

رگرسیون چند کی بیزی با تاوان لاسو سازوار برای داده‌های پانلی پویا

علی آقامحمدی^۱، سکینه محمدی^۲

تاریخ دریافت: ۱۳۹۳/۱۲/۲۳

تاریخ پذیرش: ۱۳۹۵/۱۲/۲۵

چکیده:

مدل‌های داده‌های پانلی پویا قسمت مهمی از مطالعات حوزه‌های پزشکی، اجتماعی و اقتصادی را شامل می‌شوند. ویژگی بارز این مدل‌ها وجود متغیر وابسته تأخیری به‌عنوان متغیر توصیفی است. مشکل برآورد در این مدل‌ها از همبستگی بین متغیر وابسته تأخیری و مؤلفه خطای فعلی ناشی می‌شود. اخیراً رگرسیون چند کی تاوانیده برای تحلیل داده‌های پانلی پویا مورد توجه قرار گرفته است. در این مقاله نخست مدل رگرسیون چند کی با ایجاد تاوان لاسو سازوار روی اثرهای تصادفی برای داده‌های پانلی پویا با فرض وابستگی اثرهای تصادفی و مشاهدات اولیه ارائه می‌شود. همچنین این مدل با فرض استقلال بین اثرهای تصادفی و مشاهدات اولیه نیز بررسی خواهد شد. هر دو مدل از دیدگاه آمار بیزی بیان شده، مورد تحلیل قرار می‌گیرند. چون در این دو روش، توزیع پسین پارامترها به شکل بسته قابل حصول نیست، توزیع‌های پسین شرطی کامل پارامترها محاسبه و از الگوریتم نمونه‌گیری گیبز برای استنباط استفاده می‌شود. برای مقایسه کارایی روش‌های بیزی ارائه‌شده با روش‌های متداول، مطالعه شبیه‌سازی انجام شده و در پایان نیز روش استفاده از مدل‌ها در قالب مثال کاربردی شرح داده خواهد شد.

واژه‌های کلیدی: استنباط بیزی، تاوان لاسو سازوار، توزیع لاپلاس نامتقارن، داده‌های پانلی پویا، رگرسیون چند کی، نمونه‌گیری گیبز.

۱ مقدمه

روش‌های خاص خود را می‌طلبند. معمولاً برای تحلیل این نوع داده‌ها از مدل‌های خطی پانلی پویا با اثرهای تصادفی استفاده می‌شود. اما در این مدل‌ها برآوردگرهای کمترین مربعات معمولی اغلب ناسازگار و دارای اربیبی قابل ملاحظه‌ای هستند. روش‌های برآوردی مختلفی که برآوردگرهای سازگار و ناریب در این مدل‌ها دارند نیز پیشنهاد شده است [۹]. اما در سال‌های اخیر مدل‌های رگرسیونی چند کی برای داده‌های پانلی پویا مورد توجه قرار گرفته‌اند [۷، ۶].

رگرسیون چند کی^۴ اولین بار در [۱۲] ارائه شد. این روش نسبت به روش رگرسیون میانگین^۵ دارای دو مزیت مهم است. نخست این که رگرسیون چند کی نسبت به ناپایداری واریانس و داده‌های دورافتاده حساس نیست و دوم این که اطلاعات

مطالعه داده‌های پانلی پویا^۳ طرحی در تحقیقات علوم مختلف نظیر علوم اقتصادی و اجتماعی است که به‌سبب نیاز محققان به این نوع مطالعه‌ها مورد استفاده قرار می‌گیرد. در این نوع مطالعات، متغیر پاسخ در چندین نوبت متوالی مشاهده می‌شود، یا به عبارت دیگر یک واحد آزمایشی تحت اندازه‌گیری مکرر در طول زمان قرار می‌گیرد. ویژگی بارز مدل داده‌های پانلی پویا وجود متغیر تأخیری به‌عنوان یک متغیر توصیفی است. چون متغیر پاسخ تحت تأثیر اندازه خود در دوره قبل است، این متغیر در دوره قبل به‌عنوان متغیر توصیفی برای پاسخ‌های فعلی وارد مدل می‌شود. واضح است تحلیل این نوع داده‌ها با توجه به نوع متغیر پاسخ،

^۱ استادیار گروه آمار دانشگاه زنجان، ایران

^۲ دانش‌آموخته کارشناس ارشد آمار دانشگاه زنجان، ایران

^۳ Daynamic Panel Data

^۴ qauntile regression

^۵ mean regression

را از دیدگاه آمار بیزی بررسی کردند. در این روش، با فرض این که توزیع اثرهای تصادفی لاپلاس متقارن هستند، مدل کوئنکر [۱۱] حاصل خواهد شد. گالوائو و مونتس روخاس [۷] مدل ارائه شده توسط کوئنکر [۱۱] را برای داده‌های پانلی پویا با اثرهای ثابت از دیدگاه آمار بسامدی مطالعه کردند. گالوائو [۶] رگرسیون چندکی برای داده‌های پانلی پویا را با در نظر گرفتن اثرهای تصادفی و استفاده از متغیرهای ابزاری^۱ بدون در نظر گرفتن تاوان لاسو بررسی کرد. کوباباشی و کوزومی [۱۰] نیز داده‌های پانلی پویای سانسور شده را بدون در نظر گرفتن تاوان لاسو از دیدگاه آمار بیزی بررسی کردند.

هدف این مقاله تحلیل داده‌های پانلی پویا با استفاده از روش رگرسیون چندکی با در نظر گرفتن تاوان لاسو سازوار روی اثرهای تصادفی از دیدگاه آمار بیزی است. برای این منظور، تاوان لاسو سازوار از طریق تعریف توزیع پیشین مناسب روی اثرهای تصادفی در مدل لحاظ می‌شود و مدل از دیدگاه آمار بیزی تحلیل می‌شود. بنا بر این بخش دوم، شامل مطالب کلی در مورد مدل رگرسیون چندکی و رگرسیون چندکی در داده‌های پانلی پویا است. در بخش سوم، مدل رگرسیون چندکی با تاوان لاسو سازوار از دیدگاه آمار بیزی بررسی می‌شود و در بخش چهارم نیز با روش شبیه‌سازی، کارایی مدل ارائه شده با مدل‌های دیگر مورد مقایسه قرار می‌گیرد. بخش پنجم شامل تحلیل داده‌های واقعی و ارزیابی کارایی مدل‌ها است.

۲ رگرسیون چندکی

مدل رگرسیون خطی

$$y_i = x_i' \beta + \epsilon_i, \quad i = 1, \dots, n \quad (1)$$

را در نظر بگیرید که در آن بردار x_i بردار مربوط به متغیرهای توصیفی، β یک بردار $(k \times 1)$ بعدی از پارامترها و ϵ_i مؤلفه خطا است. در مدل رگرسیونی میانگین، هدف استنباط در مورد متوسط مقدار پاسخ به‌ازای مقادیر متفاوت متغیرهای توصیفی است؛ لذا $E(y_i | x_i) = x_i' \beta$ مورد توجه است. اما در رگرسیون

جزئی‌تری نسبت به تأثیر متغیرهای توصیفی در چندک‌های مختلف توزیع متغیر پاسخ ارائه می‌کند. همان‌طور که عنوان شد در مطالعه داده‌های پانلی، علاوه بر اثرهای ثابت که تغییرات درون گروه‌ها (تغییرات در طول زمان) را کنترل می‌کنند، اثرهای تصادفی بین متغیرهای پاسخ نیز در مدل لحاظ می‌شود که تغییرات بین گروهی را کنترل می‌کنند. اما در بسیاری از مسائل، تعداد مشاهدات از متغیر پاسخ نسبت به مشاهدات انجام گرفته در طول زمان از متغیر مربوط بسیار زیاد است. بنا بر این در این مدل تعداد اثرهای تصادفی (تعداد پارامترها) نیز خیلی زیاد می‌شود و لذا تعبیر و تفسیر مدل، مشکل می‌شود و دقت آن نیز کاهش می‌یابد. برای رفع این مشکل، کوئنکر [۱۱] رگرسیون چندکی را برای تحلیل داده‌های پانلی با افزودن تاوان لاسو^۶ (تاوان با نرم l_1) روی اثرهای تصادفی بررسی کرد. در واقع این روش با ایجاد تاوان، اثرهای تصادفی کم‌اهمیت را به سمت صفر منقبض^۷ کرده و مدلی تنک^۸ ایجاد می‌کند. اولین بار تیبشیرانی [۱۷] رگرسیون لاسو را که با ایجاد تاوان لاسو روی پارامترها انجام می‌شود ارائه کرد. سپس زو [۲۰] رگرسیون لاسو سازوار^۹ را به‌منظور گسترش رگرسیون لاسو پیشنهاد کرد. این روش بر خلاف تاوان لاسو که برای همه ضرایب، تاوان یکسانی در نظر گرفته می‌شود، اندازه تاوان‌های مختلفی برای ضرایب رگرسیونی متفاوت لحاظ می‌کند. در واقع بر خلاف تاوان لاسو، این روش بر حسب میزان اهمیت هر یک از ضرایب رگرسیونی، تاوانی متناسب با آن را در نظر می‌گیرد. بررسی‌های زو [۲۰] مشخص کرد که در رگرسیون چندکی با تاوان لاسو سازوار، اریبی در برآورد پارامترها در مقایسه با روش لاسو کمتر است. الحمزوی و همکاران [۲۱] مدل رگرسیون چندکی و لنگ و همکاران [۱۴] مدل رگرسیون میانگین با تاوان لاسو سازوار را از دیدگاه آمار بیزی معرفی کردند و نشان دادند که این تاوان در مقایسه با تاوان لاسو دقت زیادی دارد.

