



## Double Penalized Mixed Effects Quantile Regression Modeling Using the Maximum Likelihood Approach

Jafari Maryaki, F. , Golalizadeh, M. 

Department of Statistics, Tarbiat Modares University, Tehran , Iran.

**Corresponding author:** M. Golalizadeh, [golalizadeh@modares.ac.ir](mailto:golalizadeh@modares.ac.ir)

**Received:** 31/10/2022 **Revised:** 16/11/2023 **Accepted and Published Online:** 18/11/2023.

### Introduction

Mixed effects is one of the powerful statistical approaches used to model the relationship between the response variable and some predictors in analyzing data with a hierarchical structure. Estimating parameters in these models is often done via either the least squares error or maximum likelihood approaches. The estimated parameters obtained through either of these approaches could be more efficient if the error distributions are non-normal. In such cases, using the mixed effects quantile regression is preferred alternatively. Moreover, when the number of variables studied increases, the penalized mixed effects quantile regression is one of the best methods to gain prediction accuracy and the model's interpretability. We propose a new tool to model the hierarchical structural data by introducing the penalty, a function of random effects and fixed effects. This new idea can simultaneously estimate the parameters and predict random effects. We can then make further statistical inferences on the parameters using the likelihood function built based on modeling such as a new penalized mixed effects quantile regression model.

### Material and Methods

In this paper, under the assumption of an asymmetric Laplace distribution for random effects, we proposed a double penalized model in which both the random and fixed effects are independently penalized. Then, invoking a standard algorithm to estimate the parameters using the likelihood approach constructs our subsequent step in modeling the data. Since the likelihood function had no tractable form, we considered the EM algorithm to approx-

imate the integrals. To do so, we also extended the previously proposed algorithm to derive the maximum likelihood estimate of the parameters.

### Results and Discussion

The performance of the new method proposed in this paper is evaluated in the simulation studies, and a discussion of the results is presented along with a comparison with some competing models. Our simulation experiments show that in addition to the robustness of the presented method, it performs better compared to the types of penalty in the model based on different penalty parameters. Applying the proposed model in real data analysis also showed that the new idea works well compared with some standard tools. The exciting feature of the double-penalized model is on confining the computation procedure to estimate only one tune parameter because the other one is optimized, using the samples and updated components while implementing the algorithm. Our results showed that considering a lower value for the tuning parameter of penalizing the fixed effects is preferred. We cannot recommend choosing a specific distribution for the random effects while using our proposed model. However, we can assure that our model does the same task as other competing models do in this regard.

### Conclusion

The model presented in this paper is a double penalized model, which is a function of random effects and model parameters. In this paper, we achieved the optimal shrinkage for the effects using the simulated and real data. Similar to some standard methods, the degree of shrinkage for the random effects in our model depends on the tuning parameter selection. We aim to study the efficiency of our proposed model while one prefers to follow a Bayesian approach. Also, correctly choosing the penalty parameters from a theoretical viewpoint is another proposal for future research.

**Keywords:** Mixed Effects Quantile Regression, Laplace Distribution, Penalty Function, Shrinkage Approach, High Dimensional.

**Mathematics Subject Classification (2010):** 62G08, 62J05.



©The Author(s). The Publisher is Iranian Statistical Society.

This is an open access article distributed under the terms and conditions of [\(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/)



مجله علوم آماری، پاییز و زمستان ۱۴۰۲

جلد ۱۷، شماره ۲، ص ۳۸۹ -- ۴۰۵

DOI: 10.52547/jss.17.2.09

مقاله پژوهشی

## مدل‌بندی رگرسیونی چندکی آمیخته تاوانیده دوگانه از طریق رویکرد درست‌نمایی

فروزان جعفری مریکی، موسی گلعلی‌زاده

گروه آمار، دانشگاه تربیت مدرس

نویسنده مسئول: موسی گلعلی‌زاده، [golalizadeh@modares.ac.ir](mailto:golalizadeh@modares.ac.ir)

تاریخ دریافت: ۱۴۰۱/۸/۹ تاریخ بازنگری: ۱۴۰۲/۸/۲۵ تاریخ پذیرش و انتشار: ۱۴۰۲/۸/۲۷

**چکیده:** مدل‌های آمیخته از جمله ابزارهای قوی آماری است که برای مدل‌بندی ارتباط بین متغیر پاسخ و متغیرهای تبیینی در تحلیل داده‌هایی با ساختار سلسله‌مراتبی به‌کار می‌رود. برآوردگرهای به‌دست آمده در این مدل‌ها با استفاده از هر یک از روش‌های کمترین توان دوم خطاها و ماکسیمم درست‌نمایی، زمانی که توزیع خطاها غیر نرمال باشد، از کارایی لازم برخوردار نیستند. در این‌گونه مواقع می‌توان از مدل رگرسیون چندکی آمیخته به‌عنوان جایگزین استفاده کرد. زمانی که تعداد متغیرهای مورد بررسی در این نوع مدل‌بندی افزایش می‌یابد، رگرسیون چندکی آمیخته تاوانیده یکی از بهترین روش‌ها برای افزایش دقت پیشگویی و تفسیرپذیری مدل است. در این مقاله با در نظر گرفتن توزیع لاپلاس نامتقارن برای اثرهای تصادفی، یک مدل تاوانیده دوگانه به‌عنوان تابعی همزمان از اثرهای تصادفی و پارامترهای مدل پیشنهاد می‌شود. سپس، عملکرد روش پیشنهادی با استفاده از مطالعه شبیه‌سازی آماری مورد ارزیابی قرار گرفته و بحث راجع به نتایج حاصل به همراه مقایسه با برخی مدل‌های رقیب ارائه می‌شود. به‌علاوه، کاربردی از آن در تحلیل یک مثال واقعی نمایش داده خواهد شد.

واژه‌های کلیدی: رگرسیون چندکی آمیخته، توزیع لاپلاس، تابع تاوان، رویکرد انقباضی، بُعد بالا.

کد موضوع‌بندی ریاضی (۲۰۱۰): 62G08 ، 62J05.



©نویسندگان). ناشر انجمن آمار ایران است.  
این مقاله با دسترسی آزاد تحت شرایط و ضوابط (CC BY-NC 4.0) توزیع شده است.

## ۱ مقدمه

مدل‌های رگرسیونی میانگین بیشترین کاربرد را برای ایجاد ارتباط بین میانگین متغیر پاسخ و هر تابعی از متغیرهای تبیینی دارند. اما در مواردی محقق علاقه‌مند است تحلیل داده‌ها را براساس مدل‌بندی معیارهای دیگری از توزیع متغیر پاسخ غیر از میانگین انجام دهد. به عنوان یک فعالیت علمی در راستای چنین نیازی، رگرسیون چندکی که حالت توسعه یافته‌ی از رگرسیون میانگین است، توسط **کوئنکر و باست** (۱۹۷۸) معرفی شد. در چند دهه اخیر، القای نقش اثرهای تصادفی در مدل‌های رگرسیون چندکی برای تحلیل داده‌های طولی مورد توجه آماردانان و محققین فعال در حوزه کاربرد آمار قرار گرفته است. به عنوان شاهد مثالی از این موضوع می‌توان به فعالیت **گراسی و بوتای** (۲۰۰۷) اشاره کرد. لازم به توضیح است که ماهیت داده‌های طولی به‌گونه‌ای است که بین مشاهدات هر آزمودنی در زمان‌های متفاوت یک تغییرپذیری درون آزمودنی رخ می‌دهد که اگر در مدل لحاظ نشود منجر برآوردهای اریب برای پارامترهای مدل می‌شود (**دیگل و همکاران**، ۲۰۰۲).

