

انتخاب متغیر در مدل‌های خطی-جزئی با اثرات آمیخته برای داده‌های طولی با بعد بالا

مژگان تعاونی و محمد آرشی

گروه آمار، دانشکده علوم ریاضی، دانشگاه صنعتی شاهرود

تاریخ دریافت: ۱۳۹۸/۰۶/۲۹ تاریخ آخرین بازنگری: ۱۳۹۹/۰۱/۲۳

چکیده: در این مقاله، مسئله برآورد و انتخاب متغیر همزمان در مدل‌های خطی-جزئی با اثرات آمیخته برای داده‌های طولی با بعد بالا در نظر گرفته شده است. مولفه ناپارامتری موجود در مدل با اسپلین‌های رگرسیونی تقریب زده شده و سپس از طریق بهینه‌سازی تابع هدف مبتنی بر تابع تاوان، برآورد و انتخاب متغیر به طور همزمان انجام می‌شود. در ادامه، رفتار حدی برآوردگرهای حاصل در چارچوب داده‌های طولی با بعد بالا که در آن تعداد پارامترها متناسب با افزایش حجم نمونه افزایش می‌یابد، مورد مطالعه قرار می‌گیرد. به منظور پیاده‌سازی روش برآورد پیشنهادی، یک الگوریتم تکراری مناسب برای انتخاب متغیرهای مهم و برآورد ضرایب غیر صفر ارائه گردیده است. در نهایت، عملکرد روش پیشنهادی با مطالعه شبیه‌سازی و تحلیل یک مجموعه داده واقعی مورد ارزیابی قرار گرفته است.

واژه‌های کلیدی: داده طولی، انتخاب متغیر، برآوردگر تاوانی، مدل خطی-جزئی، اسپلین هموارساز، بعد بالا، بیماری ایدز.

۱ مقدمه

برای در نظر گرفتن همبستگی بین پاسخ‌ها در مطالعات طولی روش‌های متفاوتی گسترش یافته است. یکی از متداول‌ترین روش‌ها مدل اثرات آمیخته (لرد و ویر، ۱۹۸۲) است که ناهمگنی بین پاسخ‌ها توسط مدل‌بندی مولفه تصادفی بررسی می‌گردد. این مدل با وجود سادگی در محاسبات و تفسیر، به دلیل فرضیات پارامتری محدود کننده از جمله اعمال رابطه خطی بین متغیرهای پاسخ و تبیینی، در مدل‌سازی بسیاری از روابط پیچیده با شکست مواجه می‌شود و مثال‌های بسیاری را می‌توان یافت که شناسایی و تشخیص اثرات رگرسیون بر اساس فرضیات بیان شده امکان‌پذیر نیست. در چنین شرایطی، آماردانان در راستای تحلیل هر چه بهتر داده‌های طولی با خلق ترکیبی از مدل‌های موجود، کلاس دیگری از مدل‌ها را ایجاد نموده‌اند که از جمله متداول‌ترین آن‌ها می‌توان به مدل خطی-جزئی^۱ با اثرات آمیخته اشاره نمود. این مدل تعمیم مدل زیگر و دیگل (۱۹۹۴) است که برای تحلیل داده‌های خود از یک مدل نیمه پارامتری با عرض از مبدأ تصادفی استفاده کرده‌اند.