کوئنکر [۱۱] و گراسی و باتای [۸] با تعریف توزیع خاص برای اثرهای تصادفی و با در نظر گرفتن توزیع لاپلاس نامتقارن برای متغیرهای پاسخ، مدل رگرسیونی چندکی در داده‌های پانلی

^۶ Lasso

^۷ shrink

^۸ spares

^۹ adaptive Lasso

^{۱۰} instrumental variables

چندکی، تابع چندک شرطی

$$Q_{\tau}(y_i|x_i) = x_i'\beta, \quad i = 1, \dots, n$$

موارد، تحلیل رگرسیون چندکی از دیدگاه آمار بیزی با استفاده از این توزیع بررسی می‌شود. اولین بار یو و مویید [۱۹] با فرض $\sigma = 1$ و تعریف توزیع پیشین ناآگاهی‌بخش برای β ، $(\pi(\beta) \propto 1)$ رگرسیون چندکی بیزی را در داده‌های مقطعی معرفی و مطالعه کردند. هدف این مقاله تعمیم این روش به مدل داده‌های پانلی پویا با اثرهای تصادفی و با فرض تاوان لاسو سازوار روی اثرهای تصادفی از دیدگاه آمار بیزی است. مدل رگرسیونی داده‌های پانلی پویا، یک مدل آمیخته خطی به صورت

$$y_{it} = x_{it}'\beta + \gamma y_{i(t-1)} + \alpha_i + \epsilon_{it}, \\ i = 1, \dots, n, \quad t = 1, \dots, T, \quad (4)$$

است که در آن y_{it} متغیر پاسخ برای i امین واحد آزمایشی در زمان t ، x_{it} بردار $(k \times 1)$ بعدی از متغیرهای توصیفی، $y_{i(t-1)}$ متغیر پاسخ تأخیری، γ و β ضرایب رگرسیونی و α_i اثر تصادفی مربوط به پاسخ y_{it} را نشان می‌دهد. توجه کنید که در این مدل α_i ها اثرهای بین گروهی و β اثرهای درون گروهی را کنترل می‌کنند. وجود $y_{i(t-1)}$ بیان‌کننده این است که متغیر پاسخ تحت تأثیر مقدار خود در دوره قبل نیز قرار دارد. با فرض این که چندک τ ϵ_{it} برابر با صفر است، تابع چندک شرطی در این مدل به صورت

$$Q_{\tau}(y_{it}|\cdot) = x_{it}'\beta + \gamma y_{i(t-1)} + \alpha_i \quad (5)$$

به دست می‌آید. کوئنکر [۱۱] مدل فوق را بدون در نظر گرفتن متغیرهای تأخیری و با فرض تاوان لاسو روی α_i ها به صورت $\lambda \sum_{i=1}^n |\alpha_i|$ مورد بررسی قرار داد. او پارامترهای مدل یعنی α و β را از مینیم کردن تابع زیان توانیده

$$\sum_{k=1}^q \sum_{i=1}^n \sum_{t=1}^T w_k \rho_{\tau_k}(y_{it} - x_{it}'\beta_{\tau_k} - \alpha_i) + \lambda \sum_{i=1}^n |\alpha_i|,$$

نسبت به α و β برآورد کرد و همچنین استدلال کرد که اگر n در مقایسه با T خیلی بزرگ باشد، تابع زیان توانیده فوق، نسبت به روش‌های غیر توانیده عملکرد بهتری دارد. در ادامه گالواتو و موتنس روخاس [۷] مدل رابطه را با فرض تاوان لاسو روی α_i ها مورد بررسی قرار دادند. در واقع آن‌ها پارامترهای α ، β و γ را از مینیم کردن تابع زیان توانیده

$$\sum_{k=1}^q \sum_{i=1}^n \sum_{t=1}^T w_k \rho_{\tau_k}(y_{it} - x_{it}'\beta_{\tau_k} - \gamma y_{i(t-1)} - \alpha_i) + \lambda \sum_{i=1}^n |\alpha_i|$$

^{۱۱} asymmetric Laplace distribution

مورد توجه است که در آن $Q_{\tau}(\cdot)$ معکوس تابع توزیع تجمعی متغیر پاسخ y_i به شرط معلوم بودن بردار x_i است. بنا بر این $x_i'\beta$ همان چندک τ ام متغیر y_i را نشان می‌دهد. لذا برای مدل رابطه ۱ در رگرسیون میانگین، فرض می‌شود که میانگین توزیع خطاها صفر است، اما در رگرسیون چندکی، چندک τ ام ϵ_i ها برابر با صفر است؛ یعنی $\int_{-\infty}^{\tau} f_{\epsilon_i}(t) dt = \tau$. در مدل رگرسیون چندکی، ضرایب رگرسیونی یعنی β از مینیم کردن عبارت

$$\sum_{i=1}^n \rho_{\tau}(y_i - x_i'\beta) \quad (2)$$

نسبت به β برآورد می‌شوند، که در آن $\rho_{\tau}(\cdot)$ نشان‌دهنده تابع زیان و به صورت $\rho_{\tau}(u) = \frac{|u| + (\tau - 1)u}{2}$ تعریف می‌شود (کوئنکر و باست [۱۲]). چون رابطه ۲ در مبدأ مشتق پذیر نیست، راه‌حل تحلیلی برای برآورد ضرایب رگرسیونی وجود ندارد. بنا بر این معمولاً از روش‌های بهینه‌سازی عددی برای برآورد β استفاده می‌شود و از این طریق هیچ‌گونه استنباط آماری نیز نمی‌توان در مورد برآوردگر β انجام داد. برای رفع این مشکل، یو و مویید [۱۹] رگرسیون چندکی را از دیدگاه آمار بیزی با استفاده از توزیع لاپلاس نامتقارن^{۱۱} مورد بحث و بررسی قرار دادند. متغیر تصادفی y را دارای توزیع لاپلاس نامتقارن گویند $(y \sim ALD(\mu, \sigma, \tau))$ هرگاه تابع چگالی آن به صورت

$$f(y) = \frac{\tau(1-\tau)}{\sigma} \exp \left\{ -\rho_{\tau} \left(\frac{y-\mu}{\sigma} \right) \right\}$$

باشد که در آن $0 < \tau < 1$ پارامتر چولگی، σ پارامتر مقیاس و $-\infty < \mu < \infty$ پارامتر مکان است. لذا با فرض $y_i \sim ALD(x_i'\beta, \sigma, \tau)$ تابع درست‌نمایی بردار $\mathbf{y} = (y_1, \dots, y_n)$ به صورت

$$L(\beta|\mathbf{y}, \mathbf{x}) \propto \sigma^{-n} \times \exp \left\{ -\sum_{i=1}^n \rho_{\tau} \left(\frac{y_i - x_i'\beta}{\sigma} \right) \right\} \quad (3)$$

به دست می‌آید. بنا بر این ما کسیم کردن رابطه ۳ در حضور پارامتر مزاحم σ با مینیم کردن رابطه ۲ نسبت به β معادل است. در واقع با در نظر گرفتن توزیع لاپلاس نامتقارن برای y_i ، تابع درست‌نمایی برای مدل رگرسیون چندکی به دست می‌آید که در تحلیل‌های بیزی ضروری است. به همین دلیل در بسیاری از