وقتی تعداد متغیرهای مورد بررسی در مدل افزایش می‌یابد، روش‌های انتخاب متغیر، با تاوانیده کردن مدل مورد نظر، باعث افزایش دقت پیشگویی و تفسیرپذیری مدل می‌شوند. در این راستا روش‌های انقباضی متفاوتی مانند لاسو (**تیشیرانی**، ۱۹۹۶)، لاسوی سازوار (**زو**، ۲۰۰۶)، انحراف مطلق برش صاف<sup>۱</sup> *SCAD* (**فن و لی**، ۲۰۰۱) و گاروت نامنفی<sup>۲</sup> (**بریمن**، ۱۹۹۵) گسترش یافت. لازم به اشاره است که با استفاده از این روش‌ها برآورد و انتخاب متغیر به‌طور همزمان انجام می‌شود. محققان زیادی از رویکردهای انقباضی برای تحلیل داده‌های طولی در رگرسیون چندکی استفاده کردند. منابع علمی نشان می‌دهند که **کوئنکر** (۲۰۰۴) برای اولین بار مدل رگرسیونی چندکی  $L_1$  تاوانیده را برای تحلیل داده‌های طولی پیشنهاد داد به‌طوری‌که برای برآورد چندک‌های توزیع با اثرهای ویژه هر آزمودنی از روش‌های تاوانیده کردن اثرهای تصادفی در تابع زیان چندکی استفاده می‌شود. سپس، **گراسی و بوتای** (۲۰۰۷) به‌جای ارزیابی تابع زیان تاوانیده از رویکرد تابع درستنمایی مبتنی بر تابع چگالی لاپلاس نامتقارن بهره بردند. آنها، با درنظر گرفتن توزیع نرمال برای اثرهای تصادفی، رویکردی معادل با روش ارائه‌شده در **کوئنکر** (۲۰۰۴) معرفی کرده و نشان دادند روش آنها در مقایسه با روش اضافه کردن تابع تاوان به تابع زیان بهتر عمل می‌کند. جالب است که این رویکرد به‌صورت خودکار سطح بهینه‌ای از پارامتر تاوان را به‌دست می‌آورد. **باندل و همکاران** (۲۰۱۰) روشی برای انتخاب همزمان اثرهای تصادفی و ثابت با استفاده از رویکرد لاسوی سازوار برای مدل‌های رگرسیونی مبتنی بر میانگین پیشنهاد دادند. **لی و همکاران** (۲۰۲۰) با تلفیق روش ارائه‌شده در **باندل و همکاران** (۲۰۱۰) و **کوئنکر** (۲۰۰۴) الگوریتمی برای به دست آوردن همزمان برآوردهای اثرهای ثابت و تصادفی ارائه دادند. در واقع، آنها، بر اساس تاوان لاسو، در مدل رگرسیونی چندکی آمیخته تاوان را هم برای اثرهای تصادفی و هم اثرهای ثابت پیشنهاد کردند. الگوریتم ارائه‌شده توسط آنها رویکردی غیر درستنمایی و بر اساس الگوریتم‌های معرفی شده در رگرسیون چندکی آمیخته تاوانیده لاسو است. به نظر می‌رسد با استفاده از رویکرد درستنمایی، در مدل رگرسیونی

<sup>1</sup>Smoothly Clipped Absolute Deviation

<sup>2</sup>Non-negative garotte

چندکی آمیخته تاوانیده که تاوان اعمال شده همزمان بر روی اثرهای تصادفی و اثرهای ثابت اعمال می‌شود رویکرد بهتری برای استنباط راجع به پارامترها فراهم شود. با چنین نگاهی، آقامحمدی و محمدی (۱۳۹۴) با کمک رویکرد بیزی و فراخوانی تاوان لاسو و لاسوی تطبیق پذیر به تحلیل داده‌های طولی توسط مدل رگرسیون چندکی پرداختند.

در مقاله حاضر با پیروی از رویکرد انتخاب متغیر معرفی شده توسط تیشیرانی (۱۹۹۶)، با تلفیق مدل ارائه شده در کوئنکر (۲۰۰۴) و رویکرد درستنمایی بحث شده در گراسی و بوتای (۲۰۰۷)، رویکردی جدید در انتخاب متغیر در مدل رگرسیون چندکی آمیخته پیشنهاد می‌شود که می‌توان به‌طور همزمان پارامترها را برآورد و اثرهای تصادفی را پیشگویی کرد. این رویکرد طوری است که به‌طور همزمان علاوه بر انقباض پارامترها، به‌طور خودکار با استفاده از داده‌ها و تابع درستنمایی آنها به درجه‌ای از انقباض اثرهای تصادفی دست می‌یابد. برای بررسی عملکرد این مدل تلفیقی، آزمایش شبیه‌سازی برای مقایسه مدل‌های متنوع در حالت‌های متفاوت انجام و نتایج مورد بحث قرار می‌گیرد. تحلیل یک مجموعه داده واقعی برای نمایش کاربست بهتر مدل پیشنهادی نیز ارائه می‌شود.

به منظور ارائه مطالب مرتبط با موضوع این مقاله، در بخش ۲ به معرفی مدل رگرسیون چندکی آمیخته پرداخته می‌شود. در بخش ۳، پس از بررسی رگرسیون چندکی آمیخته تاوانیده، مدل پیشنهادی این مقاله نیز معرفی می‌شود. نتایج حاصل از انجام شبیه‌سازی و تحلیل نتایج مرتبط با بررسی یک مثال واقعی پایان بخش مقاله حاضر است.

## ۲ خلاصه‌ای از مدل رگرسیونی چندکی آمیخته

مدل رگرسیونی چندکی آمیخته معمولاً به‌صورت

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta}_\tau + u_i + \varepsilon_{ij\tau}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, N$$

نوشته می‌شود به‌طوری‌که  $\mathbf{x}_{ij}^T$  بردار  $p$  بعدی از سطرهای ماتریس طرح معلوم  $\mathbf{X}_i$  و  $y_{ij}$ ،  $j$ امین مشاهده متغیر تصادفی پیوسته در آزمودنی  $i$ ام،  $\boldsymbol{\beta}_\tau$  برداری  $p$  بعدی از پارامترها و  $\varepsilon_{ij\tau}$  خطای تصادفی مدل است که مستقل از هم و به‌علاوه مستقل از اثرهای تصادفی  $u_i$  در نظر گرفته می‌شوند. تابع چندکی شرطی خطی آمیخته به‌صورت

$$G_{y_{ij}|u_i}(\tau|\mathbf{x}_{ij}, u_i) = \mathbf{x}_{ij}^T \boldsymbol{\beta}_\tau + u_i, \quad j = 1, \dots, n_i, \quad i = 1, \dots, N \quad (1)$$

تعریف می‌شود، که در آن  $G_{y_{ij}|u_i}(\tau|\mathbf{x}_{ij}, u_i) = F_{y_{ij}|u_i}^{-1}(\cdot)$  معکوس تابع توزیع متغیر پاسخ به شرط اثرهای تصادفی  $u_i$  است و  $0 < \tau < 1$ . می‌توان انتظار داشت که مقدار برآورد به‌دست آمده از معادله (۱) به  $\tau$  امین چندک توزیع  $y_{ij}$  وابسته است. اما، در ادامه برای سادگی مباحث، اندیس  $\tau$  حذف می‌شود. گراسی و بوتای (۲۰۰۷)

پیشنهاد دادند که اگر اثرهای تصادفی معلوم باشند، برآورد  $\beta \in \mathbb{R}^p$  کمیتی است که عبارت

$$\tilde{V}_i(\beta) = \sum_{j \in \{j: \tilde{y}_{ij} \geq \mathbf{x}_{ij}^T \beta\}}^{n_i} \tau |\tilde{y}_{ij} - \mathbf{x}_{ij}^T \beta| + \sum_{j \in \{j: \tilde{y}_{ij} < \mathbf{x}_{ij}^T \beta\}}^{n_i} (1 - \tau) |\tilde{y}_{ij} - \mathbf{x}_{ij}^T \beta| \quad (2)$$

را مینیمم کند، که در آن  $\tilde{y}_{ij} = y_{ij} - u_i$ . آن‌ها نشان دادند که با در نظر گرفتن توزیع لاپلاس نامتقارن برای متغیر پاسخ، کمینه‌کننده رابطه (۲) از طریق رویکرد تابع درست‌نمایی به دست می‌آید. لازم به اشاره است که توزیع  $y_{ij}$  به شرط معلوم بودن  $u_i$  ها، برای هر  $j = 1, \dots, n_i$ ,  $i = 1, \dots, N$  به طور مستقل توزیع لاپلاس نامتقارن است اگر

$$f(y_{ij} | \beta, u_i, \sigma) = \frac{\tau(1-\tau)}{\sigma} \exp\left\{-\rho_\tau\left(\frac{y_{ij} - \mu_{ij}}{\sigma}\right)\right\},$$

