معلوم است. در این صورت برآورد انقباضی θ_i ، متناظر با مدل سلسله مراتبی دو سطحی (۱)، به صورت

$$\hat{\theta}_i = \frac{\hat{\lambda}}{\hat{\lambda} + A_i} X_i + \frac{A_i}{\hat{\lambda} + A_i} \mu, \quad (2)$$

حاصل می‌شود که در واقع میانگین پسین توزیع شرطی θ_i به شرط X_i است. بارانچیک (۱۹۷۰) برآوردهای مینیماکس قابل قبولی برای θ_i ها ارائه داد. همچنین یک کلاس از خواص برآوردهای مینیماکس بیز توسط استرادرن (۱۹۷۱) ارائه شد. براون (۱۹۷۱) یک شرط کافی برای قابل قبول بودن برآوردهای بیز مینیماکس تعمیم‌یافته را مورد بررسی قرار داد. مقایسه بین انواع برآوردها، تحت توابع زیان مختلف، توسط آمارشناسان گوناگون مورد بحث قرار گرفت که از آن جمله می‌توان به برگر و استرادرن (۱۹۹۶) و منابع داخل آن اشاره نمود.

در راستای پرداختن به خواص برآوردهای انقباضی (۶)، نحوه برآورد ابرپارامترهای μ و λ نیز در عمل از اهمیت فراوانی برخوردار است. این مهم توسط شی و همکاران (۲۰۱۲) با مقایسه برآوردهای مختلفی که از روش‌های مختلف به دست می‌آمدند انجام شد. آنها ثابت کردند که برآوردهای SURE ابرپارامترها منجر به مینیمم مخاطره مجانبی برای برآوردهای انقباضی (۶) می‌شوند. برای بررسی خواص گوناگون دیگر از برآوردهای انقباضی مذکور شامل برآوردهای انقباضی در خانواده توزیع‌های با واریانس درجه دو (از میانگین) و روش بیزی با استفاده از چگالی پیشین دیریکله-لاپلاس می‌توان به منابع شی و همکاران (۲۰۱۶) و قریشی (۲۰۱۷) اشاره نمود. با توجه به تمام تحقیقاتی که تاکنون در ارتباط با بررسی خواص برآوردهای انقباضی در مدل‌های سلسله مراتبی (۱) صورت گرفته است، از مشکلاتی که دقت برآوردها را تحت تأثیر قرار می‌دهد پذیرش فرض نرمال بودن برای پارامترهای θ_i در

سطح دوم این مدل هاست. به عبارتی هرگاه ضریب کشیدگی θ_i نسبت به ضریب متناظر در توزیع نرمال کمتر یا بیشتر باشد در این صورت خواص برآوردهای انقباضی دستخوش تغییر می‌شوند. از اینرو در این مقاله به بررسی تأثیر کم یا زیاد پراکنش توزیع سطح دوم مدل سلسله مراتبی (۱) با تعریف مدل‌های سلسله مراتبی نیم-پارامتری^۱ پرداخته و از روش ماکسیمم درست‌نمایی تجربی توأم مقید^۲ (REML) به عنوان روش مکمل برای برآورد ابرپارامترها μ و λ استفاده می‌شود. همچنین به روش شبیه‌سازی نشان داده می‌شود که برای داده‌های دور افتاده، برآوردهای انقباضی حاصل از روش REML عملکرد بهتری نسبت به روش‌های موجود نظیر روش گشتاوری و روش SURE خواهند داشت.

بخش ۲ به معرفی مدل‌های سلسله مراتبی نیم-پارامتری اختصاص دارد. در بخش ۳ برآوردهای انقباضی پارامترهای θ_i تحت سه روش گشتاوری، SURE و REML به دست می‌آیند. بخش ۴ به یک مطالعه شبیه‌سازی با سناریوهای مختلف اختصاص دارد. در بخش ۵ از روش معرفی شده در این مقاله برای تحلیل داده‌های «تعداد تأخیر در پروازهای شرکت‌های هواپیمایی» استفاده خواهد شد. بخش آخر به نتیجه‌گیری اختصاص یافته است.