با میانگین $\frac{\sigma}{\tau(1-\tau)}$ و مستقل از هم باشند،

$$u \stackrel{d}{=} k_1 e + \sqrt{2\sigma e z}$$

که در آن $k_1 = (1 - 2\tau)$ است [۱۳]. با استفاده از این خاصیت توزیع لاپلاس نامتقارن و با فرض این که $\epsilon_{it} \sim ALD(0, \sigma, \tau)$ باشد، مدل رابطه ۴ را می‌توان به صورت

$$y_{it} = \mathbf{x}'_{it}\boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i + k_1 e_{it} + \sqrt{2\sigma e_{it}} z_{it} \quad (۸)$$

بازنویسی کرد، که در آن e_{it} و z_{it} به ترتیب دارای توزیع نرمال استاندارد و نمایی با میانگین $\frac{\sigma}{\tau(1-\tau)}$ و مستقل از هم هستند. بنا بر این توزیع شرطی y_{it} به شرط e_{it} ، عضو خانواده توزیع‌های نرمال است و لذا تحلیل‌ها را می‌توان با استفاده از رگرسیون خطی میانگین انجام داد. به همین طریق برای سادگی، توزیع پیشین روی α_i ها را که توزیع لاپلاس متقارن است به صورت توزیعی آمیخته از دو توزیع نرمال و نمایی به صورت

$$\alpha_i | s_i \sim N(0, s_i), \quad s_i | \nu_i \sim \text{Exp}\left(\frac{\nu_i}{\lambda}\right) \quad (۹)$$

در نظر می‌گیریم که در آن $\text{Exp}\left(\frac{\nu_i}{\lambda}\right)$ توزیع نمایی با میانگین $\frac{\lambda}{\nu_i}$ همان توزیع آمیختگی است [۱۶].

در برازش مدل‌های رگرسیونی پویا با اثرهای تصادفی به داده‌های پانلی پویا، وابستگی یا استقلال اثرهای تصادفی از مشاهدات اولیه موضوعی اساسی است که به آن مسئله شرایط اولیه گویند [۱۴ و ۱۸]. در این مقاله مدل را با هر دو فرض استقلال و وابستگی در نظر می‌گیریم و از دیدگاه آمار بیزی بررسی قرار می‌کنیم.

۱.۳ مدل با فرض استقلال α_i و y_i .

در این بخش فرض می‌کنیم که مشاهدات اولیه از اثرهای تصادفی مستقل هستند. چون هدف تحلیل بیزی است، برای ابرپارامتر ν_i در رابطه ۹ توزیع پیشین مزدوج را نمایی به صورت $\pi(\nu_i | \phi) = \phi \exp\{-\phi \nu_i\}$ و توزیع‌های پیشین برای پارامترهای σ و β را نیز مزدوج به ترتیب، معکوس گاما و نرمال چندمتغیره در نظر می‌گیریم. توزیع پیشین ϕ را به منظور کم کردن تأثیر آن بر استنباط‌ها، توزیع پیشین پخ^{۱۲} به صورت $\pi(\phi) \propto \frac{1}{\phi}$ فرض می‌کنیم. حال با توجه به مدل رابطه ۸ و ۹ و توزیع‌های پیشین

نسبت به این پارامترها برآورد کردند. آنها نشان دادند این روش، زمانی که متغیر پاسخ تحت تأثیر مقادیر دوره قبل است، عملکرد بهتری دارد. هدف این مقاله تحلیل مدل رابطه ۴ با تاوان لاسو سازوار روی اثرهای تصادفی از دیدگاه آمار بیزی است. با در نظر گرفتن این تاوان روی اثرهای تصادفی، تابع زیان تاوانیده به صورت

$$\sum_{i=1}^n \sum_{t=1}^T \rho_{\tau}(y_{it} - \mathbf{x}'_{it}\boldsymbol{\beta}_{\tau} - \gamma y_{i(t-1)} - \alpha_i) + \sum_{i=1}^n \lambda_i |\alpha_i| \quad (۶)$$

به دست می‌آید که در آن عبارت $\sum_{i=1}^n \lambda_i |\alpha_i|$ تاوان لاسو سازوار است و به λ_i پارامتر تنظیم گویند. توجه کنیم که در این روش برای هر α_i اندازه تاوان متفاوت λ_i در نظر گرفته شده است. هرچه مقدار این پارامتر بزرگتر باشد، انقباض به سمت صفر نیز بیشتر بوده، اثرهای تصادفی کم‌اهمیت از مدل حذف خواهند شد.

۳ مدل بیزی با تاوان لاسو سازوار

همان‌طور که عنوان شد در مدل رگرسیون چندکی با تاوان لاسو سازوار در داده‌های پانلی پویا، تابع زیان تاوانیده با تاوان لاسو سازوار روی اثرهای تصادفی، به صورت رابطه ۶ است. لذا اگر فرض کنیم در رابطه ۴، $\epsilon_{ij} \sim ALD(0, \sigma, \tau)$ و توزیع پیشین برای $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)$ را به صورت

$$\pi(\boldsymbol{\alpha} | \lambda, \sigma) = \prod_{i=1}^n \frac{\lambda_i}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{\sum_{i=1}^n \lambda_i |\alpha_i|}{\sqrt{\sigma}}\right\}$$

در نظر بگیریم، توزیع پسین بردار $\boldsymbol{\alpha}$ به صورت

$$\pi(\boldsymbol{\alpha} | \boldsymbol{\beta}, \mathbf{y}, \mathbf{x}, \sigma, \lambda) \propto \exp\left\{-\sum_{i=1}^n \sum_{t=1}^T \frac{1}{\sigma} \rho_{\tau}(y_{it} - \mathbf{x}'_{it}\boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i)\right\} \times \exp\{-\sum_{i=1}^n \nu_i |\alpha_i|\}, \quad (۷)$$

به دست می‌آید که در آن $\nu_i = \frac{\lambda_i}{\sqrt{\sigma}}$ است. بنا بر این مینیمم کردن تابع هدف رابطه ۶ معادل ماکسیمم کردن توزیع پسین $\boldsymbol{\alpha}$ در رابطه ۷ در حضور پارامتر مزاحم σ است. لذا می‌توان توزیع لاپلاس نامتقارن را برای این مدل مورد استفاده قرار داد. چون استفاده از تابع درست‌نمایی در توزیع لاپلاس نامتقارن به دلیل وجود قدر مطلق آسان نیست، در تحلیل‌های بیزی معمولاً از شکل آمیخته این توزیع استفاده می‌شود. اگر $u \sim ALD(0, \sigma, \tau)$ و متغیرهای z و e به ترتیب دارای توزیع نرمال استاندارد و نمایی

تعریف شده، مدل سلسله‌مراتبی را می‌توان به صورت

$$y_{it}|e_{it}, \beta, \sigma, \alpha_i \sim N(x'_{it}\beta + \gamma y_{i(t-1)} + \alpha_i + (1 - \tau)e_{it}, \tau\sigma e_{it})$$

به قسمی که

$$\begin{aligned} y_i|\alpha_i &\sim N(\phi_0 + \phi_1\alpha_i, \sigma^2), & \alpha_i|s_i &\sim N(\cdot, s_i), \\ s_i|\nu_i &\sim \text{Exp}\left(\frac{\nu_i}{\tau}\right), & \gamma|\psi^2 &\sim N(\cdot, \psi^2), \psi^2 \sim IG(a, b), \\ e_{it}|\sigma &\sim \text{Exp}\left(\frac{\tau(1-\tau)}{\sigma}\right), & \beta &\sim N_k(\mathbf{b}, \mathbf{B}), \\ \nu_i|\phi &\sim \text{Exp}(\phi), & \pi(\phi) &\propto \frac{1}{\phi}, \sigma|c., d. \sim IG(c., d.), \\ \phi_0 &\sim N(\mu_{\phi_0}, \sigma_{\phi_0}^2), & \phi_1 &\sim N(\mu_{\phi_1}, \sigma_{\phi_1}^2), \\ \sigma^2 &\sim IG\left(\frac{v_0}{\tau}, \frac{\delta_0}{\tau}\right), \end{aligned}$$

نوشت. در این مدل نیز توزیع پسین کامل پارامترها و ابرپارامترها به صورت بسته در دسترس نیست. اما با اندکی محاسبات جبری، توزیع پسین شرطی کامل پارامترها و ابرپارامترها قابل حصول است، لذا می‌توان از روش نمونه‌گیری گیبز برای استنباط بیزی استفاده کرد. توزیع‌های پسین پارامترهای $\beta, \gamma, \sigma, \delta, \nu, \phi, e$ و ψ^2 مانند بخش ۱.۳ است؛ اما توزیع پسین شرطی کامل سایر پارامترها در پیوست ب ارائه شده است.