یا به صورت نمادین،  $y_{ij} | u_i \sim ALD(\mu_{ij}, \phi, \tau)$ ، که در آن  $\mu_{ij} = \mathbf{x}_{ij}^T \beta + u_i$  به عنوان پارامتر مکانی، معرف پیشگوی خطی  $\tau$  امین چندک بوده،  $\sigma \in (0, \infty)$  پارامتر مقیاس و  $\tau \in (0, 1)$  پارامتر چولگی،  $\rho_\tau(\nu) = \nu(\tau - I(\nu \leq 0))$  تابع چک (به دلیل نمایش هندسی آن) و  $I(\cdot)$  تابع نشانگر است. اگر بردار مشاهدات هر آزمودنی به صورت  $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})$  نوشته شود، آنگاه به شرط معلوم بودن اثرهای تصادفی داریم  $i = 1, \dots, N$  به صورت  $f(\mathbf{y}_i | u_i, \beta, \sigma) = \prod_{j=1}^{n_i} f(y_{ij} | u_i, \beta, \sigma)$  به دلیل نامعلوم بودن  $u_i$  ها تابع چگالی توام  $(\mathbf{y}_i, u_i)$  برای آن  $f(\mathbf{y}_i, u_i | \eta) = f(\mathbf{y}_i | u_i, \beta, \sigma) f(u_i | \varphi)$  نوشته می‌شود، که در آن  $\eta = (\beta, \sigma, \varphi)$  پارامترهای مدل بوده و  $f(u_i | \varphi)$  تابع چگالی اثرهای تصادفی است. تابع چگالی توام  $(\mathbf{y}, \mathbf{u})$  برای  $N$  آزمودنی به صورت  $f(\mathbf{y}, \mathbf{u} | \eta) = \prod_{i=1}^N f(\mathbf{y}_i | u_i, \beta, \sigma) f(u_i | \varphi)$  است که در آن  $\mathbf{y} = (y_1, \dots, y_N)$  و  $\mathbf{u} = (u_1, \dots, u_N)$  در نتیجه تابع چگالی حاشیه‌ای  $\mathbf{y}$  با انتگرال‌گیری نسبت به اثرهای تصادفی به صورت  $f(\mathbf{y} | \eta) = \int_{\mathbb{R}^N} f(\mathbf{y}, \mathbf{u} | \eta) d\mathbf{u}$  بدست می‌آید، که منجر به تابع درست‌نمایی به صورت  $L(\eta; \mathbf{y}) = \sum_{i=1}^N L(\eta; \mathbf{y}_i)$  می‌شود. این تابع مبنای مناسبی برای انجام استنباط آماری راجع به پارامترهای مدل است که در بخش بعد به آن پرداخته می‌شود.

### ۳ برآورد درست‌نمایی مدل رگرسیونی چندکی آمیخته تاوانیده

رویکرد رگرسیون چندکی تاوانیده توسط کوئنکر (۲۰۰۴) با اضافه کردن تابع تاوان لاسو به اثرهای تصادفی مدل به صورت

$$\min_{\mathbf{u}, \beta} \sum_{i=1}^N \sum_{j=1}^{n_i} \omega_{ij} \rho_\tau(y_{ij} - \mathbf{x}_{ij}^T \beta - u_i) + \lambda \sum_{i=1}^N |u_i|$$

معرفی شد، که در آن  $w_{ij}$  وزنی است که برای کم کردن اثر داده‌های دورافتاده و هم‌چنین چندک  $\tau$  ام در پیشگویی  $u_i$  استفاده می‌شود تا برآورد حاصل به استواری<sup>۱</sup> برسد و پارامتر  $\lambda$  به‌عنوان پارامتر تاوان در نظر گرفته می‌شود. همانطور که ملاحظه می‌شود در این رویکرد تاوان به‌کار برده شده تنها تابعی از اثرهای تصادفی است و انقباضی روی پارامتر اثرهای ثابت، یعنی  $\beta$  صورت نمی‌گیرد. به علاوه، می‌توان ملاحظه کرد که رویکرد پیشنهادی کوئنکر (۲۰۰۴) توانایی مدل‌بندی داده‌های بُعد بالا را ندارد چرا که مدل ایشان کوششی در انتخاب متغیر ندارد. بنابراین، پیشنهاد مقاله حاضر آن است که مدنظر قرار دادن عبارت

$$\min_{u, \beta} \sum_{i=1}^N \sum_{j=1}^{n_i} \omega_{ij} \rho_{\tau}(y_{ij} - \mathbf{x}_{ij}^T \beta - u_i) + \lambda_1 \sum_{i=1}^N |u_i| + \lambda_2 \sum_{l=1}^p |\beta_l| \quad (۳)$$

که در آن  $\lambda_i \geq 0$ ، به ازای  $i = 1, 2$  پارامترهای تاوان هستند، رویکرد مناسبی برای مدل‌بندی داده‌های بُعد بالا از طریق مدل رگرسیون چندکی آمیخته است. با قبول چنین پیشنهادی، هدف اولیه ما این است که استنباط در مورد پارامترهای مدل را براساس تابع درستنمایی مرتبط با این رویکرد انجام دهیم. می‌توان انتظار داشت که تابع درستنمایی حاصل شامل انتگرال یا انتگرال‌هایی است که شکل صریحی برای حل آن‌ها وجود ندارد. نکته امیدوارکننده این است که، معمولاً برای تقریب اینگونه انتگرال‌ها (ها) روش الگوریتم مونت کارلوی  $EM$  (بوث و هابرت، ۱۹۹۹) مورد استفاده قرار می‌گیرد. به علاوه، توزیع اثرهای تصادفی نقش مهم و اساسی را در برآورد درستنمایی پارامترهای مدل بازی می‌کند. بنابراین در ادامه به بررسی نقش توزیع‌های متفاوت برای اثرهای تصادفی پرداخته می‌شود.

اگر اثرهای تصادفی دارای توزیع لاپلاس نامتقارن با پارامتر مکان صفر، پارامتر مقیاس  $\phi$  و پارامتر چولگی  $\sigma/5$  نمایش داده شده با نماد اختصاری  $ALD(\sigma, \phi, \sigma/5)$  به ازای  $u_i \sim$   $i = 1, \dots, N$  باشد. آنگاه، تابع چگالی توام  $(y_i, u_i)$  برای  $i$  امین آزمودنی به‌صورت

$$\begin{aligned} f(y_i, u_i | \eta) &= f(u_i | \varphi) \prod_{j=1}^{n_i} f(y_{ij} | u_i, \beta, \sigma) \\ &= \left\{ \frac{\tau(1-\tau)}{\sigma} \right\}^{n_i} \frac{1}{4\varphi} \exp \left\{ - \sum_{j=1}^{n_i} \left\{ -\rho_{\tau} \left( \frac{y_{ij} - \mu_{ij}}{\sigma} \right) \right\} - \frac{|u_i|}{2\varphi} \right\} \end{aligned} \quad (۴)$$

به دست می‌آید. با توجه به این که توزیع شرطی  $f(y_i, u_i | \eta) \propto f(u_i | y_i, \eta)$  تابعی لگ-مقعر (گیلکز، ۱۹۹۵) است، می‌توان از طریق روش نمونه‌گیری گیبز یا الگوریتم نمونه‌گیری ردی تطبیقی برای تولید نمونه‌های تصادفی به‌عنوان تحقق‌هایی از اثرهای تصادفی استفاده کرد. لذا، اگر فرض شود نمونه  $m_i$  تایی  $v_i = (v_{i1}, \dots, v_{im_i})$  رخ داده‌هایی از اثرهای تصادفی تولید شده از تابع چگالی شرطی اثرهای تصادفی به شرط متغیر پاسخ باشند برآورد

<sup>1</sup>Robust

ماکسیم درستنمایی  $\phi$  به صورت

$$\hat{\phi} = \frac{1}{N} \sum_{i=1}^N \frac{1}{m_i} \sum_{k=1}^{m_i} \rho_{\sigma, \lambda} (v_{ik}) \quad (5)$$

به دست می‌آید. توجه شود که تابع چگالی توام (۴) را می‌توان به صورت

$$f(\mathbf{y}_i, u_i | \eta) = \frac{\tau^{n_i} (1 - \tau)^{n_i}}{(\lambda \varphi^2)^{n_i} \lambda_1^{n_i}} \exp \left\{ -\frac{1}{\sigma} \left[ \sum_{j=1}^{n_i} c_{ij} |y_{ij} - \mu_{ij}| + \lambda_1 |u_i| \right] \right\} \quad (6)$$

بازنویسی کرد، که در آن  $\lambda_1 = \frac{\sigma}{\tau \varphi}$  و  $c_{ij} = [(1 - \tau) I(y_{ij} \leq \mu_{ij}) + \tau I(y_{ij} > \mu_{ij})]$  توجه شود که در این حالت همین مقدار  $\lambda_1$  در رابطه (۳) هم استفاده می‌شود. پس از آن، برای به دست آوردن برآورد ماکسیم درستنمایی پارامتر  $\eta$  برای  $\tau$  امین چندک مدل، از الگوریتم معرفی شده در (گراسی و بوتای، ۲۰۰۷) که به صورت زیر توسعه داده شد، استفاده می‌شود.