۲ مدل‌های سلسله مراتبی نیم-پارامتری

مدل‌های سلسله مراتبی نیم-پارامتری، معادل مدل سلسله مراتبی (۱)، به صورت

$$X_i | \theta_i \sim N(\theta_i, A_i), \quad \theta_i \sim \pi_{\mu, \lambda}(\theta_i), \quad i = 1, \dots, n, \quad (3)$$

تعریف می‌شود که در آن توزیع دلخواه در سطح دوم مدل سلسله مراتبی است. در این مدل فرض می‌شود θ_i ها به ترتیب دارای میانگین و واریانس ثابت μ و λ هستند. همچنین فرض می‌کنیم A_1, \dots, A_n کمیت‌های معلومی باشند. به راحتی می‌توان تحقیق کرد که توزیع حاشیه‌ای X_i نامعلوم با میانگین μ و $E(X_i) = \mu$ و واریانس $\text{Var}(X_i) = \lambda + A_i$ است. دو چالش در کاربرد مدل سلسله مراتبی (۳) وجود دارد: الف- ساختار برآورد انقباضی θ_i در عمل نامعلوم است. ب- به دلیل نامعلوم بودن توزیع حاشیه‌ای X_i ، از بعضی روش‌های موجود، مانند روش ماکسیمم درست‌نمایی، نمی‌توان برای برآورد ابر پارامترهای μ و λ استفاده نمود.

در این مقاله از برآورد انقباضی (۶) برای برآورد پارامترهای θ_i استفاده می‌شود. با تأکید بر اینکه این برآوردها تنها وقتی قابل استفاده هستند که هر دو سطح مدل سلسله مراتبی نرمال باشند. با این وجود برای مدل‌های سلسله مراتبی نیم-پارامتری نیز ترجیح داده می‌شود که مجدداً از این برآوردها برای برآورد θ_i استفاده شود، زیرا این برآوردها دارای ترکیب موزونی از مشاهده‌های X_i و میانگین θ_i ها (یعنی μ) هستند. همچنین برای برآورد ابر پارامترهای μ و λ علاوه بر دو روش گشتاوری و SURE از روش REML استفاده خواهد شد. علاوه بر اینکه در

¹Semi-parametric hierarchical models

²Restricted Empirical Maximum Likelihood

بخش شبیه‌سازی، کارآمدی روش REML، در مقایسه با روش‌های گشتاوری و SURE مورد بحث و ارزیابی قرار می‌گیرد. دلیل انتخاب سه روش گشتاوری، SURE و REML برای برآورد ابر پارامترهای μ و λ به شرح زیر است:

- ۱- هر سه روش مذکور برای مدل‌های سلسله مراتبی نیم-پارامتری (۳) قابل استفاده‌اند.
- ۲- روش REML برای زمانی که داده‌های X_i پراکندگی زیاد دارند و حتی زمانی که نمودار بافت‌نگار متناظر آنها متقارن نیست، برآوردهای بهتری نسبت به روش‌های گشتاوری و SURE ارائه می‌کند.
- ۳- اگر چه سه روش مذکور برای مدل‌های سلسله مراتبی ناپارامتری، که در آن حتی توزیع سطح اول X_i نیز نرمال نباشد، کاربرد دارد، لیکن این مقاله بر فرض نرمال برای مشاهده‌های X_i و یک توزیع منعطف برای θ_i تمرکز دارد. این فرض برای تحلیل داده‌هایی که بر اساس یک تبدیل مناسب نرمال شده‌اند قابل استفاده‌اند، براون (۲۰۰۸).

۳ برآوردهای انقباضی و خواص آنها

سه روش گشتاوری، SURE و REML برای برآورد ابر پارامترهای مدل (۳) و در نهایت برآورد پارامترهای انقباضی (۶) معرفی می‌شوند.