۴ مطالعه شبیه‌سازی

در این بخش برای مقایسه کارایی روش‌های بیزی ارائه شده با روش‌های متداول دیگر به مطالعه شبیه‌سازی، می‌پردازیم. در این شبیه‌سازی مدل را به صورت

$$y_{it} = x_{it}^{(1)}\beta_1 + x_{it}^{(2)}\beta_2 + \gamma y_{i(t-1)} + \alpha_i + \epsilon_{it},$$

$$i = 1, \dots, n, \quad t = 1, \dots, T$$

در نظر می‌گیریم که در آن $n = 100$ و $T = 5$ است. توجه کنیم که تعداد متغیرهای پاسخ در مقایسه با متغیرهای توصیفی زیاد است. لذا استفاده از تاوان لاسو سازوار منجر به بهبود برآورد پارامترها و دقت مدل خواهد شد. متغیرهای توصیفی را به صورت $x_{it}^{(j)} \sim N(\cdot, 1)$ و $j = 1, 2$ و α_i ها را مانند لیو و همکاران [۱۵] از توزیع $N(\cdot, 4)$ تولید کرده، مؤلفه‌های خطا را نیز از سه توزیع به صورت

$$\epsilon_{it} \sim N(\cdot, 1), \quad \epsilon_{it} \sim t(3), \quad \epsilon_{it} \sim \text{unif}(-2, 2)$$

$$y_{it}|e_{it}, \beta, \sigma, \alpha_i \sim N(x'_{it}\beta + \gamma y_{i(t-1)} + \alpha_i + (1 - \tau)e_{it}, \tau\sigma e_{it})$$

به قسمی که

$$\begin{aligned} \gamma|\psi^2 &\sim N(\cdot, \psi^2), & \psi^2 &\sim IG(a, b), \\ e_{it}|\sigma &\sim \text{Exp}\left(\frac{\tau(1-\tau)}{\sigma}\right), & \beta &\sim N_k(\mathbf{b}, \mathbf{B}), \\ \alpha_i|s_i &\sim N(\cdot, s_i), & s_i|\nu_i &\sim \text{Exp}\left(\frac{\nu_i}{\tau}\right), \\ \nu_i|\phi &\sim \text{Exp}(\phi), & \pi(\phi) &\propto \frac{1}{\phi}, \sigma|c., d. \sim IG(c., d.) \end{aligned}$$

نوشت، که در آن $IG(a, b)$ نشان‌دهنده توزیع گامای معکوس با پارامترهای a و b است. با توجه به مدل سلسله‌مراتبی فوق، توزیع پسین همه پارامترها و ابرپارامترها به صورت بسته قابل حصول نیست. اما می‌توان توزیع‌های پسین شرطی کامل کلیه پارامترها و ابرپارامترها را محاسبه کرد که در پیوست الف ارائه شده است. همه توزیع‌های پسین شرطی کامل پارامترها و ابرپارامترها، توزیع‌های شناخته شده‌ای هستند و تولید نمونه‌های تصادفی از آنها نیز با استفاده از بسته‌های نرم‌افزار R در دسترس اند؛ لذا با استفاده روش نمونه‌گیری گیبز می‌توان استنباط‌های بیزی را انجام داد.

۲.۳ مدل با فرض وابستگی α_i و y_i

در این بخش، مدلی را در نظر می‌گیریم که در آن، مشاهدات اولیه و اثرهای تصادفی وابسته هستند. در این مدل، تابع درست‌نمایی به صورت تابع چگالی توأم از مشاهدات اولیه و آتی خواهد بود. لذا توزیع مشاهدات اولیه به شرط اثرهای تصادفی را به پیروی از [۱۸] به صورت

$$y_i|\alpha_i \sim N(\phi_0 + \phi_1\alpha_i, \sigma^2)$$

در نظر می‌گیریم. حال برای تحلیل بیزی باید توزیع پیشین برای پارامترها و ابرپارامترها مشخص شود. توزیع‌های پیشین پارامترهای β و σ را مانند بخش ۱.۳ در نظر می‌گیریم. توزیع‌های پیشین پارامترهای ϕ_0 ، ϕ_1 و σ^2 را نیز مانند [۱۰] به صورت

$$\phi_1 \sim N(\mu_{\phi_1}, \sigma_{\phi_1}^2), \phi_0 \sim N(\mu_{\phi_0}, \sigma_{\phi_0}^2), \sigma^2 \sim IG\left(\frac{v_0}{\tau}, \frac{\delta_0}{\tau}\right)$$

در نظر می‌گیریم. حال با توجه به مدل رابطه‌های ۸ و ۹ و توزیع‌های پیشین تعریف شده، مدل سلسله‌مراتبی را می‌توان به صورت

به صورت

$$\text{RMSE}(\hat{\beta}_l) = \sqrt{\frac{1}{100} \sum_{r=1}^{100} (\hat{\beta}_l^r - \beta_l)^2}, \quad l = 1, \dots, k,$$

$$\text{Bias}(\hat{\gamma}) = \frac{1}{100} \sum_{r=1}^{100} (\hat{\gamma}^r - \gamma),$$

$$\text{RMSE}(\hat{\gamma}) = \sqrt{\frac{1}{100} \sum_{r=1}^{100} (\hat{\gamma}^r - \gamma)^2}.$$

واضح است هرچه اندازه این معیارها کوچک‌تر باشند، حاکی از کارایی زیاد مدل است. چون نوشتن همه مقادیر با توجه به وجود پنج مدل، دو روش برآورد برای سه نوع توزیع خطا و وجود سه پارامتر در هر شبیه‌سازی و در مجموع ۱۸ پارامتر، در قالب جدول امکان‌پذیر نیست، بنا بر این نتایج در قالب نمودارهای جعبه‌ای برای چندک‌های $\tau = (0/25, 0/5, 0/75)$ ، خلاصه شده و نتایج در شکل‌های ۱ و ۲ و ۳ ارائه شده است. هرچه نمودارهای جعبه‌ای کوچک‌تر و نزدیک به صفر باشند، نشان‌دهنده اریبی کمتر و برآورد بهتر پارامترها است. شکل ۱ عملکرد روش‌ها را از نظر اریبی و RMSE در چندک $\tau = 0/25$ مقایسه می‌کند. با توجه به این شکل در روش‌های بیزی BI و BD اریبی پارامتر بسیار نزدیک به صفر است. از نظر معیار RMSE نیز به غیر از توزیع نرمال استاندارد که روش PQ1 عملکرد نزدیک به روش‌های بیزی دارد، روش‌های بیزی عملکرد بهتری نسبت به سایر روش‌ها دارند. شکل‌های ۲ و ۳ نیز عملکرد روش‌ها را به لحاظ معیارهای ذکر شده در چندک‌های $\tau = 0/5$ و $\tau = 0/75$ مقایسه می‌کنند. با توجه به این دو شکل نیز مشاهده می‌شود که مانند شکل ۱، در این نمودارها نیز روش‌های بیزی برآورد بهتری از پارامترها در همه توزیع‌های فرض شده برای خطا دارند. همچنین توجه کنیم روش QR به دلیل این که اثرهای تصادفی را در مدل لحاظ نمی‌کند، نسبت به روش‌های دیگر که اثرات تصادفی را در مدل لحاظ می‌کنند، عملکرد خوبی ندارد.

۵ تحلیل داده‌های واقعی

در این بخش به مطالعه و تحلیل داده‌های واقعی می‌پردازیم. این داده‌ها شامل رشد اندازه وزن ۲۰۰ کودک از بدو تولد تا دوسالگی است که از سال ۱۳۹۱ تا سال ۱۳۹۲ (دو سال) از درمانگاه‌های شهرستان قم به تصادف انتخاب و گردآوری شده است. هدف بررسی تأثیر برخی عوامل بر میزان رشد وزن کودکان است. این

تولید می‌کنیم که در آن $t = -49, \dots, T$ است. همچنین برای تولید متغیرهای پاسخ فرض می‌کنیم $y_{i,-50} = 0$ و پاسخ‌های y_{it} را برای $i = 1, \dots, n$ و $t = -49, \dots, T$ تولید می‌کنیم. سپس ۵۰ مشاهده اول را کنار می‌گذاریم و مشاهدات مربوط به $t = 0, \dots, T$ را برای تحلیل مورد استفاده قرار می‌دهیم. توجه کنید که این روش تولید داده‌ها برای اطمینان از عدم تأثیر مقدار اولیه بر داده‌های مورد تحلیل است. مقدار واقعی پارامتر β را برابر با (۱، ۲) قرار می‌دهیم و شبیه‌سازی را برای مقادیر مختلف γ به صورت $\gamma = (0, 0/1, 0/3, 0/5, 0/8, 0/9)$ انجام می‌دهیم. توزیع‌های پیشین را برای پارامترها و ابرپارامترها تقریباً ناآگاهی بخش به صورت