**الگوریتم ۱.** الگوریتم برآورد ماکسیم درستنمایی پارامتر  $\eta$  برای  $\tau$  امین چندک:

گام ۱- ابتدا قرار دهید  $t = 0$  و برای پارامترهای  $(\beta^{(t)}, \sigma^{(t)}, \varphi^{(t)}) = \eta^{(t)}$  مقادیر اولیه‌ای در نظر بگیرید و آن‌ها را در تابع چگالی  $f(u_i | \mathbf{y}_i, \eta^{(t)})$  جای‌گذاری کنید.

گام ۲- با استفاده از روش‌های نمونه‌گیر به‌طور مثال نمونه‌گیر گیبز (گیلکز و همکاران، ۱۹۹۵) از تابع چگالی،  $f(u_i | \mathbf{y}_i, \eta^{(t)})$  برای هر آزمودنی  $i = 1, \dots, N$  به‌طور مستقل نمونه‌هایی تولید کنید و آن‌ها را در بردار  $\mathbf{v}_i^{(t)} = (v_{i1}^{(t)}, \dots, v_{im_i}^{(t)})$  قرار دهید.

گام ۳- عبارت  $\sum_{l=1}^p |\beta_l| + \lambda_2 \sum_{i=1}^N \frac{1}{m_i} \sum_{k=1}^{m_i} \sum_{j=1}^{n_i} \rho_{\tau} (\tilde{y}_{ikj}^{(t)} - \mathbf{x}_{ij}^T \beta)$  را نسبت به  $\beta$  مینیمم کنید، که در

آن  $\tilde{y}_{ikj}^{(t)} = y_{ij} - v_{ik}^{(t)}$  و  $\sum_{l=1}^p |\beta_l|$  تاوان لاسو در تابع هدف است. توجه شود که بعد از انجام این عمل  $\beta^{(t+1)}$ ها به دست می‌آیند. در نتیجه می‌توان بلافاصله واریانس خطاها را از طریق رابطه

$$\sigma^{(t+1)} = \frac{1}{\sum_{i=1}^N n_i} \sum_{i=1}^N \frac{1}{m_i} \sum_{k=1}^{m_i} \sum_{j=1}^{n_i} \rho_{\tau} (\tilde{y}_{ikj}^{(t)} - \mathbf{x}_{ij}^T \beta^{(t+1)})$$

محاسبه کرد. در نهایت قرار دهید  $(\mathbf{v}_1^{(t)}, \dots, \mathbf{v}_N^{(t)}) = \varphi^{(t+1)}$ ، که در آن  $\hat{\phi}$  برآورد ماکسیم درستنمایی  $\phi$  به شرط  $\mathbf{y}$  است که از طریق عبارت (۵) قابل محاسبه است. شایان ذکر است که، در این مرحله،  $\lambda_1$  در عبارت (۳) به صورت  $\lambda_1^{(t+1)} = \frac{\sigma^{(t+1)}}{\tau \varphi^{(t+1)}}$  قابل محاسبه است.

گام ۴- حال قرار دهید  $t = t + 1$  و گام ۱ تا ۳ را تکرار کنید.



گام ۵- الگوریتم را تا جایی تکرار کنید که برآوردهای  $\eta^{(t)}$  به ازای دنباله متناهی از  $t$  به همگرایی قابل قبولی برسد.

در الگوریتم ذکر شده، در هر تکرار  $\eta^{(t+1)}$  برآورد ماکسیمم درستمایی  $\eta$  به شرط  $\eta^{(t)}$  است. مرحله نمونه‌گیری در اجرای این الگوریتم مستلزم محاسبه همگرایی زنجیره مارکوف است (کولز و کارلین، ۱۹۹۶) و بنابراین برآورد به‌دست آمده با استفاده از این روش به تعداد نمونه‌ای که از تابع چگالی تولید می‌شود و تعداد دورریز بستگی دارد. ذکر این نکته ضروری است که فرآیند اجرای الگوریتم طوری است که شخص به‌طور خودکار به درجه بهینه‌ای از انقباض اثرهای تصادفی دست می‌یابد. بنابراین، برای اجرای این الگوریتم تنها لازم است پارامتر  $\lambda_2$  از پیش تعیین شود. این بدان معنی است که مقدار بهینه پارامتر  $\lambda_1$  هنگام اجرای الگوریتم خود به خود به دست می‌آید.

## ۴ آزمایش شبیه‌سازی و تحلیل نتایج آن

در این بخش قصد داریم کارکرد مدل پیشنهادی را از طریق آزمایش‌های شبیه‌سازی بررسی کنیم. واضح است که می‌توان رویکرد شبیه‌سازی پیشنهادی را برای چندک‌های متفاوت توزیع متغیر پاسخ به‌کار گرفت. اما، به دلیل محدودیت فضای مقاله، خود را محدود به چندک خاص از توزیع متغیر پاسخ خواهیم کرد. شایان ذکر است که در تمام موارد شبیه‌سازی، ستون‌های ماتریس طرح مرکزی شده‌اند و در نتیجه مدل نهایی شامل عرض از مبدا نیست. همچنین توجه شود که برای تولید نمونه تصادفی با استفاده از روش نمونه‌گیر گیبز در این بخش و بخش مثال کاربردی از کتابخانه *armspp* در نرم افزار *R* و تابع *arms()* استفاده شده است. برای ارزیابی مدل پیشنهاد شده در بخش قبل، آزمایش شبیه‌سازی با تعداد آزمودنی  $N = 10$  انجام می‌شود که هر یک دارای  $n = 23$  مشاهده هستند. همچنین فرض می‌کنیم  $\tau = 0/25$ . برای تولید داده از مدل

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + u_i + \varepsilon_{ij} \quad i = 1, \dots, 10, \quad j = 1, \dots, 23 \quad (7)$$

به قسمی که  $\mathbf{x}_{ij}^T = (x_{ij1}, \dots, x_{ij8})$  و  $\boldsymbol{\beta} = (3, 1.5, 0, 0, 2, 0, 0, 0)^T$  و  $x_{ij}$ ها از توزیع نرمال ۸ متغیره استاندارد تولید شده‌اند طوری که همبستگی بین هر  $x_{ijk}$  و  $x_{ijl}$  به صورت

$$\text{Corr}(x_{ij\ell}, x_{ijk}) = (0.5)^{|\ell-k|}, \quad 1 \leq \ell, k \leq 8$$

تعریف شده است. تعداد تکرار در شبیه‌سازی برای هر یک از حالت‌های مورد مطالعه را برابر  $10^6$  در نظر گرفتیم. در هر تکرار تعداد نمونه تولیدشده در نمونه گیر گیبز برابر  $m_i = 10000$  به ازای  $i = 1, \dots, 10$  اختیار شد و سپس تعداد گام ۳ و تعداد دورریز در هر تکرار برابر با ۲۵۰۰ مدنظر گرفته شده است. لازم به اشاره است که نحوه انتخاب توزیع آماری برای  $\varepsilon_{ij}$ ها و  $u_i$ ها در ادامه بحث تشریح می‌شوند. برای مقایسه آماری مناسب مدل‌های رقیب از نرم

$L_1$  و  $L_2$  استفاده شده است که به ترتیب با پانویس ۱ و ۲ برای ارزیابی بردار شیب به صورت

$$\|\hat{\beta} - \beta\|_1 = \sqrt{\sum_{\ell=1}^p |\hat{\beta}_\ell - \beta_\ell|}, \quad \|\hat{\beta} - \beta\|_2 = \sqrt{\sum_{\ell=1}^p (\hat{\beta}_\ell - \beta_\ell)^2}$$