الف- روش گشتاوری: در مدل سلسله مراتبی نیم-پارامتری (۳) داریم

$$\begin{aligned} E(X_i) &= E[E(X_i|\theta_i)] = \mu \\ \text{Var}(X_i) &= E[\text{Var}(X_i|\theta_i)] + \text{Var}[E(X_i|\theta_i)] = A_i + \lambda \end{aligned} \quad (۴)$$

در نتیجه برآوردهای گشتاوری λ و μ به ترتیب به صورت $\hat{\lambda}^M = \left[\frac{\sum_{i=1}^n [(X_i - \bar{X})^2 - A_i]}{n} \right]_+$ و $\hat{\mu}^M = \bar{X}$ خواهند بود، که در آن $[a]_+$ برابر a است هرگاه $a > 0$ و برابر 0 است اگر $a \leq 0$. راحتی در محاسبات از امتیاز برجسته این روش آماری برای برآورد ابر پارامترها است. با توجه به برآوردهای فوق برآورد انقباضی متناظر θ_i به صورت $\hat{\theta}_i^M = \frac{\hat{\lambda}^M}{\hat{\lambda}^M + A_i} X_i + \frac{A_i}{\hat{\lambda}^M + A_i} \bar{X}$ خواهد بود که در واقع میانگین پسین توزیع شرطی θ_i به شرط X_i می‌باشد.

ب- روش SURE: با تابع زیان توان دوم $L(\theta, \hat{\theta}) = \frac{1}{n} \sum_{i=1}^n (\theta - \hat{\theta})^2$ تابع مخاطره متناظر با

$$R(\theta, \hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \frac{A_i}{(A_i + \lambda)^2} (A_i(\theta_i - \mu)^2 + \lambda^2),$$

است و برآورد ناریب این تابع مخاطره به صورت

$$\text{SURE}(\mu, \lambda) = \frac{1}{n} \sum_{i=1}^n \frac{A_i}{(A_i + \lambda)^2} (A_i(X_i - \mu)^2 + \lambda^2 - A_i^2).$$

حاصل می‌شود. با مینیمم کردن $SURE(\mu, \lambda)$ نسبت به μ و λ ، برآوردهای SURE ابرپارامترها به دست می‌آید. در این صورت می‌توان برآورد انقباضی θ_i را به صورت $\hat{\theta}_i^S = \frac{\hat{\lambda}^S}{\hat{\lambda}^S + A_i} X_i + \frac{A_i}{\hat{\lambda}^S + A_i} \hat{\mu}^S$ ارائه داد. به دلیل اینکه سطح اول مدل سلسله مراتبی (۳) نرمال است می‌توان نتیجه گرفت که در این مدل‌ها برآوردهای $\hat{\theta}_i^S$ دارای خاصیت بینه مخاطره، تحت شرایط عمومی، هستند (شی و همکاران، ۲۰۱۲).

ج- روش REML: با در نظر گرفتن مدل سلسله مراتبی (۳) و با توجه به روابط (۴)، توابع برآورد ساز^۱ ابر پارامترهای μ و λ به صورت $g_1(X_i; \mu, \lambda) = (X_i - \mu)^\top - A_i - \lambda$ و $g_2(X_i; \mu, \lambda) = (X_i - \mu)$ داده می‌شوند و برای این دو برآورد ساز داریم $\sum_{i=1}^n p_i g_j(X_i; \mu, \lambda) = 0$ ، $j = 1, 2$ ، که در آن p_i ها احتمال توزیع تجربی مقید هستند که با سه شرط $\sum_{i=1}^n p_i = 1$ ، $p_i > 0$ و $\sum_{i=1}^n p_i g_j(X_i; \mu, \lambda) = 0$ از ماکسیمم کردن تابع درستنمایی تجربی مقید^۲ $L(\mu, \lambda) = \prod_{i=1}^n p_i$ به صورت $L(\mu, \lambda) = \prod_{i=1}^n [1 + w_1 g_1^\top(X_i; \mu, \lambda) + w_2 g_2^\top(X_i; \mu, \lambda)]^{-1}$ به دست می‌آیند (بیاتی و همکاران، ۲۰۲۱)، به عبارتی تابع درستنمایی تجربی مقید برابر خواهد بود با

$$L(\mu, \lambda) = \prod_{i=1}^n [1 + w_1 g_1^\top(X_i; \mu, \lambda) + w_2 g_2^\top(X_i; \mu, \lambda)]^{-1}. \quad (5)$$