$\phi_1 \sim N(0, 100)$ و $\phi_2 \sim IG(0/01, 0/01)$ ، $\sigma \sim IG(0/01, 0/01)$ ، $\beta \sim N(0, 100I)$ ، $\phi_1 \sim N(0, 100)$ و $\sigma^2 \sim IG(0/01, 0/01)$ در نظر می‌گیریم. برای برآورد پارامترها در مدل‌های بیزی با استفاده از روش نمونه‌گیری گیبز، ۸۰۰۰ نمونه از توزیع‌های پسین شرطی کامل پارامترها تولید کرده، ۲۰۰۰ نمونه اول را به عنوان دوره همگرایی مطلوب کنار می‌گذاریم. برای ارزیابی همگرایی زنجیر و تعیین تقریبی تعداد نمونه‌ها به عنوان همگرایی مطلوب از آماره \hat{R} که در [۵] ارائه شده است، استفاده می‌کنیم. روش‌های بیزی، یعنی رگرسیون چندکی با تاوان لاسو سازوار از دیدگاه آمار بیزی با فرض وابستگی اثرات تصادفی و مشاهدات اولیه (BD) و با فرض استقلال اثرهای تصادفی و مشاهدات اولیه (BI) را با روش رگرسیون چندکی معمولی بدون در نظر گرفتن اثرهای تصادفی (QR) و روش رگرسیون چندکی در داده‌های پانلی با در نظر گرفتن تاوان لاسو روی اثرات تصادفی از دیدگاه آمار بسامدی (PQR) که توسط کوئنکر [۱۱] ارائه شده است، مقایسه می‌کنیم. همانند گراسی و باتای [۸] دو نوع تاوان $\lambda = 0/01$ ، $\lambda = 5$ و (PQ1) و (PQ2) را برای روش کوئنکر [۱۱] در نظر می‌گیریم. برای مقایسه کارایی روش‌ها و به منظور لحاظ کردن تغییرپذیری نتایج حاصل از شبیه‌سازی، روند شبیه‌سازی را ۱۰۰ مرتبه تکرار می‌کنیم و دو معیار اریبی^{۱۳} پارامترها و مجذور میانگین توان‌های دوم خطا^{۱۴} را طبق روابط زیر برای پارامتر β و γ در نظر می‌گیریم.

$$\text{Bias}(\hat{\beta}_l) = \frac{1}{100} \sum_{r=1}^{100} (\hat{\beta}_l^r - \beta_l), \quad l = 1, \dots, k,$$

^{۱۳} bias

^{۱۴} root mean square errors

به این جدول برآورد پارامترها در همه روش‌ها تقریباً یکسان است. همان‌طور که انتظار می‌رفت، ضریب متغیر تأخیری برای همه روش‌ها مقدار بزرگ (نزدیک یک) است که حاکی از وابستگی متغیر پاسخ به اندازه آن در زمان قبل است.

بحث و نتیجه‌گیری

در داده‌های پانلی پویا متغیر پاسخ تحت تأثیر مقدار خود در دوره قبل است. در این نوع مطالعات، متغیر تأخیری به‌عنوان متغیر توصیفی وارد مدل می‌شود. در مدل‌های پانلی پویا علاوه بر اثرهای ثابت که تأثیر هر یک از متغیرهای توصیفی را بر متغیر پاسخ بیان می‌کند، اثرهای تصادفی نیز در مدل لحاظ می‌شود. این اثرها تغییرات بین متغیرهای پاسخ (تغییرات بین گروهی) را کنترل می‌کنند. در مطالعاتی که تعداد مشاهدات از متغیر پاسخ زیاد است، اثرهای تصادفی نیز زیاد می‌شود و پارامترهای زیادی وارد مدل شده، بنا بر این دقت مدل کم و تفسیر آن مشکل می‌شود. در این مقاله با ایجاد تاوان لاسو سازوار روی اثرهای تصادفی، اثرهای تصادفی را به سمت صفر منقبض کرده و اثرهای کم‌اهمیت از مدل حذف شدند. تاوان لاسو سازوار با تعریف توزیع پیشین مناسب روی اثرهای تصادفی، در مدل لحاظ شده و از دیدگاه آمار بیزی مورد بررسی قرار گرفت. نتایج حاصل از شبیه‌سازی نشان دادند که روش بیزی ارائه‌شده در تمامی چندک‌ها نسبت به روش‌های متداول دیگر در رگرسیون چندکی برای داده‌های پانلی پویا با اثرات تصادفی عملکرد بهتری دارد. نتایج حاصل از تحلیل داده‌های واقعی نیز مشخص کرد که در تحلیل داده‌های رشد وزن کودک، استفاده از مدل داده‌های پانلی پویا مناسب‌تر است. همچنین دو معیار AIC و BIC حاکی از برازش خوب مدل بیزی ارائه‌شده به داده‌ها، در مقایسه با روش‌های بسامدی است.

عوامل، یعنی متغیرهای توصیفی، عبارت‌اند از جنسیت کودک (β_1)، محل سکونت خانوار (β_2)، شاغل بودن مادر (β_3)، سن مادر (β_4)، تحت مراقبت ویژه بودن مادر (β_5)، وزن مادر (β_6)، قد مادر (β_7)، سن جنینی نوزاد (β_8)، چندقلو بودن نوزاد (β_9)، مدت شیردهی انحصاری مادر (β_{10})، مدت شیردهی مادر در کنار غذای کمکی (β_{11}) و متغیر پاسخ نیز اندازه وزن کودک است که در فواصل زمانی بدو تولد و ۲، ۴، ۶، ۹، ۱۲، ۱۸ و ۲۴ ماهگی اندازه‌گیری شده است. دو مدل برای تحلیل داده‌ها به‌صورت زیر در نظر می‌گیریم.

مدل ۱: مدل داده‌های پانلی بدون در نظر گرفتن متغیر تأخیری؛ یعنی

$$y_{it} = x_{it}^{(1)}\beta_1 + x_{it}^{(2)}\beta_2 + \dots + x_{it}^{(11)}\beta_{11} + \alpha_i + \epsilon_{it},$$

$$i = 1, \dots, 200, \quad t = 1, \dots, 8.$$

مدل ۲: مدل داده‌های پانلی پویا با در نظر گرفتن متغیرهای تأخیری؛ یعنی

$$y_{it} = x_{it}^{(1)}\beta_1 + \dots + x_{it}^{(11)}\beta_{11} + \gamma y_{i(t-1)} + \alpha_i + \epsilon_{it},$$

$$i = 1, \dots, 200, \quad t = 1, \dots, 8.$$

معیارهای اطلاع آکائیکه^{۱۵} (AIC) و بیزی^{۱۶} (BIC) برای روش‌های مختلف عنوان شده و دو مدل ۱ و ۲ محاسبه شده است که نتایج آن برای چندک‌های $\tau = (0/25, 0/5, 0/75)$ در جدول ۱ ارائه‌شده است. نتایج جدول ۱ نشان می‌دهد که مقدار این دو معیار برای مدل دوم در همه روش‌ها کوچک‌تر است، لذا مدل ۲ در همه چندک‌های ذکر شده بهتر از مدل اول به داده‌ها برازنده شده است. همچنین با توجه به این جدول، دو روش بیزی در مدل دوم نیز نسبت به روش‌های بسامدی، بهتر به داده‌ها برازنده شده‌اند. جدول ۲ نیز برآورد ضرایب رگرسیونی را در مدل ۲ برای روش‌های مختلف در چندک ۰/۵ نشان می‌دهد. با توجه