تعریف می‌شوند، که در آن  $\hat{\beta}_1$  برآورد چندکی ضریب رگرسیونی و  $\beta_1$  مقدار واقعی آن پارامتر متناظر است. ذکر این نکته ضروری است که با توجه به روش ارائه شده در این مقاله، پارامتر تاوان برای اثرهای تصادفی ( $\lambda_1$ ) در الگوریتم معرفی شده به صورت بهینه انتخاب می‌شود. به عبارتی دیگر، این پارامتر تاوان به‌طور خودکار از طریق نمونه‌های تولید شده و تابع درست‌نمایی متناظر با آن نمونه‌ها به دست می‌آید و نیازی به اعلام مقدار آن قبل از شروع شبیه‌سازی نیست. لذا، با توجه به روش ارائه شده در این مقاله کفایت برای تولید نمونه‌ها پارامتر تاوان برای ضرایب رگرسیونی ( $\lambda_2$ ) از پیش تعیین شود. انتظار می‌رود پس از مدل‌بندی داده‌ها با مدل‌های مورد مطالعه و با استفاده از معیارهای موجود مقدار بهینه آن توسط روش بهینه‌سازی مورد محاسبه قرار گیرد. برای تولید نمونه‌های شبیه‌سازی، پنج مقدار متفاوت  $\lambda_2 = (0.5, 1, 5)$  برای پارامتر تاوان  $\lambda_2$  در نظر گرفته شده است. تجربه محاسباتی ما نشان داد که برای مقادیر بزرگتر، به طور مثال  $(0.5, 1, 5) = \lambda_2$  تمام ضرایب معادله (۷) برابر صفر برآورد می‌شود که مطلوب ما در انجام شبیه‌سازی نیست. در نهایت، برای تکمیل شبیه‌سازی، توزیع اثرهای تصادفی به صورت  $ALD(\phi, \sigma, \delta)$  اختیار شده است. نتایج حاصل از آزمایش شبیه‌سازی در این حالت در جدول ۱ آمده است که در ادامه به تشریح نتایج پرداخته خواهد شد. طبق جدول ۱ شبیه‌سازی برای سه توزیع تصادفی برای خطای مدل به صورت تاوان دوگانه رگرسیون چندکی آمیخته<sup>۱</sup> ( $DLQR$ )، رگرسیون ضرایب رگرسیونی تاوانیده<sup>۲</sup> ( $PFQR$ )، رگرسیون اثرهای تصادفی تاوانیده<sup>۳</sup> ( $PRQR$ ) و رگرسیونی چندکی آمیخته معمولی<sup>۴</sup> ( $QR$ ) انجام شده است.

با توجه به نتایج به دست آمده ملاحظه می‌شود که انتخاب توزیع خطاها در عملکرد نهایی هر یک از روش‌های به کار برده شده تاثیرگذار است. برای پیگیری سراسر نتایج، ابتدا به ازای هر روش مورد استفاده، مدل با  $\lambda_2$  های متفاوت با هم و سپس مدل‌ها بر اساس توزیع‌های متفاوت در نظر گرفته شده با هم مقایسه می‌شوند. می‌توان ملاحظه کرد، زمانی که  $\lambda_2$  برابر با  $0.2$  و  $0.4$  باشد، وقتی که توزیع خطا نرمال و  $t(4)$  اختیار شوند مدل پیشنهادی عملکرد مناسبی نداشته است. البته باید اشاره شود که اگر چه در تمامی تکرارهای مرتبط با شبیه‌سازی برآورد پارامترهایی که واقعاً صفر بودند دقیقاً صفر شد اما میزان  $L_1$  و  $L_2$  به دست آمده در این دو حالت تفاوت چشم‌گیری نسبت به حالت  $\lambda_2$  های دیگر داشتند. اتفاق چنین امری به این دلیل است که برآوردهای غیر صفر پارامترها تفاوت معنی‌داری با مقدار واقعی پارامترهای متناظر غیر صفر در شروع تولید نمونه‌های شبیه‌سازی داشتند. به عبارتی دیگر، مدل مورد اشاره توانسته است پارامترهایی که دقیقاً صفر بوده‌اند را به صورت دقیق صفر برآورد کند. اما به همان اندازه

<sup>1</sup>Double Lasso Penalized Quantile Regression

<sup>2</sup>Penalized Fixed Effect Quantile Regression

<sup>3</sup>Penalized Random Effect Quantile Regression

<sup>4</sup>General Linear Quantile Regression

برآوردهای پارامترهای غیرصفر را بسیار منقبض کرد، و می‌توان گفت مدل به حالت کم برازش رسید. باید اشاره شود که وقتی برای مولفه خطا توزیع  $ALD(0, 1, 0.25)$  اختیار شد، تفاوت مشاهده شده خیلی محسوس نبود. طبق نتایج به دست آمده به نظر می‌رسد که برای  $\lambda_2$  های متفاوت، نتایج مربوط به  $\lambda_2 = 0.5$  برای توزیع نرمال در مدل  $DLQR$  عملکرد قابل قبولی داشته است. حال آن‌که، در صورت اختیار  $t(4)$  برای توزیع خطا، اگر شخص خود را محدود به  $\lambda_2 = 0.1$  کند نتایج حاصل نسبت به  $\lambda_2$  های دیگر رضایت بخش‌تر است. نباید از این نکته غافل شد که هنگام برازش مدل  $PFQR$ ، انتخاب متغیر بدون اعمال تاوان روی اثرهای تصادفی صورت می‌گیرد. این در حالی است که مدل  $PRQR$  تاوان را فقط روی اثرهای تصادفی اعمال می‌کند و مدل  $QR$  چیزی جز برازش یک مدل رگرسیون چندکی آمیخته عادی نیست. اگر شخص خود را محدود به مدل  $PFQR$  کند، به نظر می‌رسد اختیار  $\lambda_2 = 0.4$  منجر به عملکرد بهتری نسبت به بقیه  $\lambda_2$  ها خواهد شد. ملاحظه می‌شود دو مدل  $PRQR$  و  $QR$  هیچ انتخاب متغیری نداشتند. نکته حائز اهمیت این است که در حالت  $PRQR$  که به درجه بهینه از پارامتر تاوان برای اثرهای تصادفی رسیدیم نتایج حاصل از اجرای آن زمانی که توزیع خطاها  $t(4)$  و  $ALD(0, 1, 0.25)$  اختیار شود با نتایج حاصل از اجرای مدل  $QR$  تقریباً برابر است.

طبق نتایج به دست آمده در این بخش و در ادامه پژوهش مرتبط با موضوع این مقاله می‌توان امیدوار بود که هنگام برازش مدل  $DLQR$  بتوان  $\lambda_2$  ای به دست آورد که عملکرد بهتری نسبت به بقیه مدل‌های دیگر داشته باشد. این امر نیاز به بررسی عمیق‌تر (از هر دو دیدگاه نظری و کاربردی) دارد و در فعالیت‌های آتی دنبال خواهد شد.