از آنجا که در تابع درستنمایی تجربی مقید (۵) از ترکیب $1 + w_1 g_1^\top(X_i; \mu, \lambda) + w_2 g_2^\top(X_i; \mu, \lambda)$ استفاده می‌شود، ممکن است در آن اثر یکی از توابع برآورد ساز (در اینجا g_1 به عنوان تابعی از پارامتر μ) تحت تأثیر تابع برآورد ساز دیگر (g_2 به عنوان تابعی از میانگین μ و واریانس‌های حاشیه‌ای $\lambda + A_i$) قرار گیرد و در عمل نقشی در برآورد ابرپارامترها نداشته باشد. بنابراین تابع درستنمایی تجربی جدید به صورت

$$L_R(\mu, \lambda) = \prod_{i=1}^n [1 + w_1 g_1^\top(X_i; \mu, \lambda)]^{-1} [1 + w_2 g_2^\top(X_i; \mu, \lambda)]^{-1} \times \prod_{i=1}^n [1 + w_1 g_1^\top(X_i; \mu, \lambda)]^{-1} \prod_{i=1}^n [1 + w_2 g_2^\top(X_i; \mu, \lambda)]^{-1}, \quad (6)$$

تعریف می‌شود که با استفاده از حاصل ضرب $(1 + w_1 g_1^\top(X_i; \mu, \lambda))(1 + w_2 g_2^\top(X_i; \mu, \lambda))$ نقش هر تابع برآورد ساز در برآورد پارامترها به طور مجزا شرکت داده می‌شود. چون تابع درستنمایی (۶) متشکل از ضرب دو تابع درستنمایی تجربی مقید جزئی است، به آن تابع درستنمایی تجربی توأم مقید گوییم.

قضیه ۰۱. برای توابع درستنمایی (۵) و (۶) برای هر μ و λ داریم، $L_R(\mu, \lambda) \leq L(\mu, \lambda)$.

برهان: با فرض $a_{1i} = w_1 g_1^\top(X_i; \mu, \lambda)$ و $a_{2i} = w_2 g_2^\top(X_i; \mu, \lambda)$ و استفاده از نامساوی $(1 + a_{1i})(1 + a_{2i})$

¹Estimating functions

²Restricted empirical likelihood

$$a_{\gamma_i} \geq 1 + a_{\gamma_i} + a_{\gamma_i} \text{ قضیه ثابت می‌شود.}$$

طبق نامساوی در قضیه ۱ مقادیری از μ و λ که تابع $L_R(\mu, \lambda)$ را ماکسیم می‌کنند تابع $L(\mu, \lambda)$ را نیز به سمت مقدار ماکسیم خود سوق می‌دهند.

برآوردهای ماکسیم درست‌نمایی تجربی μ و λ به صورت $(\hat{\mu}^E, \hat{\lambda}^E) = \arg \min_{\mu, \lambda} L_R(\mu, \lambda)$ حاصل می‌شوند. جذابیت استفاده از روش REML در این است که برای داده‌های پرت احتمال کوچک‌تری برآورد می‌شود. به عبارتی تأثیر داده‌های پرت در برآورد پارامترهای مدل کمتر می‌شود (بیاتی و همکاران، ۲۰۲۱). با این وجود یکی از چالش‌های اساسی در محاسبه برآوردهای ابر پارامترهای ماکسیم درست‌نمایی تجربی $(\hat{\mu}^E, \hat{\lambda}^E)$ ، وابستگی آنها به کمیت‌های ثابت w_1 و w_2 است که مجهول هستند. در عمل می‌توان به دو روش بیزی (با فرض چگالی پیشین برای w_1 و w_2) و روش کلاسیک، اعتبارسنجی متقابل^۱، این دو کمیت ثابت را برآورد نمود. پس از محاسبه برآورد ابر پارامترها، برآورد انقباضی θ_i به روش REML برابر $\hat{\theta}_i^E = \frac{\hat{\lambda}^E}{\hat{\lambda}^E + A_i} X_i + \frac{A_i}{\hat{\lambda}^E + A_i} \hat{\mu}^E$ خواهد بود.