^{۱۵} Akaike Information Criterion

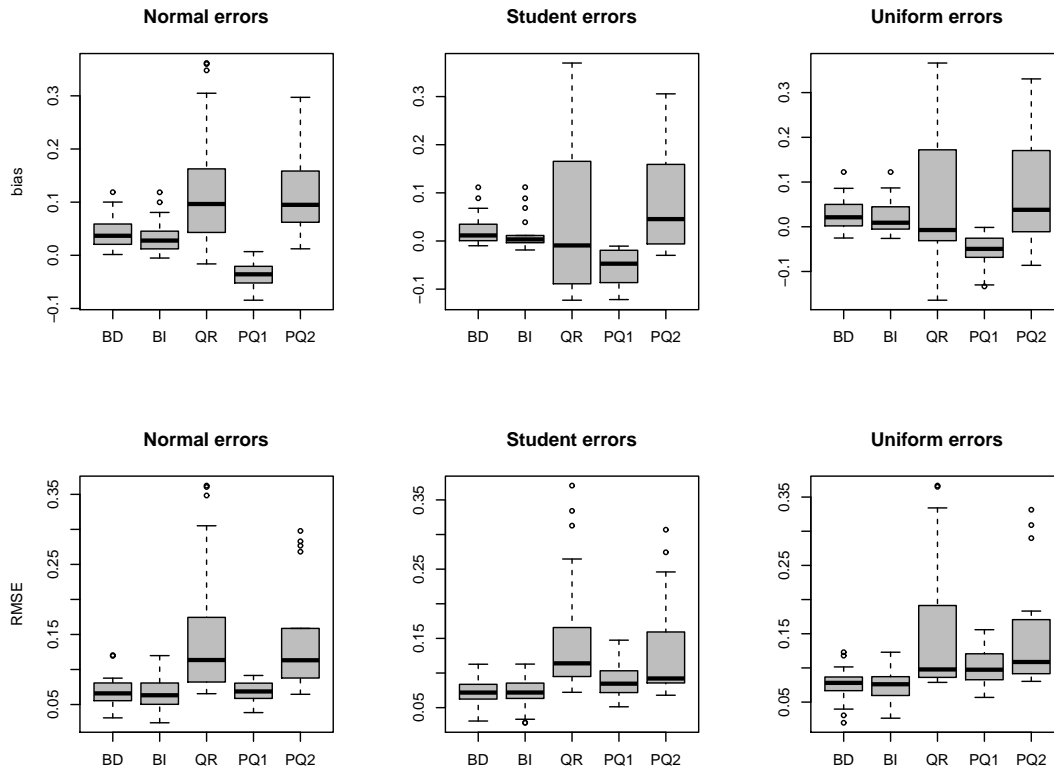
^{۱۶} Bayesian Information Criterion

جدول ۱. برآورد معیارهای AIC و BIC برای دو مدل عنوان شده

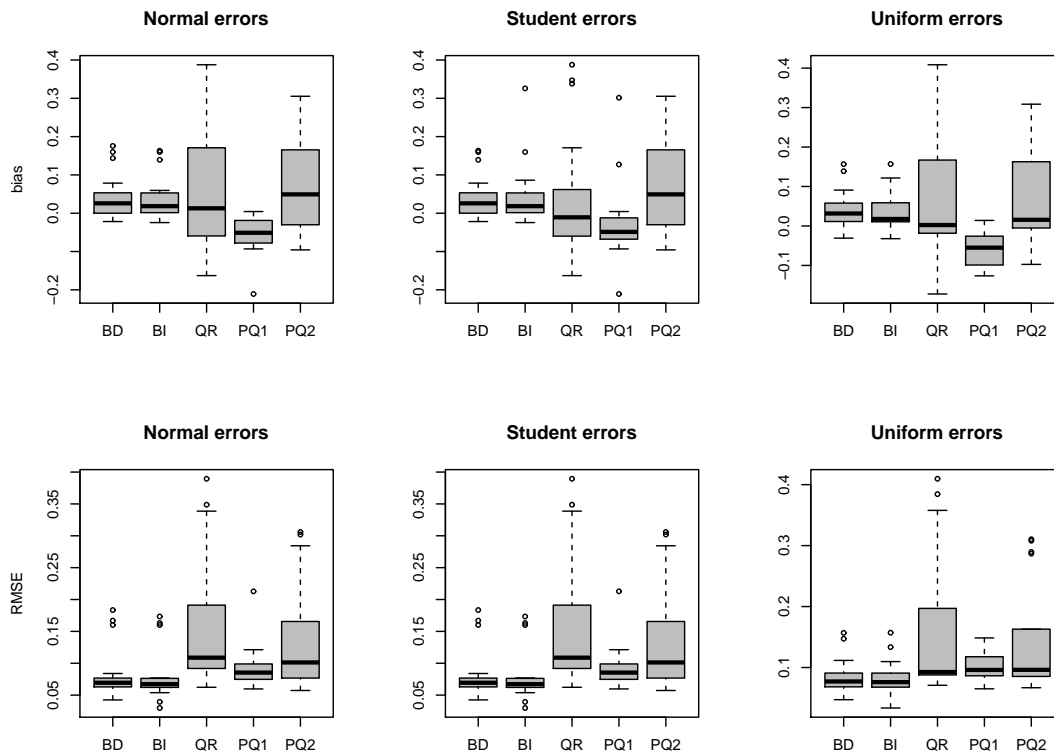
مدل ۲		مدل ۱		روش	چندک
BIC	AIC	BIC	AIC		
۰/۴۰۸۷	۰/۴۶۵۳	۲/۸۳۶۸	۲/۷۸۳۹	BD	۰/۲۵
۰/۵۰۵۲	۰/۴۶۹۲	۲/۸۳۷۴	۲/۷۸۴۳	BI	
۱/۲۰۷۶	۱/۱۶۷۶	۳/۰۶۱۱	۳/۰۰۸۲	PQ1	
۱/۲۰۳۷	۱/۱۷۱۵	۲/۸۸۷۸	۲/۸۳۴۹	PQ2	
۰/۴۱۸۸	۰/۴۴۶۴	۴/۰۹۴۷	۴/۰۲۲۱	BD	۰/۵
۰/۴۸۸۹	۰/۴۵۲۸	۴/۱۰۶۲	۴/۰۳۳۵	BI	
۱/۰۴۵۳	۱/۰۰۹۲	۴/۲۹۷۴	۴/۲۲۴۸	PQ1	
۱/۱۹۶۸	۱/۱۶۰۷	۴/۱۸۱۹	۴/۱۰۹۳	PQ2	
۰/۴۰۷۶	۰/۴۸۹۶	۱/۳۹۶۶	۱/۳۶۸۳	BD	۰/۷۵
۰/۵۳۲۳	۰/۴۹۶۲	۱/۳۹۷۴	۱/۳۶۹۱	BI	
۱/۲۵۱۲	۱/۲۱۵۱	۱/۴۴۵۳	۱/۴۱۷۰	PQ1	
۱/۲۱۳۶	۱/۷۷۵۱	۲/۲۵۵۷	۲/۲۲۷۴	PQ2	

جدول ۲. برآورد ضرایب رگرسیونی برای مدل ۲

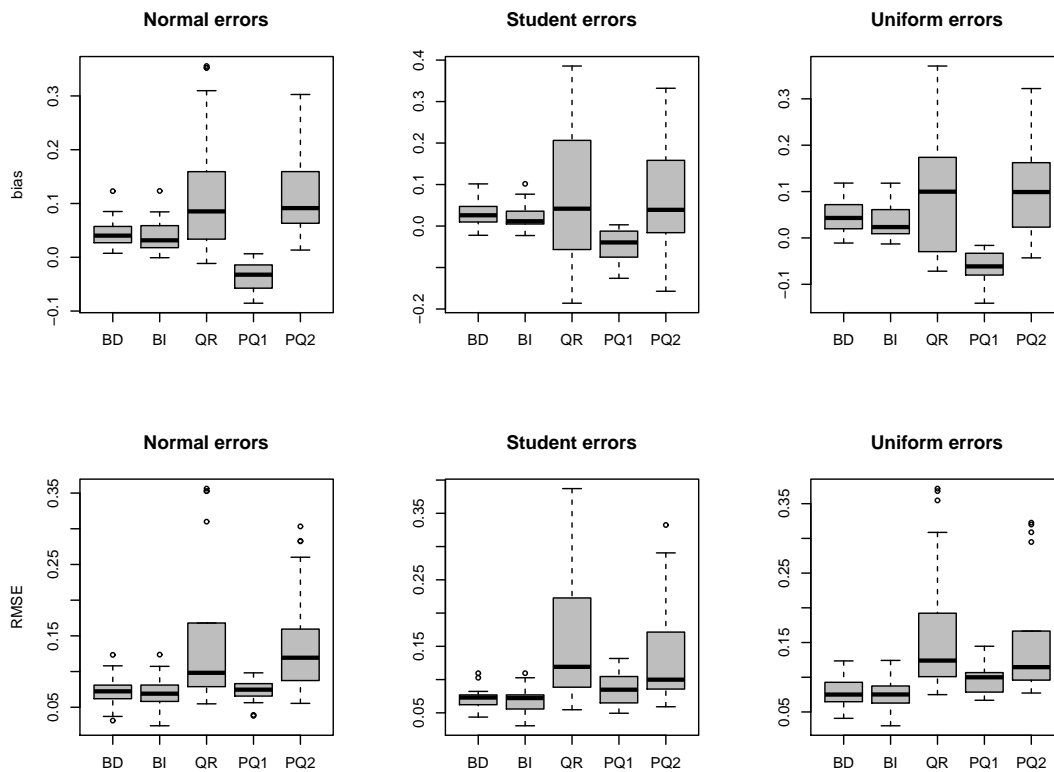
PQ2	PQ1	BI	BD	پارامتر
۰/۱۶۴۸	۰/۰۲۷۲	-۰/۰۰۰۲	-۰/۰۰۲۴	β_1
۰/۰۰۷۵	۰/۰۰۷۸	۰/۰۰۶۷	۰/۰۰۵۳	β_2
۰/۰۰۰۰	۰/۰۰۰۰	-۰/۰۰۳۹	-۰/۰۰۳۷	β_3
-۰/۰۲۳۵	-۰/۰۲۵۰	-۰/۰۲۳۵	-۰/۰۲۴۰	β_4
۰/۰۴۳۱	۰/۰۴۴۱	۰/۰۹۴۹	۰/۰۹۶۵	β_5
۰/۰۰۲۷	۰/۰۰۲۸	۰/۰۰۳۹	۰/۰۰۴۰	β_6
-۰/۰۰۰۲	-۰/۰۰۰۳	-۰/۰۰۲۲	-۰/۰۰۲۱	β_7
۰/۰۱۸۲۳	۰/۰۱۸۳	۰/۰۲۵۹	۰/۰۲۵۵	β_8
۰/۱۷۹۶	۰/۱۷۷۵	۰/۱۲۹۹۱	۰/۱۳۱۳	β_9
-۰/۱۳۸۱	۰/۰۰۰۰	۰/۰۲۹۴	۰/۰۲۹۱	β_{10}
-۰/۰۰۰۲۰	-۰/۰۰۰۲۱	۰/۰۰۰۰۵	۰/۰۰۰۰۱	β_{11}
۰/۹۹۴۲	۰/۹۸۶۲	۰/۹۹۴۴	۰/۹۹۴۴	γ