امروزه توسعه و گسترش سریع علوم در زمینه‌های مختلف و افزایش سرعت ثبت و تحلیل اطلاعات، آماردانان را با حجم عظیمی از اطلاعات پیچیده مواجه نموده است که بکارگیری روش‌های موجود را ناممکن ساخته است. رخداد مساله بعد بالا^۲ در داده‌های طولی نیز اجتناب ناپذیر است. به منظور سهولت در امر تجزیه و تحلیل، اغلب منطقی و مفید است که فرض شود تنها تعداد کمی از متغیرها برای مدل‌سازی متغیر پاسخ مناسب هستند. این فرض، رویکرد نوینی تحت عنوان رگرسیون تاوانیده^۳ را گسترش داد. پس از آن که فو (۲۰۰۳) و فن و لی (۲۰۰۴) به تعمیم رگرسیون تاوانیده به ترتیب در مدل‌های خطی تعمیم یافته و مدل‌های خطی-جزئی با داده‌های طولی پرداختند، مساله انتخاب متغیر در زمینه داده‌های طولی با پیشرفت چشمگیری روبرو شد. در این میان می‌توان به تحقیقات کانتونی و همکاران (۲۰۰۵)، وانگ و کوای (۲۰۰۹) و سوای و همکاران (۲۰۱۰) در مدل‌های خطی تعمیم یافته، باندل و همکاران (۲۰۱۰) در مدل اثرات آمیخته، نی و همکاران (۲۰۱۰) و ما و همکاران (۲۰۱۳) در مدل‌های خطی-جزئی اشاره نمود. با این وجود، مساله انتخاب متغیر برای داده‌های طولی با بعد بالا به واسطه چالش‌های تحمیل شده از جمله وجود همبستگی درون گروهی پیشینه‌چندانی نداشته و تنها می‌توان فعالیت‌های سو و همکاران (۲۰۱۳) و وانگ و همکاران (۲۰۱۲) را نام برد. اخیراً کاظمی و همکاران (۱۳۹۷)، در مدل‌های خطی-جزئی با بعد بالا مساله انتخاب متغیر را بررسی کرده است. در عین حال، تنها تحقیق موجود درباره انتخاب

¹Partially linear model²High dimension³Penalized regression

متغیر در داده‌های طولی با بعد بالا، **تعاونی و آرشی** (۲۰۱۹) است، که تفاوت عمده مقاله حاضر با آن در فرضیات پایه‌ای و در نتیجه روش متفاوت برآورد پارامترها در توزیع نرمال بوده که بسیار ساده‌تر از مرجع اخیر است. در این مقاله برآورد و انتخاب متغیر در مدل خطی-جزئی با اثرات آمیخته مبتنی بر رهیافت توابع تاوان، در نظر گرفته شد. روش پیشنهادی به علت وجود توام اثرات تصادفی و مولفه ناپارامتری در مدل، نسبتاً با تحقیقات انجام شده متفاوت است. بخش ۲ مقاله به معرفی نمادگذاری متداول مدل می‌پردازد. در مرحله اول ابتدا مولفه ناپارامتری موجود در مدل با اسپلاین‌های رگرسیونی تقریب زده شده و سپس در بخش ۳ از طریق بهینه‌سازی تابع هدف مبتنی بر تابع تاوان مانند ماکسیمم درستنمایی تاوانیده^۱ برآورد و انتخاب متغیر انجام می‌گیرد. از آنجایی که برآوردگر معرفی شده در تعامل با ماتریس کوواریانس پاسخ‌ها و اثرات تصادفی هستند، به منظور بهبود کارایی برآوردگر ماکسیمم درستنمایی تاوانیده، از الگوریتم تکراری امیدگیری ماکسیمم‌سازی^۲ (EM) استفاده می‌شود. علاوه بر این، نحوه انتخاب پارامترهای تابع تاوان در مورد بررسی قرار می‌گیرد. بخش ۴ به مطالعه خواص مجانبی برآوردگرها در چارچوب داده‌های بعد بالا که در آن تعداد پارامترها متناسب با افزایش حجم نمونه افزایش می‌یابد، می‌پردازد. برای بررسی عملکرد روش معرفی شده، مطالعه‌ای شبیه‌سازی و همچنین تحلیل داده‌های واقعی به ترتیب در بخش‌های ۵ و ۶ فراهم آمده است. بحث و نتیجه‌گیری بخش پایانی مقاله است.