قضیه ۲. تحت شرایط نظم برای توابع $g_1(X_i; \mu, \lambda)$ و $g_2(X_i; \mu, \lambda)$ ، توزیع مجانبی توأم $(\hat{\mu}^E, \hat{\lambda}^E)$ نرمال دو متغیره است.

برهان: از درست‌نمایی $L_R(\mu, \lambda) = \prod_{i=1}^n [1 + w_1 g_1^*(X_i; \mu, \lambda)]^{-1} [1 + w_2 g_2^*(X_i; \mu, \lambda)]^{-1}$ داریم:

$$\ell(\mu, \lambda) = \ln L_R(\mu, \lambda) = - \sum_{i=1}^n [\ln(1 + w_1 g_1^*(X_i; \mu, \lambda)) + \ln(1 + w_2 g_2^*(X_i; \mu, \lambda))].$$

با محاسبه ماتریس هسین متناظر با تابع $\ell(\mu, \lambda)$ و استفاده از قضیه حد مرکزی، توزیع $(\hat{\mu}^E, \hat{\lambda}^E)$ برای $n \rightarrow \infty$ نرمال دو متغیره با میانگین (μ, λ) و ماتریس کوواریانس زیر خواهد بود

$$I^{-1}(\mu, \lambda) = \begin{pmatrix} \frac{\partial^2 \ell(\mu, \lambda)}{\partial \mu^2} & \frac{\partial^2 \ell(\mu, \lambda)}{\partial \mu \partial \lambda} \\ \frac{\partial^2 \ell(\mu, \lambda)}{\partial \mu \partial \lambda} & \frac{\partial^2 \ell(\mu, \lambda)}{\partial \lambda^2} \end{pmatrix}^{-1}$$

^۱Cross validation

۴ مطالعه شبیه‌سازی

در بخش قبل سه روش گشتاوری، SURE و REML برای یافتن برآوردهای انقباضی در مدل‌های سلسله مراتبی نیم-پارامتری استفاده گردید. در آنجا توضیح داده شد که برآوردگر SURE دارای خاصیت مینیمم مخاطره است. اما اینکه در عمل کدامیک از روش‌های فوق عملکرد بهتری را به نمایش می‌گذارد نیاز به تحقیقات وسیع‌تری دارد که در این بخش به روش شبیه‌سازی به آن پرداخته خواهد شد. در اینجا چهار حالت مختلف زیر در نظر گرفته شد.

حالت اول: فرض می‌شود علاوه بر سطح اول، توزیع سطح دوم مدل سلسله مراتبی (۳) نیز نرمال است، یعنی $\theta_i \sim N(\mu, \lambda)$ و $X_i | \theta_i \sim N(\theta_i, A_i)$.

حالت دوم: توزیع سطح دوم مدل سلسله مراتبی (۳)، نسبت به توزیع نرمال پراکندگی بیشتری دارد و از توزیع لاپلاس پیروی می‌کند، به عبارتی داریم $X_i | \theta_i \sim N(\theta_i, A_i)$ و $\theta_i \sim L(\mu, 2\lambda)$.

حالت سوم: توزیع سطح دوم مدل سلسله مراتبی (۳) نسبت به توزیع نرمال پراکندگی کمتری داشته و از توزیع یکنواخت پیروی می‌کند، یعنی $X_i | \theta_i \sim N(\theta_i, A_i)$ و $\theta_i \sim U(\mu - \sqrt{3}\lambda, \mu + \sqrt{3}\lambda)$.

حالت چهارم: با فرض آنکه X دارای داده پرت است، ابتدا توزیع سطح دوم مدل سلسله مراتبی (۳) نرمال فرض می‌شود سپس ۱۰ درصد از داده‌ها با مقدار ثابت ۲۰ جایگزین می‌گردد، به عبارتی داریم $X_i | \theta_i \sim N(\theta_i, A_i)$ که در آن

$$\begin{aligned} \theta_i &\sim N(20, \lambda), & i = 1, \dots, \lfloor \frac{n}{10} \rfloor \\ \theta_i &\sim N(\mu, \lambda), & i = \lfloor \frac{n}{10} \rfloor + 1, \dots, n, \end{aligned}$$

و [۰] نشان دهنده جزء صحیح است.