شکل ۱. عملکرد روش‌ها به لحاظ معیار اریبی و RMSE در چندک $\tau = 0.25$.



شکل ۲. عملکرد روش‌ها به لحاظ معیار اریبی و RMSE در چندک $\tau = 0.5$.



شکل ۳. عملکرد روش‌ها به لحاظ معیار اریبی و RMSE در چندک $\tau = 0.75$.

پیوست

الف) توزیع‌های پسین شرطی کامل پارامترها و ابرپارامترها با فرض استقلال α_i و y_i .

$$\begin{aligned} \pi(\alpha_i|\cdot) &\propto f(\mathbf{y}_i|\boldsymbol{\beta}, \alpha_i, \mathbf{e}_i, \gamma, \sigma, s_i, \nu^\lambda, \phi, \cdot) \times \pi(\alpha_i|s_i) \implies \\ \alpha_i|\cdot &\sim N(\bar{\mu}_i, \tilde{\sigma}_i), \quad i = 1, \dots, n \\ \bar{\mu}_i &= \tilde{\sigma}_i \left[\frac{1}{\gamma\sigma} \sum_{t=1}^T \frac{(y_{it} - \mathbf{x}'_{it}\boldsymbol{\beta} - \gamma y_{i(t-1)} - k_1 e_{it})}{e_{it}} \right] \\ \tilde{\sigma}_i &= \left(\frac{1}{s_i} + \frac{1}{\gamma\sigma} \sum_{t=1}^T \frac{1}{e_{it}} \right)^{-1}. \end{aligned}$$

$$\begin{aligned} \pi(\boldsymbol{\beta}|\cdot) &\propto f(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{e}, \sigma, \mathbf{s}, \nu^\lambda, \phi, \cdot) \pi(\boldsymbol{\beta}) \implies \boldsymbol{\beta}|\cdot \sim N_k(\mathbf{B}\mathbf{b}, \mathbf{B}) \\ \mathbf{b} &= \frac{1}{\gamma\sigma} \sum_{i=1}^n \sum_{t=1}^T \frac{\mathbf{x}_{it}(y_{it} - \alpha_i - \gamma y_{i(t-1)} - k_1 e_{it})}{e_{it}} + \mathbf{B}_0^{-1} \mathbf{b}_0 \\ \mathbf{B}^{-1} &= \frac{1}{\gamma\sigma} \sum_{i=1}^n \sum_{t=1}^T \frac{\mathbf{x}_{it}\mathbf{x}'_{it}}{e_{it}} + \mathbf{B}_0^{-1}. \end{aligned}$$

$$\begin{aligned} \pi(\sigma|\cdot) &\propto f(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{e}, \sigma, \mathbf{s}, \nu^\lambda, \phi, \cdot) \pi(\sigma) \implies \sigma|\cdot \geq \sim IG(\nu_1, w_1) \\ w_1 &= \sum_{i=1}^n \sum_{t=1}^T \left\{ \frac{(y_{it} - \mathbf{x}'_{it}\boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i - k_1 e_{it})^2}{\gamma e_{it}} + \tau(1 - \tau)e_{it} \right\} + d. \\ \nu_1 &= \frac{\gamma N}{\gamma} + c. \end{aligned}$$

$$\pi(e_{it}|\cdot) \propto f(y_{it}|\boldsymbol{\beta}, \alpha_i, e_{it}, \sigma, s_i, \nu^\lambda, \phi, \cdot) f(e_{it}),$$

$$\begin{aligned} & \propto (e_{it})^{-1} \exp \left\{ -\frac{1}{\nu} (c_{it}^{\nu} e_{it}^{-1} + d_{it}^{\nu} e_{it}) \right\} \Rightarrow \\ e_{it} | \cdot & \sim GIG(\cdot, c_{it}, d_{it}), \quad i = 1, \dots, n, \quad t = 1, \dots, T, \\ c_{it}^{\nu} & = \frac{(y_{it} - \mathbf{x}'_{it} \boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i)^{\nu}}{\nu \sigma}, \quad d_{it}^{\nu} = \frac{(1 - \nu \tau)^{\nu}}{\nu \sigma} + \frac{\nu \tau (1 - \tau)}{\sigma}, \end{aligned}$$

$GIG(r, c, d)$ نشان‌دهنده توزیع گاوسی وارون تعمیم‌یافته^{۱۷} با تابع چگالی

$$f(x|r, c, d) \propto x^{r-1} \exp\{-\frac{1}{\nu}(c^{\nu} x^{-1} + d^{\nu} x)\}, x > 0, -\infty < r < +\infty, \quad m, n > 0.$$

است. برای تولید نمونه تصادفی از این توزیع می‌توان از تابع $rgig()$ استفاده کرد که در بسته $ghyp$ در نرم‌افزار R موجود است [۳].

$$\begin{aligned} \pi(\gamma | \cdot) & \propto f(\mathbf{y} | \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{e}, \sigma, \mathbf{s}, \nu^{\nu}, \phi, \gamma) \pi(\gamma | \psi) \\ & \propto (e_{it})^{-1} \exp \left\{ -\sum_{i=1}^n \sum_{t=1}^T \frac{1}{\nu \sigma e_{it}} (y_{it} - \mathbf{x}'_{it} \boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i - k_{\nu} e_{it})^{\nu} \right\} \\ & \quad \times \exp \left\{ -\frac{\gamma^{\nu}}{\nu \psi^{\nu}} \right\} \Rightarrow (\gamma | \cdot) \sim N(\bar{\mu}_{\gamma}, \tilde{\sigma}_{\gamma}), \\ \bar{\mu}_{\gamma} & = \tilde{\sigma}_{\gamma} \left[\frac{1}{\nu \sigma} \sum_{i=1}^n \sum_{t=1}^T \frac{y_{i(t-1)} (y_{it} - \mathbf{x}'_{it} \boldsymbol{\beta} - \alpha_i - k_{\nu} e_{it})}{e_{it}} \right], \\ \tilde{\sigma}_{\gamma} & = \left[\frac{1}{\nu \sigma} \sum_{i=1}^n \sum_{t=1}^T \frac{y_{i(t-1)}^{\nu}}{e_{it}} + \frac{1}{\psi^{\nu}} \right]^{-1}. \end{aligned}$$

$$\begin{aligned} \psi^{\nu} | \cdot & \sim IG(a + \frac{1}{\nu}, b + \frac{\gamma^{\nu}}{\nu}), \quad s_i | \cdot \sim GIG(\frac{1}{\nu}, \alpha_i, \nu_i), \quad i = 1, \dots, n, \\ \phi | \cdot & \sim G(n, \sum_{i=1}^n \nu_i^{\nu}), \quad \nu_i^{\nu} | \cdot \sim G(\nu, \frac{s_i}{\nu} + \phi), \end{aligned}$$