۲ مدل خطی-جزئی با اثرات آمیخته

فرض کنید نمونه‌ای تصادفی از n آزمودنی وجود داشته باشد، به طوری که نمونه i ام دارای n_i مشاهده مکرر در طول زمان است. لذا تعداد کل مشاهدات برابر با $N = \sum_{i=1}^n n_i$ است. همچنین، فرض کنید $(i = 1, \dots, n; j = 1, \dots, n_i)$ ، پاسخ آزمودنی i ام در زمان t_{ij} باشد. مدل خطی-جزئی با اثرات آمیخته به صورت

$$y_i(t_{ij}) = \mathbf{x}_i^T(t_{ij})\boldsymbol{\beta} + g(t_{ij}) + \mathbf{z}_i^T(t_{ij})\mathbf{b}_i + \varepsilon_i(t_{ij}),$$

است، که در آن $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{p_n})^T$ برداری $1 \times p_n$ بعدی از ضرایب ثابت رگرسیونی وابسته به متغیرهای تبیینی $\mathbf{x}_i(t_{ij})$ ؛ $g(\cdot)$ تابعی نامعلوم، هموار و دارای مشتق دوم؛ $\mathbf{b}_i = (b_{i_1}, \dots, b_{i_q})^T$

¹Penalized maximum likelihood

²Expectation maximization (EM)

برداری $q \times 1$ بعدی از ضرایب تصادفی، مستقل از هم، دارای توزیع نرمال با میانگین \circ و ماتریس کوواریانس D_i و وابسته به متغیرهای تبیینی $z_i(t_{ij})$ ؛ $\varepsilon_i(\cdot)$ خطای تصادفی نرمال و مستقل از b_i ها با $E\{\varepsilon_i(\cdot)\} = \circ$ هستند. بدون کاستن از کلیت مساله می‌توان فرض کرد که t_{ij} ها در بازه $[0, 1]$ هستند.

مشابه فن و همکاران (۲۰۰۷)، فرض می‌شود $\sigma^2(t_{ij}) = \text{Var}\{\varepsilon_i(t_{ij})\}$ ، یعنی واریانس خطاهای تصادفی تابعی ناپارامتری و هموار نسبت به زمان هستند. از طرف دیگر فرض می‌شود تابع همبستگی بین $\varepsilon_i(t_{ij})$ و $\varepsilon_i(t_{ik})$ دارای شکل تابعی $\text{Corr}\{\varepsilon_i(t_{ij}), \varepsilon_i(t_{ik})\} = \rho(t_{ij}, t_{ik}, \theta)$ است، به طوری که $\rho(t_{ij}, t_{ik}, \theta)$ تابعی معین مثبت و θ بردار پارامترهای همبستگی است.

فرض کنید $\mathbf{g}(t_i) = \mathbf{X}_i = (x_i^\top(t_{i1}), \dots, x_i^\top(t_{in_i}))^\top$ ، $\mathbf{y}_i = (y_i(t_{i1}), \dots, y_i(t_{in_i}))^\top$ ، $\boldsymbol{\varepsilon}_i = (\varepsilon_i(t_{i1}), \dots, \varepsilon_i(t_{in_i}))^\top$ ، $\mathbf{Z}_i = (z_i^\top(t_{i1}), \dots, z_i^\top(t_{in_i}))^\top$ ، $(g(t_{i1}), \dots, g(t_{in_i}))^\top$ که در آن بردار پاسخ n_i بعدی، \mathbf{X}_i ماتریس $n_i \times p$ بعدی از اثرات ثابت، \mathbf{Z}_i ماتریس $n_i \times q$ بعدی از اثرات تصادفی، $\mathbf{g}(t_i)$ بردار n_i بعدی از توابع نامعلوم و $\boldsymbol{\varepsilon}_i$ بردار n_i بعدی از خطاهای مدل است. مدل ماتریسی که عناصر آن براساس آزمودنی i ام مشخص شده‌اند به صورت

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{g}(t_i) + \mathbf{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad (1)$$