به منظور مقایسه حالت‌های مختلف، مطالعه شبیه‌سازی برای حجم نمونه‌های ۵۰۰، ۱۰۰۰ و ۵۰۰۰ انجام شد. در هر مرحله شبیه‌سازی از مقادیر $A_i \sim U(0.1, 1)$ ، $\mu = 2$ ، $\lambda = 0.4$ ، $w_1 = w_2 = 0.5$ برای استخراج نتایج استفاده گردید. شبیه‌سازی به تعداد $M = 500$ مرتبه انجام شد و برای ارزیابی کارآمدی روش‌های مختلف، از معیار $MSE = \frac{1}{n} \sum_{i=1}^n (\hat{\theta}_i - \theta_i)^2$ استفاده شد. در هر مرحله از شبیه‌سازی MSE مقادیر واقعی θ_i با برآورد اراکل^۱ متناظر، $\hat{\theta} = \frac{\lambda}{\lambda + A_i} Y_i + \frac{A_i}{\lambda + A_i} \mu$ ، که دارای کمترین مقدار میانگین توان دوم خطا است، نیز به دست آمد. این مقدار برای ارزیابی دقیق‌تر روش‌های مختلف برآورد مورد نیاز است. نتایج مطالعه شبیه‌سازی در جدول ۱ آمده است. در این جدول به ازای n ‌های مختلف مقدار MSE برآوردهای انقباضی $\hat{\theta}_i$ تحت حالت‌های مختلف گزارش شده است. براساس نتایج این جدول ملاحظه می‌شود تحت حالت‌های دوم و چهارم که داده‌ها پراکندگی بیشتری نسبت به توزیع نرمال دارند برآوردهای ماکسیمم درست‌نمایی تجربی عملکرد بهتری نسبت به برآوردهای حاصل از روش‌های مختلف، به‌ویژه زمانی که داده‌ها شامل تعدادی داده‌های پرت هستند، دارند.

¹Oracle

جدول ۰۱. مقادیر MSE برای حالت‌ها و روش‌های مختلف

حالت	روش	n		
		۵۰۰	۱۰۰۰	۵۰۰۰
اول	Oracle	۰/۳۲۱۴	۰/۳۳۹۸	۰/۳۲۵۱
	SURE	۰/۴۵۹۳	۰/۴۶۲۸	۰/۴۶۳۶
	REML	۰/۳۹۸۹	۰/۴۱۸۵	۰/۴۲۴۳
	گشتاوری	۰/۴۴۹۷	۰/۴۵۸۳	۰/۴۷۲۹
دوم	Oracle	۰/۳۱۷۱	۰/۳۰۶۸	۰/۳۴۶۱
	SURE	۰/۴۵۹۴	۰/۴۶۰۳	۰/۴۶۳۹
	REML	۰/۳۶۳۴	۰/۳۹۱۰	۰/۴۳۸۱
	گشتاوری	۰/۴۳۷۷	۰/۴۴۹۰	۰/۴۷۵۸
سوم	Oracle	۰/۳۱۷۰	۰/۳۰۶۹	۰/۳۴۶۱
	SURE	۰/۴۵۹۹	۰/۴۶۰۸	۰/۴۶۳۹
	REML	۰/۴۰۶۵	۰/۴۰۹۴	۰/۴۱۴۰
	گشتاوری	۰/۴۴۰۹	۰/۴۵۰۸	۰/۴۷۶۵
چهارم	Oracle	۰/۲۷۲۹	۰/۲۳۰۵۶	۰/۳۴۰۸
	SURE	۰/۰۱۱۳	۰/۱۱۳	۰/۰۱۱۲
	REML	۰/۰۰۴۲	۰/۰۰۴۲	۰/۰۰۴۲
	گشتاوری	۰/۰۱۱۳	۰/۰۱۱۱	۰/۰۱۱۲