## ۵ تحلیل مثال واقعی

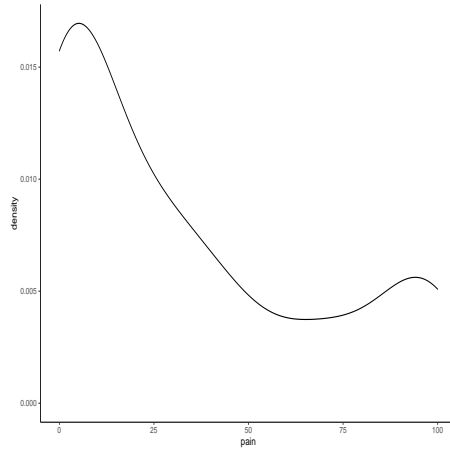
داده‌های آزمایش درد زایمان که توسط داویس (۱۹۹۱) گزارش و توسط جانگ (۱۹۹۶) تحلیل شده است، شامل گزارشات ثبت شده از میزان درد زایمان در زمان‌های متفاوت از آزمودنی‌های مورد بررسی است. تعداد آزمودنی مورد بررسی برابر با  $N=34$  زن در آزمایشگاه است. نحوه ثبت داده‌ها به‌گونه‌ای است که ۱۵ آزمودنی به تصادف انتخاب می‌شوند و در گروه بیماران مصرف‌کننده داروهای درد قرار می‌گیرند و مابقی از آنها دارونما مصرف می‌کنند. میزان درد برای هر آزمودنی، هر ۳۰ دقیقه یک بار اندازه گرفته می‌شود طوری که در نهایت ۶ بار میزان درد ثبت می‌شود. میزان درد (به عنوان متغیر پاسخ) اعدادی بین صفر تا ۱۰۰ اختیار می‌کند که صفر نمایانگر عدم وجود درد و ۱۰۰ نشان‌دهنده میزان درد شدید است. میزان چارک اول، میانه، میانگین و چارک سوم متغیر پاسخ برای گروه اول به ترتیب برابر با ۱۲/۱۸، ۴۲، ۵۰/۹۹، ۹۲/۳۸ و برای گروه دوم که دارونما مصرف کرده اند به ترتیب برابر با ۵، ۹/۵۰، ۱۷/۵۹، ۲۶ است. نمودار جعبه‌ای و برآورد تجربی تابع چگالی برای کل داده‌ها صرف نظر از نوع داروی مصرفی در شکل ۱ رسم شده است. به علاوه، چارک‌های اول، دوم و سوم برای همه زمان‌ها در نمودار جعبه‌ای رسم شد تا نمایش مناسب از توزیع داده‌ها ارائه شود. مشابه توضیحات ارائه شده برای شکل ۱، در شکل ۲ نیز نمایش تغییرات متغیر پاسخ به تفکیک هر دو گروه رسم شده است. همانطور که ملاحظه می‌شود در دوره زمانی اندازه گرفته شده رفتار توزیع متغیر پاسخ (میزان درد) در زمان‌های اندازه گرفته شده تا حد زیادی نامتقارن است. به زبان علمی، در برآورد

جدول ۰۱. نتایج شبیه‌سازی براساس مدل‌های متفاوت و توزیع‌های مختلف خطا.

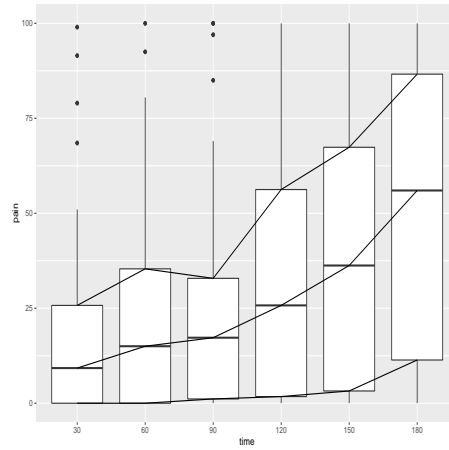
معیارهای ارزیابی مدل		$\lambda_2$	مدل	توزیع خطا
$\ \beta - \beta\ _2$	$\ \beta - \beta\ _1$			
۲/۳۴۵۵	۲/۹۹۵۶	۰/۰۱	DLQR	$N(0, 1)$
۱۴/۶۲۸۷	۷/۵۶۲۰	۰/۰۲		
۱۵/۲۸۶۲	۷/۵۹۸۷	۰/۰۴		
۰/۱۵۸۵	۰/۵۳۰۹	۰/۰۵		
۰/۸۶۶۳	۱/۲۰۴۵	۰/۰۱		
۰/۳۸۷۱	۰/۷۶۶۲	-	PRQR	
۰/۳۱۲۰	۰/۷۴۳۳	۰/۰۱		
۱/۲۴۴۱	۱/۹۰۳۳	۰/۰۲		
۰/۱۱۶۳	۰/۵۱۱۱	۰/۰۴	PFQR	
۰/۳۹۳۵	۰/۷۲۳۰	۰/۰۵		
۰/۸۴۶۸	۱/۱۸۸۸	۰/۰۱		
۰/۳۸۵۳	۰/۷۶۴۹	-	QR	
۰/۲۹۹۹	۰/۸۱۲۰	۰/۰۱		
۱۲/۸۴۲۸	۷/۱۹۱۲	۰/۰۲		
۱۴/۳۷۹۹	۷/۱۰۶۲	۰/۰۴	DLQR	
۰/۵۹۹۰	۱/۰۲۰۳	۰/۰۵		
۱/۳۶۴۲	۱/۵۳۵۷	۰/۰۱		
۰/۵۷۷۷	۱/۲۳۶۳	-	PRQR	$t(4)$
۰/۳۰۷۱	۰/۸۱۹۴	۰/۰۱		
۲/۲۲۱۹	۲/۶۲۴۹	۰/۰۲		
۰/۲۱۳۰	۰/۸۰۷۲	۰/۰۴	PFQR	
۰/۵۸۳۹	۱/۰۰۳۸	۰/۰۵		
۱/۳۰۹۹	۱/۵۰۵۳	۰/۰۱		
۰/۶۱۹۱	۱/۲۴۲۱	-	QR	
۰/۰۸۸۷	۰/۵۸۹۶	۰/۰۱		
۰/۰۸۶۲	۰/۵۹۲۶	۰/۰۲		
۰/۱۸۱۲	۰/۷۲۳۹	۰/۰۴	DLQR	
۰/۲۷۸۲	۰/۸۹۹۱	۰/۰۵		
۱/۲۳۸۱	۱/۴۵۰۲	۰/۰۱		
۰/۴۵۶۲	۰/۸۶۸۱	-	PRQR	$ALD(0, 1, 0.25)$
۰/۲۹۳۵	۰/۷۶۷۴	۰/۰۱		
۰/۴۹۰۸	۱/۳۰۱۰	۰/۰۲		
۰/۱۸۶۴	۰/۷۳۶۶	۰/۰۴	PFQR	
۰/۵۳۷۶	۰/۹۵۵۷	۰/۰۵		
۱/۲۵۲۰	۱/۴۵۹۳	۰/۰۱		
۰/۴۵۴۸	۰/۸۷۲۴	-	QR	

تجربی چگالی‌ها، چوله بودن توزیع میزان درد قابل ملاحظه است. چنین استنباطی از طریق بررسی چندک‌های متغیر پاسخ نیز قابل درک است. نمایش نمودار برای دو گروه نشان می‌دهند که رفتار مورد اشاره برای آزمودنی‌هایی که دارو مصرف کردند مشهودتر است و این یعنی برآورد تجربی چگالی توزیع متغیر پاسخ چوله است. برای داشتن معیار عددی برای میزان چولگی، مقدار عددی آن بدون تفکیک گروه‌ها محاسبه و عدد  $0.83$  به دست آمد که نشان‌دهنده چوله به راست بودن توزیع داده‌ها است. علاوه بر آن، میزان چولگی به تفکیک دو گروه نشان داد که این معیار برای گروه مصرف‌کننده دارونما  $0.05$  و برای مصرف دارو  $1/48$  شد. اعداد گزارش شده نشان‌دهنده این واقعیت هستند

که در گروه اول چولگی توزیع متغیر پاسخ مشاهده نشده است، اما در گروه دوم (مصرف کننده دارو) توزیع متغیر پاسخ به شدت چوله است. این موضوع نشان می‌دهد که رگرسیون چندکی گزینه مناسب‌تری نسبت به رگرسیون مبتنی بر میانگین برای نمایش تغییرات متغیر پاسخ براساس تغییرات ناشی از متغیرهای تبیینی است.



(ب)



(الف)

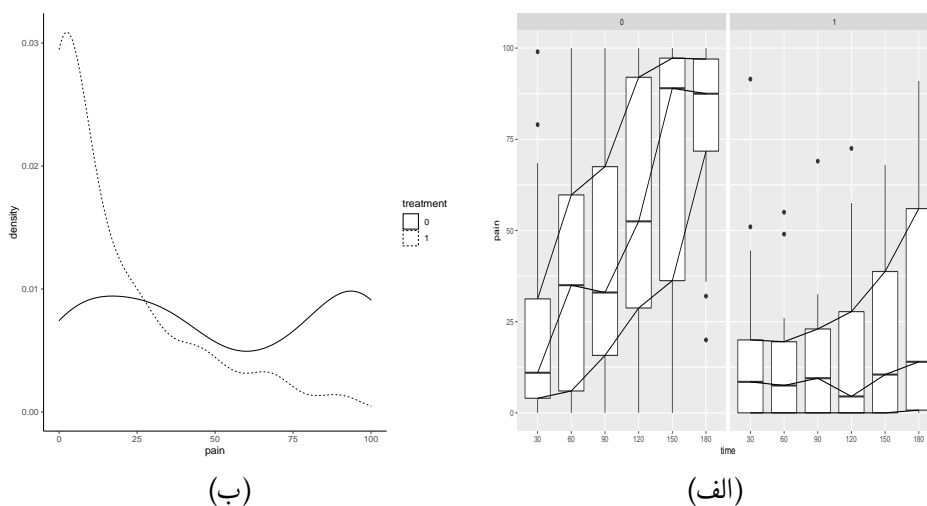
شکل ۱. الف- نمودار جعبه‌ای بمیزان درد آزمودنی‌ها به تفکیک زمان اندازه‌گیری و ب- برآورد تجربی چگالی میزان درد تمام مشاهدات بدون تفکیک گروه‌بندی.