نمایش داده می‌شود. فرض کنید $\mathbf{V}_i = \mathbf{Z}_i^\top \mathbf{D}_i \mathbf{Z}_i + \boldsymbol{\Sigma}_i$ ماتریس کوواریانس $n_i \times n_i$ بعدی مربوط به پاسخ \mathbf{y}_i باشد که در آن \mathbf{D}_i ماتریس کوواریانس $q \times q$ بعدی مربوط به \mathbf{b}_i ها و $\boldsymbol{\Sigma}_i$ ماتریس کوواریانس $n_i \times n_i$ بعدی مربوط به $\boldsymbol{\varepsilon}_i$ ها است. به تبعیت از لیانگ و زیگر (۱۹۸۶)، $\boldsymbol{\Sigma}_i$ به صورت $\boldsymbol{\Sigma}_i = \mathbf{A}_i \mathbf{R}_i(\boldsymbol{\theta}) \mathbf{A}_i^\top$ در نظر گرفته می‌شود، به طوری که $\mathbf{A}_i = \text{diag}(\sigma(t_{i1}), \dots, \sigma(t_{in_i}))$ یک ماتریس قطری $n_i \times n_i$ بعدی با عناصر قطری $\sigma(t_{ij})$ و $\mathbf{R}_i(\boldsymbol{\theta})$ ماتریس همبستگی بین $\varepsilon_i(t_{ij})$ و $\varepsilon_i(t_{ik})$ است که عضو (j, k) ام آن به صورت $\rho(t_{ij}, t_{ik}, \boldsymbol{\theta})$ تعریف می‌شود.

مدل (۱) را می‌توان به صورت مدل ماتریسی کلی $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{g}(t) + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}$ بازنویسی کرد، به طوری که $\mathbf{Z} = \text{diag}(\mathbf{Z}_1, \dots, \mathbf{Z}_n)^\top$ ، $\mathbf{X} = (\mathbf{X}_1^\top, \dots, \mathbf{X}_n^\top)^\top$ ، $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_n^\top)^\top$ ، $\mathbf{b} = (\mathbf{b}_1^\top, \dots, \mathbf{b}_n^\top)^\top$ و $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1^\top, \dots, \boldsymbol{\varepsilon}_n^\top)^\top$ $\mathbf{V} = \text{diag}(\mathbf{V}_1, \dots, \mathbf{V}_n)$ همچنین $\mathbf{D} = \text{diag}(\mathbf{D}_1, \dots, \mathbf{D}_n)$ ماتریس کوواریانس $N \times N$ بعدی مربوط به \mathbf{y} بوده که در آن $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_n)$ ماتریس کوواریانس $nq \times nq$ بعدی مربوط به بردار اثر تصادفی \mathbf{b} و $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_n)$ ماتریس کوواریانس $N \times N$ بعدی مربوط به $\boldsymbol{\varepsilon}$ است. سپس $\mathbf{g}(t)$ با اسپلاین‌های رگرسیونی به صورت

$g(t) = B(t)^T \alpha$ تقریب می‌شود، که در آن $\alpha = (\alpha_1, \dots, \alpha_{h_n})^T$ بردار ضرایب رگرسیونی،

$$B(t) = (1, t, \dots, t^d, (t - t_1)_+^d, \dots, (t - t_{L_n})_+^d)^T,$$

بیان نمود، که در آن $B_i(t_i) = (B^T(t_{i1}), \dots, B^T(t_{in_i}))^T$ یک ماتریس $n_i \times h_n$ بعدی است و اثرات ثابت و اثرات اسپیلین و $\theta = (\beta^T, \alpha^T)^T$ بردار پارامتر ادغام شده $1 \times (p_n + h_n)$ بعدی هستند.

$$\begin{aligned} y_i &= X_i \beta + B_i(t_i) \alpha + Z_i b_i + \varepsilon_i \\ &= \tilde{X}_i \theta + Z_i b_i + \varepsilon_i, \end{aligned} \quad (2)$$

فرض کنید $\Theta = (\theta, D, \Sigma)$ مجموعه تمام پارامترهای مدل است. اگر اثرات تصادفی b_i را به عنوان مقادیر گمشده و مجموعه $\{(y_i, b_i), i = 1, \dots, n\}$ داده‌های کامل در نظر گرفته شود، آنگاه لگاریتم تابع درست‌نمایی بدون در نظر گرفتن جمله ثابت به صورت