۵ مثال کاربردی

طبق گزارش خبرگزاری جمهوری اسلامی ایران میزان تأخیر بیست شرکت هواپیمایی در چهارماه ابتدای سال ۱۴۰۲ در جدول ۲ آمده است. این داده‌ها شامل تعداد کل پروازهای شرکت هواپیمایی و تعداد پروازها با تأخیرشان می‌باشد. N نمایانگر تعداد پروازها و X مشخص‌کننده تعداد پروازها با تأخیر است. توزیع X_i به صورت $X_i \sim B(N_i, p_i)$ است که در آن $i = 1, \dots, 20$ مشخص‌کننده تعداد شرکت هواپیمایی است. همانند براون (۲۰۰۸) از تبدیل ثابت‌سازی واریانس برای داده‌ها استفاده شد تا مقادیر Y به صورت $Y_i = \text{Arcsin} \sqrt{\frac{X_i + 0.25}{N_i + 0.5}}$ به دست آیند، که دارای توزیع نرمال به صورت $Y_i \sim N(\theta_i, \frac{1}{4N_i})$ می‌باشد و در آن $\theta_i = \text{Arcsin}(\sqrt{p_i})$. سپس با استفاده از تابع درستمایی تجربی توأم مقید برای تحلیل داده‌های Y_i استفاده شد و برآوردهای $\hat{\theta}_i$ از رابطه $\hat{\theta}_i = \frac{\hat{\lambda}}{\hat{\lambda} + \frac{1}{4N_i}} Y_i + \frac{\frac{1}{4N_i}}{\hat{\lambda} + \frac{1}{4N_i}} \hat{\mu}$ به دست آمدند. نتایج در جدول ۲ گزارش شده و برآوردهای ابرپارامترهای متناظر با $\hat{\mu} = 0.525$ و $\hat{\lambda} = 0.16$ به دست آمدند.

بحث و نتیجه‌گیری

از روش‌های گشتاوری، SURE و روش پیشنهادی REML برای برآورد پارامترهای انقباضی مدل سلسله مراتبی نیم-پارامتری استفاده شد. در یک مطالعه شبیه‌سازی نشان داده شد که روش REML به‌ویژه وقتی تعدادی داده پرت در جمع داده‌ها وجود دارد عملکرد بهتری نسبت به سایر روش‌ها دارد. این مهم به دلیل آن است که برخلاف

جدول ۲. داده‌های واقعی، تبدیل شده و برآورد θ_i

θ_i	Y_i	پرواز خروجی تأخیردار	پرواز خروجی	شرکت هواپیمایی
۰٫۷۱۱۰	۰٫۷۱۲۸	۶۸۲	۱۵۹۵	تابان
۰٫۵۶۴۵	۰٫۵۶۴۶	۱۶۰۶	۵۶۱۰	آنا ایر
۰٫۵۷۰۶	۰٫۵۷۱۳	۳۲۰	۱۰۹۵	چابهار ایر
۰٫۶۲۴۲	۰٫۶۲۵۱	۵۶۹	۱۶۶۲	ساها
۰٫۶۱۶۳	۰٫۶۱۶۸	۹۱۶	۲۷۳۸	ایرتور
۰٫۶۰۳۱	۰٫۶۰۳۳	۱۷۸۶	۵۵۴۹	آسمان
۰٫۶۶۷۰	۰٫۶۶۷۶	۱۳۹۶	۳۶۴۲	کیش ایر
۰٫۵۵۷۶	۰٫۵۵۷۸	۹۸۸	۳۵۲۷	وارش
۰٫۵۷۵۶	۰٫۵۷۵۷	۱۵۱۲	۵۱۰۱	ایران ایر
۰٫۶۴۱۲	۰٫۶۴۱۷	۱۲۷۳	۳۵۵۳	کازرون
۰٫۵۴۰۶	۰٫۵۴۰۷	۷۸۳	۲۹۵۶	کاسپین
۰٫۵۳۸۱	۰٫۵۳۸۱	۷۵۳	۲۸۶۷	قشم ایر
۰٫۴۵۲۰	۰٫۴۵۱۳	۳۲۷	۱۷۲۰	پارس ایر
۰٫۵۰۵۶	۰٫۵۰۵۵	۶۸۵	۲۹۲۲	سیهران
۰٫۵۶۸۷	۰٫۵۶۹۴	۲۹۴	۱۰۱۲	معراج
۰٫۴۳۳۸	۰٫۴۳۳۱	۳۶۸	۲۰۹۰	زاگرس
۰٫۳۷۸۲	۰٫۳۷۶۰	۱۴۱	۱۰۴۷	پویا ایر
۰٫۲۲۶۹	۰٫۲۱۷۱	۲۲	۴۷۹	آساجت
۰٫۵۰۸۹	۰٫۵۰۸۹	۳۶۳۲	۱۵۳۰۳	فلای پرشیا
۰٫۲۵۹۳	۰٫۲۵۸۶	۳۹۳	۶۰۱۲	ماهان