با توجه به این که برای یک آزمودنی در زمان‌های متفاوتی متغیر پاسخ اندازه گرفته شده است، بنابراین در داده‌های مورد بررسی بین مشاهدات هر آزمودنی یک همبستگی ذاتی وجود دارد که با اضافه کردن اثرهای تصادفی به مدل این همبستگی لحاظ می‌شود. لذا، برای لحاظ این واقعیت، از مدل رگرسیونی پیشنهادی توسط **گراسی و بوتای (۲۰۰۷)** به صورت

$$G_{y_{ij}|u_i}(\tau|\mathbf{x}_{ij}, u_i) = \mathbf{x}_{ij}^T \boldsymbol{\beta} + u_i, \quad j = 1, \dots, n_i, \quad i = 1, \dots, 34 \quad (8)$$

استفاده شد. با پیروی از ایشان، متغیرهای تبیینی و پاسخ و پارامترهای مدل در ساختار بندی مبتنی بر رابطه (۸) به صورت

$$\mathbf{X}_i = \begin{bmatrix} 1 & R_i & T_{i1} & R_i T_{i1} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & R_i & T_{in_i} & R_i T_{in_i} \end{bmatrix}, \quad \boldsymbol{\beta} = (\beta_0, \dots, \beta_3)^T, \quad \mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^T$$



شکل ۲. الف- نمودار جعبه‌ای و ب- برآورد تجربی چگالی میزان درد دو گروه آزمودنی دریافت‌کننده دارونما (۰) و دریافت‌کننده دارو (۱).

هستند، که در آن  $y_{ij}$  میزان درد  $i$  امین بیمار در زمان  $j$  است. توجه شود که  $R_i$  نشان‌دهنده مصرف دارو یا دارونما است. در پیاده‌سازی مدل طوری عمل شد که اگر آزمودنی دارو مصرف کند  $R_i$  مقدار ۱ را اختیار می‌کند. به علاوه،  $T_{ij}$  حاصل تقسیم زمان اندازه‌گیری متغیر پاسخ برای هر آزمودنی بر ۳۰ دقیقه است و  $n_i=6$  زمان آخرین اندازه‌گیری برای  $i$  امین آزمودنی است. بنا به دلایلی که در تحلیل اکتشافی (توزیع) داده‌ها مطرح شد یکی از روش‌های پیشنهادی برای تحلیل این داده‌ها مدل رگرسیون چندکی است. از طرفی دیگر، همانطور که در مقدمه تشریح شد برای برازش اولیه مدل رگرسیون چندکی نیازی به لحاظ فرض نرمال بودن توزیع خطاها نیست. برای این منظور، مدل رگرسیونی چندکی (۱) مد نظر قرار می‌گیرد.

برای تحلیل مدل (۱)، گراسی و بوتای (۲۰۰۷) توزیع نرمال را برای توزیع اثرهای تصادفی در نظر گرفتند. بنا به رویکرد مطرح شده در این مقاله، ما فرض می‌کنیم:  $u_i \sim ALD(\phi, \sigma, \tau)$ . سپس از رویکردهای  $DLQR$  و  $PFQR$  برای تاوانیده کردن پارامترهای شیب استفاده می‌شود. علاوه بر آن، رویکرد  $PRQR$  و  $QR$  نیز در برازش مدل (۸) دنبال خواهد شد. نتایج حاصل از برازش مدل‌های مورد اشاره در جدول ۲ آمده است. شایان ذکر است که معیار مناسب برای انتخاب مدل در اینجا  $-2 \log L(\hat{\eta} | \mathbf{y}, \mathbf{u})$  است. طبق نتایج جدول ۲، گزارش برازش مدل برای سه چندک  $\tau = (0.25, 0.5, 0.75)$  ارائه شده است. همانطور که ملاحظه می‌شود هنگام برازش مدل‌های  $DLQR$  و  $PFQR$  انتخاب متغیر صورت گرفته است و برای  $\tau = (0.25, 0.75)$  طبق معیار در نظر گرفته شده به نظر می‌رسد  $\lambda_2=5$  نتیجه بهتری نسبت به  $\lambda_2$  های دیگر دارد. در این مدل‌ها برای این مقدار  $\lambda_2$  میزان  $\hat{\beta}_1$  صفر برآورد شده است و این یعنی مصرف دارو به تنهایی در میزان درد تاثیر ندارد. اما برای  $\tau = 0.5$ ، مقدار  $\hat{\beta}_3$  برابر با صفر برآورد شده است که نشان می‌دهد اثر متقابل داروی مصرفی و زمان اندازه‌گیری در میزان درد اثری ندارد.

این نتیجه را می‌توان اینگونه تفسیر کرد که شخص بیمار دارو را در چه زمانی از ملاحظه درد مصرف کند تاثیر قابل ملاحظه‌ای در میزان کاهش درد او ندارد. در صورت برازش مدل  $PRQR$ ، به نظر می‌رسد برای وقتی که  $\tau = 0.25$  برآوردهای حاصل مناسب‌تر از برآوردهای به دست آمده با انتخاب چندک‌های دیگر باشد. در عوض، هنگام برازش مدل  $QR$ ، نتایج مرتبط با  $\tau = 0.75$  بهتر از چندک‌های دیگر است. اما نکته حائز اهمیت این است که خواه تاوان بر روی اثرات تصادفی اعمال شود خواه نشود، انتخاب متغیر در مدل‌های  $PRQR$  و  $QR$  صورت نمی‌گیرد. در حالت کلی به نظر می‌رسد که داروی مصرفی روی میزان درد اثر معناداری ندارد. به عبارت دیگر بین دو گروه تفاوت معناداری بر اساس میزان درد وجود ندارد. اما زمان اندازه‌گیری میزان درد و اثر متقابل آن با دارو می‌تواند در میزان درد اثر معناداری داشته باشد.

جدول ۰۲. نتایج برازش مدل‌های متفاوت در تحلیل مثال واقعی.