۳ الگوریتم EM برای ماکسیمم درست‌نمایی تاوانیده

برای ارزیابی امید شرطی (۳) به شرط داده‌های مشاهده شده $\{y_i, i = 1, \dots, n\}$ و برآوردهای جاری $(\hat{\theta}^{(r)}, \hat{D}^{(r)}, \hat{\Sigma}^{(r)})$ ، از لم زیر استفاده می‌شود که نتیجه‌ای از قضیه ۳ در وانگ و فن (۲۰۱۱) است.

$$\ell_c(\Theta) = \sum_{i=1}^n \frac{1}{p} \{ \log |\Sigma_i^{-1}| + \log |D_i^{-1}| - [\varepsilon_i^T \Sigma_i^{-1} \varepsilon_i + b_i^T D_i^{-1} b_i] \}, \quad (3)$$

برآورددهای جاری $(\hat{\theta}^{(r)}, \hat{D}^{(r)}, \hat{\Sigma}^{(r)})$ ، از لم زیر استفاده می‌شود که نتیجه‌ای از قضیه ۳ در وانگ و فن (۲۰۱۱) است.

لم ۱. تحت مفروضات این بخش و مدل (۱)، داریم

$$\begin{bmatrix} \mathbf{y}_i \\ \mathbf{b}_i \end{bmatrix} \sim N_{n_i+q} \left(\begin{bmatrix} \tilde{\mathbf{X}}_i \boldsymbol{\theta} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{V}_i & \mathbf{Z}_i \mathbf{D}_i \\ \mathbf{D}_i \mathbf{Z}_i^\top & \mathbf{D}_i \end{bmatrix} \right),$$

$$\mathbf{b}_i | (\mathbf{y}_i) \sim N_q (\mathbf{D}_i \mathbf{Z}_i^\top \mathbf{V}_i^{-1} (\mathbf{y}_i - \tilde{\mathbf{X}}_i \boldsymbol{\theta}), (\mathbf{D}_i^{-1} + \mathbf{Z}_i^\top \boldsymbol{\Sigma}_i^{-1} \mathbf{Z}_i)^{-1}).$$

مراحل اصلی الگوریتم EM به شرح زیر انجام می‌شود:

گام E: محاسبه تابع $Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)}) = E(\ell_c(\boldsymbol{\Theta}) | \mathbf{y}, \hat{\boldsymbol{\Theta}}^{(r)})$ که برای $i = 1, \dots, n$ مجموع عبارت $Q_i(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)}) = E(\ell_c^{[i]}(\boldsymbol{\Theta}) | \mathbf{y}_i, \hat{\boldsymbol{\Theta}}^{(r)})$ است و

$$Q_i(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)}) = \frac{1}{\gamma} \{ \log |\boldsymbol{\Sigma}_i^{-1}| + \log |\mathbf{D}_i^{-1}| - \text{tr}(\mathbf{D}_i^{-1} \hat{\mathbf{B}}_i^{(r)}) - \text{tr}(\boldsymbol{\Sigma}_i^{-1} \hat{\boldsymbol{\Psi}}_i^{(r)}(\boldsymbol{\theta})) \},$$