توابع درست‌نمایی مرسوم، در ساختار این تابع، توابع برآوردساز به صورت مجزا نقش آفرینی می‌کنند. بنابراین پیشنهاد می‌شود که از این تابع به عنوان یک روش مفید برای تحلیل مدل‌های سلسله مراتبی ناهمگن با ساختار نیم-پارامتری استفاده شود.

تقدیر و تشکر

نویسندگان مقاله مراتب قدردانی و سپاس خود را از پیشنهادات ارزنده داوران، سردبیر و ویراستار محترم مجله که باعث افزایش سطح کیفی مقاله شده است، اعلام می‌دارند.

مراجع

کرمی‌کبیر، ح. و آرشی، م. (۱۳۹۳)، برآوردگر انقباضی در توزیع نرمال چند متغیره تحت فضای پارامتر محدود، مجله علوم آماری، ۸، ۷۵-۹۲.

Baranchik A. J. (1970), A Family of Minimax Estimators of the Mean of a Multivariate Normal Distribution, *The Annals of Mathematical Statistics*, **41**, 642-645.

- Bayati, M., Ghoreishi, S. K., and Wu, J. (2021), Bayesian Analysis of Restricted Penalized Empirical Likelihood, *Computational Statistics*, **36**, 1321-1339.
- Berger, J. and Strawderman, W. E. (1996), Choice of Hierarchical Priors: Admissibility in Estimation of Normal Means, *The Annals of Statistics*, **24**, 931-951.
- Brown, L. D. (1971), Admissible Estimators, Recurrent Diffusions, and Insoluble Boundary Value Problems, *The Annals of Mathematical Statistics*, **42**, 855-903.
- Brown, L. D. (2008), In-season Prediction of Batting Average: A Field Test of Empirical Bayes and Bayes Methodologies, *The Annals of Applied Statistics*, **2**, 113-152.
- Ghoreishi, S. K. (2017), Bayesian Analysis of Hierarchical Heteroscedastic Linear Models Using Dirichlet-Laplace Priors, *Journal of statistical Theory and Applications*, **16**, 53-64.
- James, W. and Stein, C. M. (1961), Estimation With Quadratic Loss. *Proceedings of the 4th Berkeley Symposium on Probability and Statistics*, **I**, 367-379.
- Stein, C. M. (1962), Confidence Sets for the Mean of a Multivariate Normal Distribution (With Discussion), *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **24**, 265-296.
- Strawderman, W. E. (1971), Proper Bayes Minimax Estimators of the Multivariate Normal Mean, *The Annals of Mathematical Statistics*, **42**, 385-388.
- Xie, X., Kou, S. C. and Brown, L. D. (2012), SURE Estimates for a Heteroscedastic Hierarchical Model, *Journal of the American Statistical Association*, **107**, 1465-1479.
- Xie, X., Kou, S. C. and Brown, L. D. (2016), Optimal Shrinkage Estimation of Mean Parameters in Family of Distributions With Quadratic Variance, *The Annals of Statistics*, **44**, 564.