برآورد پارامتر				$-\tau \log L(\hat{\eta} \mathbf{y}, \mathbf{u}.)$	$\lambda_{\tau}$	$\tau$	مدل
$\beta_{\tau}$	$\beta_{\tau}$	$\beta_{\gamma}$	$\beta_{\circ}$				
۰/۰۰۰۰	۰/۰۰۴۲	۰/۰۷۰۴۲	۲/۴۹۷۲	۶۹۵۲۶/۴۷	۰/۰۱		
۰/۰۰۰۰	۰/۰۰۰۰	۰/۰۰۰۰	۲/۴۱۵۴	۶۹۵۴۷/۹۵	۰/۰۵		
۰/۰۰۰۰	۰/۰۰۰۰	۰/۰۰۰۰	۲/۴۱۵۰	۶۶۹۵۵/۵۳	۱	۰/۲۵	
۰/۰۰۷۵	۰/۰۱۴۰	۰/۰۰۰۰	۲/۳۳۷۶	۶۹۳۳۲/۳۵	۵		
۰/۰۰۰۹۳	۰/۰۲۶۹	۰/۰۰۰۱	۱/۹۶۳۵	۳۰۷۵۹/۹۸	۰/۰۱		
۰/۰۰۰۰۱	۰/۰۱۳۰	۰/۰۰۰۱	۲/۰۳۲۸	۳۰۷۹۴/۶۶	۰/۰۵		
۰/۰۰۰۰	۰/۰۰۸۰	۰/۰۰۰۱	۲/۰۶۶۶	۳۰۸۰۲/۹۸	۱	۰/۵	DLQR
۰/۰۰۰۰	۰/۰۰۷۷	۰/۰۱۵	۲/۱۴۴۸	۳۰۷۹۲/۲۹	۵		
۰/۰۰۰۰	۶/۳۹۹۵	۰/۶۰۳۶	۰/۶۰۵۸	۲۳۳۵۱/۲۴	۰/۰۱		
۰/۰۱۶۶۷	۶/۳۳۵۱	۰/۰۰۰۰	۰/۳۳۹۲	۸۵۵۳۲/۹۳	۰/۰۵		
۰/۰۱۶۶۷	۶/۳۳۵۱	۰/۰۰۰۰	۰/۳۳۹۲	۲۳۰۹۳/۰۱	۱	۰/۷۵	
۰/۰۱۶۶۷	۶/۳۳۵۱	۰/۰۰۰۰	۰/۳۳۹۲	۸۸۱۳۷/۸	۵		
۰/۰۹۸۷۴	۰/۰۰۴۱	۰/۰۰۰۰	۲/۴۹۱۴	۷۲۵۲۹/۹۱	۰/۰۱		
۰/۰۴۲۴	۰/۰۰۱۵	۰/۰۰۰۰	۲/۳۳۵۶	۷۰۱۰۹/۶۶	۰/۰۵		
۰/۳۶۳۷	۰/۰۰۷۸	۰/۰۰۰۰	۲/۳۹۸۷	۷۰۲۰۵/۸۸	۱	۰/۲۵	
۰/۰۲۰۵	۰/۰۰۰۵	۰/۰۰۰۰	۲/۳۲۸۹	۷۰۰۲۲/۰۶	۵		
۰/۱۰۹۴	۰/۰۱۸۸	۰/۰۰۰۱	۱/۹۵۷۴	۳۱۴۰۹/۲۳	۰/۰۱		
۰/۰۰۰۰	۰/۰۱۴۴	۰/۱۷۳۶	۲/۰۲۴۴	۳۱۴۱۶/۸۹	۰/۰۵		
۰/۰۰۰۰	۰/۰۰۵۱	۰/۰۰۰۰	۲/۰۴۷۰	۳۱۴۲۶/۲۱	۱	۰/۵	PFQR
۰/۰۰۰۰	۰/۰۰۰۰	۰/۳۹۵۹	۲/۱۱۶۲	۳۱۴۲۸/۴۸	۵		
۰/۰۰۰۰	۶/۳۹۸۳۱	۰/۲۳۷۴	۰/۶۱۳۴	۲۴۱۲۴/۶۰	۰/۰۱		
۰/۰۴۱۳	۶/۳۸۶۴۵	۰/۰۰۰۰	۰/۵۶۹۸	۲۴۰۶۴/۸۵	۰/۰۵		
۰/۰۰۰۰	۶/۳۵۷۲۳	۰/۰۰۰۰	۰/۶۶۱۰	۲۴۰۴۲/۰۱	۱	۰/۷۵	
۰/۰۰۰۲	۰/۰۰۰۱۴	۰/۰۶۸۷	۱/۸۷۶۹	۱۸۱۶۹/۷۰	۵		
۰/۰۸۳۹۵	۰/۰۰۳۷۹	۰/۰۹۳۰۵	۲/۳۳۱۸	۱۶۷۲۵۳/۷۴	-	۰/۲۵	
۰/۰۵۲۵۵	۰/۰۱۹۳۴	۱/۸۹۳۷	۱/۹۸۳۰	۳۱۲۵۵/۵۷	-	۰/۵	PRQR
۰/۰۵۲۵۵	۰/۰۱۹۳۴	۱/۸۹۳۷	۱/۹۸۳۰	۳۵۵۷۴/۹۲	-	۰/۷۵	
۰/۱۱۹۹۰	۰/۰۰۴۵	۰/۰۲۷۴۷	۲/۴۹۶۰	۷۳۱۳۰/۸۵	-	۰/۲۵	
۰/۰۷۵۱۰	۰/۰۲۰۰	۲/۴۰۲۹	۱/۹۶۱۹	۳۲۱۲۷/۸۰	-	۰/۵	QR
۲/۵۰۹۹	۶/۴۰۰۱	۰/۹۱۱۵	۰/۶۱۳۳	۲۶۰۷۶/۴۹	-	۰/۷۵	

## بحث و نتیجه‌گیری

مدل ارائه شده در این مقاله یک مدل تاوانیده دوگانه است که تاوان علاوه بر اینکه تابعی از اثرهای تصادفی است، تابعی از پارامترهای مدل هم هست. در مقایسه با روش ارائه شده در **کوئنکر (۲۰۰۴)** که درجه انقباض اثرها وابسته به انتخاب پارامتر تاوان بود در این مقاله به درجه بهینه‌ای از انقباض اثرها با استفاده از داده‌های شبیه‌سازی شده و داده‌های واقعی دست یافتیم که در نوع خود پژوهش جدیدی در حوزه مورد مطالعه مرتبط با مدل بندی رگرسیونی چندکی آمیخته است. آنچه در ادامه پژوهش مطالعه خواهد شد بررسی نحوه انتخاب پارامتر تاوان مناسب از هر دو دیدگاه نظری و کاربردی است.

## تقدیر و تشکر

نویسندگان مقاله از پیشنهادات و نظرات داوران و ویراستار محترم مجله که در بهبود کیفیت مقاله بسیار مؤثر واقع شد، تقدیر و تشکر می‌نمایند.

## مراجع

آقامحمدی، ع. و محمدی، س. (۱۳۹۴)، رگرسیون چندکی بیزی با تاوان لاسو و لاسوی تطبیق‌پذیر برای داده‌های طولی دودویی. *مجله علوم آماری*، ۹، ۱۶۷-۱۴۹.

Booth, J. G. and Hobert, J. P. (1999), Maximizing Generalized Linear Mixed Model Likelihoods with an Automated Monte Carlo EM Algorithm. *Journal of the Royal Statistical Society*, **61**, 265-285.

Bondell, H. D, Krishna, A. and Ghosh, S. K. (2010), Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effect Models, *Biometrics*, **66**, 1069-1077.

Breiman, L. (1995), Better Subset Selection Using Nonnegative Garrote, *Techonometrics*, **37**, 373-384.

Cowles, M. K. , and Carlin, B. P. (1996), Markov Chain Monte Carlo Convergence Diagnostics: A Comparative Review, *Journal of the American Statistical Association*, **91**, 883-904.



- Davis, C. S. (1991), Semi-Parametric and Non-Parametric Methods for the Analysis of Repeated Measurements with Applications to Clinical Trials., *Statistics in Medicine*, **10**, 1959–1980.
- Diggle, P. J., Heagerty, P., Liang, K. Y. and Zeger, S. L. (2002), *Analysis of Longitudinal Data*, 2nd ed. Oxford: Oxford University Press.
- Fan, J. and Li, R. (2001), Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties, *Journal of the American Statistical Association*, **96**, 1348–1360.
- Geraci, M. , and Bottai, M. (2007), Quantile Regression for Longitudinal Data Using the Asymmetric Laplace Distribution. *Biostatistics*, **8**, 140-154.
- Gilks, W. R. (1995), Derivative-Free Adaptive Rejection Sampling for Gibbs Sampling. *Bayesian Statistics*, **4**, 641–649.
- Gilks, W. R., Best, N. G. and Tan, K. K. C. (1995), Adaptive Rejection Metropolis Sampling within Gibbs Sampling, *Applied Statistics*, **44**, 455-472.
- JUNG, S. (1996), Quasi-Likelihood for Median Regression Models, *Journal of the American Statistical Association*, **91**, 251–257.
- Koenker, R., and Bassett, G. (1978), Regression Quantiles. *Econometrica*, **46**, 33-50.
- Koenker, R. (2004), Quantile Regression for Longitudinal Data. *Journal of Multivariate Analysis*, **91**, 74-89.
- Li, H., Liu, Y. and Luo, Y. (2020), Double Penalized Quantile Regression for the Linear Mixed Effects Model, *Journal of Systems Science and Complexity*, **33**, 2080–2102.
- Tibshirani, R. (1996), Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society*, **58**, 267–88.
- Zou, H. (2006), The Adaptive Lasso and Its Oracle Properties, *Journal of the American Statistical Association*, **101**, 1418–1429.