(ب): توزیع‌های پسین شرطی کامل پارامترها و ابرپارامترها با فرض وابستگی α_i و y_i .

$$\begin{aligned} \pi(\alpha_i | \cdot) & \propto f(\mathbf{y}_i | \boldsymbol{\beta}, \alpha_i, \mathbf{e}_i, \gamma, \sigma, s_i, \nu^{\nu}, \phi) f(y_i | \alpha_i, \phi, \cdot, \sigma^{\nu}) \pi(\alpha_i | s_i) \\ & \propto \exp \left\{ -\frac{1}{\nu \sigma} \sum_{t=1}^T \frac{(y_{it} - \mathbf{x}'_{it} \boldsymbol{\beta} - \gamma y_{i(t-1)} - \alpha_i - k_{\nu} e_{it})^{\nu}}{e_{it}} - \frac{\alpha_i^{\nu}}{\nu s_i} \right\} \\ & \quad \times \exp \left\{ -\frac{1}{\sigma^{\nu}} (y_i - \phi - \phi_{\nu} \alpha_i)^{\nu} \right\} \propto \exp \left\{ -\frac{\tilde{\sigma}_i^{\nu}}{\nu} (\alpha_i - \bar{\mu}_i)^{\nu} \right\} \Rightarrow \\ \alpha_i | \cdot & \sim N(\bar{\mu}_i, \tilde{\sigma}_i), \quad i = 1, \dots, n \\ \bar{\mu}_i & = \tilde{\sigma}_i \left[\frac{1}{\nu \sigma} \sum_{t=1}^T \frac{(y_{it} - \mathbf{x}'_{it} \boldsymbol{\beta} - \gamma y_{i(t-1)} - k_{\nu} e_{it})}{e_{it}} + \frac{(y_i - \phi) \phi_{\nu}}{\sigma^{\nu}} \right], \\ \tilde{\sigma}_i & = \left(\frac{1}{s_i} + \frac{1}{\nu \sigma} \sum_{t=1}^T \frac{1}{e_{it}} + \frac{\phi_{\nu}^{\nu}}{\sigma^{\nu}} \right)^{-1}. \end{aligned}$$

$$\begin{aligned} \pi(\phi_{\nu} | \cdot) & \propto \exp \left\{ -\frac{1}{\nu \sigma^{\nu}} \sum_{i=1}^n (y_i - \phi - \phi_{\nu} \alpha_i)^{\nu} \right\} \times \exp \left\{ -\frac{1}{\sigma_{\phi_{\nu}}^{\nu}} (\phi_{\nu} - \mu_{\phi_{\nu}})^{\nu} \right\} \Rightarrow \\ \phi_{\nu} | \cdot & \sim N(\bar{\mu}_{\phi_{\nu}}, \tilde{\sigma}_{\phi_{\nu}}) \\ \bar{\mu}_{\phi_{\nu}} & = \tilde{\sigma}_{\phi_{\nu}} \left[\frac{\sum_{i=1}^n \alpha_i (y_i - \phi)}{\sigma^{\nu}} + \frac{\mu_{\phi_{\nu}}}{\sigma_{\phi_{\nu}}^{\nu}} \right], \quad \tilde{\sigma}_{\phi_{\nu}} = \left[\frac{\sum_{i=1}^n \alpha_i^{\nu}}{\sigma^{\nu}} + \frac{1}{\sigma_{\phi_{\nu}}^{\nu}} \right]^{-1}. \end{aligned}$$

$$\begin{aligned} \pi(\phi | \cdot) & \propto \exp \left\{ -\frac{1}{\nu \sigma^{\nu}} \sum_{i=1}^n (y_i - \phi - \phi_{\nu} \alpha_i)^{\nu} \right\} \times \exp \left\{ -\frac{1}{\sigma_{\phi}^{\nu}} (\phi - \mu_{\phi})^{\nu} \right\} \Rightarrow \\ \phi | \cdot & \sim N(\bar{\mu}_{\phi}, \tilde{\sigma}_{\phi}), \end{aligned}$$

^{۱۷} Generalized Inverse Gaussian

$$\bar{\mu}_{\phi} = \bar{\sigma}_{\phi} \left[\frac{\sum_{i=1}^n (y_i - \phi \alpha_i)}{\sigma_{\phi}^2} + \frac{\mu_{\phi}}{\sigma_{\phi}^2} \right], \quad \bar{\sigma}_{\phi} = \left[\frac{n}{\sigma_{\phi}^2} + \frac{1}{\sigma_{\phi}^2} \right]^{-1}.$$

$$\pi(\sigma_{\phi} | \cdot) \propto \left(\frac{1}{\sigma_{\phi}^2} \right)^{\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma_{\phi}^2} \sum_{i=1}^n (y_i - \phi - \phi \alpha_i)^2 \right\} \times \left(\frac{1}{\sigma_{\phi}^2} \right)^{\left(\frac{v}{2} + 1 \right)} \exp \left\{ -\frac{\delta}{2\sigma_{\phi}^2} \right\} \Rightarrow$$

$$\sigma_{\phi}^2 | \cdot \sim IG \left(\frac{n + v}{2}, \frac{\sum_{i=1}^n (y_i - \phi - \phi \alpha_i)^2}{2} + \frac{\delta}{2} \right)$$

مراجع

- [1] Alhamzawi, R., Yu, K. and Benoit, D. F. (2012). Bayesian adaptive Lasso quantile regression, *Statistical Modeling*, **12**, 279-297.
- [2] Alhamzawi, R., Benoit, D. F. and Yu, K. (2013). Binary quantile regression with adaptive lasso: A Bayesian approach, working paper.
- [3] Breyman, W. and Luthi, D. (2014). ghyp: A package on the generalized hyperbolic distribution and its special cases, R Package, Version 1.5.6, URL: <http://www.r-project.org>
- [4] Crouchley, R. and Davies, R. B. (2001). A Comparison of GEE and random effects models for distinguishing heterogeneity nonstationarity and state dependent in collection of short binary event series, *Statistical Modelling*, **1**, 271-285.
- [5] Gelman, A. Carlin, J. B. Stern, H. S. and Rubin, D. B. (1995). *Bayesian Data Analysis*, Chapman and Hall, London.
- [6] Galvao, A. F. (2011). Quantile regression for dynamic panel data with fixed effects, *Journal of Applied Econometric*, **164**(1), 142-157.
- [7] Galvao, A. F. and Montes-Rojas, G. V. (2010). Penalized quantile regression for dynamic panel data, *Journal of Statistical Planning and Inference*, **140**, 3476-3497.
- [8] Geraci, M. and Bottai, M. (2007). Quantile regression for longitudinal data using the asymmetric Laplace distribution, *Biostatistics*, **8**, 140-154.
- [9] Hsiao, C. (2003). *Analysis of Panel Data*, 2nd ed. Cambridge University Press, Cambridge.
- [10] Kobayashi, G., and Kozumi, H. (2012). Bayesian analysis of quantile regression for censored dynamic panel data, *Journal of Statistical Computation and Simulation*, **27**, 359-380.
- [11] Koenker, R. (2004). Quantile regression for longitudinal data, *Journal of Multivariate Analysis*, **91**, 74-89.
- [12] Koenker, R. and Bassett, G. (1978). Quantile regression, *Econometrica*, **46**, 33-50.
- [13] Kozumi, H. and Kobayashi, G. (2011). Gibbs sampling methods for Bayesian quantile regression, *Journal of Statistical Computation and Simulation*, **81**, 1565-1578.

- [14] Leng, C., Tran, M. N. and Nott, D. A. (2014). Bayesian adaptive lasso, *Annals of the Institute of Statistical Mathematics*, **66**, 221-224.
- [15] Luo, Y., Lian, H. and Tian, M. (2012), Bayesian quantile regression for longitudinal data model, *Journal of Statistical Computation and Simulation*, **82**, 1635-1649.
- [16] Park, T. and Casella, T. (2008), The bayesian lasso. *Journal of the American Statistical Assosation*, **103**, 681-686.
- [17] Tibshirani, R. (1996), Regression shrinkage and selection via the Lasso, *Journal of the Royal Statistical Society, Series B*, **58**, 267 – 288.
- [18] Wooldridge, J. M. (2005). Simple solutions to the initial conditions problem in dynamic nonlinear panel data models with unobserved heterogeneity, *Journal of Applied Econometrics*, **20**, 30-54
- [19] 13. Yu, K. and Moyeed, R.A. (2001). Bayesian quantile regression, *Statistics and Probability Letters*, **54**, 437-447.
- [20] Zou, H. (2006). The adaptive Lasso and its oracle properties. *Journal of the American Statistical Association*, **101**, 1418-1429.