که در آن

$$\begin{aligned} \hat{\mathbf{B}}_i^{(r)} &= E(\mathbf{b}_i \mathbf{b}_i^\top | \mathbf{y}_i, \hat{\boldsymbol{\Theta}}^{(r)}) = \hat{\mathbf{b}}_i^{(r)} \hat{\mathbf{b}}_i^{(r)\top} + \hat{\mathbf{V}}_{\mathbf{b}_i}^{(r)}, \\ \hat{\boldsymbol{\Psi}}_i^{(r)}(\boldsymbol{\beta}) &= E(\mathbf{e}_i \mathbf{e}_i^\top | \mathbf{y}_i, \hat{\boldsymbol{\Theta}}^{(r)}) \\ &= (\mathbf{y}_i - \tilde{\mathbf{X}}_i \boldsymbol{\theta} - \mathbf{Z}_i \hat{\mathbf{b}}_i^{(r)}) (\mathbf{y}_i - \tilde{\mathbf{X}}_i \boldsymbol{\theta} - \mathbf{Z}_i \hat{\mathbf{b}}_i^{(r)})^\top + \mathbf{Z}_i \hat{\mathbf{V}}_{\mathbf{b}_i}^{(r)} \mathbf{Z}_i^\top, \\ \hat{\mathbf{b}}_i^{(r)} &= E(\mathbf{b}_i | \mathbf{y}_i, \hat{\boldsymbol{\Theta}}^{(r)}) = \hat{\mathbf{D}}_i^{(r)} \mathbf{Z}_i^\top \hat{\mathbf{V}}_i^{(r)-1} (\mathbf{y}_i - \tilde{\mathbf{X}}_i \hat{\boldsymbol{\theta}}^{(r)}), \\ \hat{\mathbf{V}}_{\mathbf{b}_i}^{(r)} &= \text{Cov}(\mathbf{b}_i | \mathbf{y}_i, \hat{\boldsymbol{\Theta}}^{(r)}) = (\hat{\mathbf{D}}_i^{(r)} + \mathbf{Z}_i^\top \hat{\boldsymbol{\Sigma}}_i^{(r)-1} \mathbf{Z}_i)^{-1}. \end{aligned}$$

گام M: به هنگام کردن $\hat{\boldsymbol{\Theta}}^{(r)}$ توسط $\hat{\boldsymbol{\Theta}}^{(r+1)} = \max_{\boldsymbol{\Theta}} Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)})$.

برای انتخاب متغیرهای مهم و برآورد ضرایب آن‌ها، به لگاریتم تابع درستنمایی (۳) به صورت

$$\ell_c^{\text{Pen}}(\boldsymbol{\Theta}) = \ell_c(\boldsymbol{\Theta}) - n \sum_{k=1}^p p_{\lambda_n}(|\beta_k|). \quad (۴)$$

تاوانیده می‌شود، که در آن p_{λ_n} را تابع تاوان و λ_n را پارامتر تاوان می‌نامند. به طور مشابه تابع Q تاوانیده در گام E به صورت $Q_i^{\text{Pen}}(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)}) = Q_i(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(r)}) - n \sum_{k=1}^{p_n} p_{\lambda_n}(|\beta_k|)$ در نظر

گرفته می‌شود. فن ولی (۲۰۰۱) یک تابع تاوان مقعر به نام اسکد^۱ را معرفی کردند که دارای هر سه ویژگی مطلوب تنکی، نارایی و پیوستگی است. این تابع تاوان دارای مشتق مرتبه اول به صورت

$$p'_{\lambda_n}(|\beta|) = \lambda_n \left\{ I(|\beta| \leq \lambda_n) + \frac{(a\lambda_n - |\beta|)_+}{(a-1)\lambda_n} I(|\beta| > \lambda_n) \right\}; \quad a > 2.$$

است. برای حل مسئله بهینه‌سازی در گام M، از روش تکراری تقریب موضعی توان دوم (LQA) استفاده می‌شود (فن ولی، ۲۰۰۱). در این روش با استفاده از بسط تیلور و داشتن یک مقدار اولیه $\beta_{\circ,k}$ ، $|\beta_{\circ,k}| > 0$ مشتق تابع تاوان به صورت

$$[p_{\lambda_n}(|\beta_k|)]' = q_{\lambda_n}(|\beta_k|) \text{sign}(\beta_k) \approx \frac{q_{\lambda_n}(|\beta_{\circ,k}|)}{|\beta_{\circ,k}|} \beta_k,$$

تقریب زده می‌شود، که در آن $\text{sign}(a) = I(a > 0) - I(a < 0)$ و $q_{\lambda_n}(|\beta_k|) = p'_{\lambda_n}(|\beta_k|)$ مشتق مرتبه اول $p_{\lambda_n}(\cdot)$ است. هانتز ولی (۲۰۰۵) خاصیت همگرایی الگوریتم LQA را بررسی و برای محاسبه برآوردگر تاوانیده، الگوریتم MM را پیشنهاد دادند. در این الگوریتم برای کنترل عدم شناساپذیری برآوردگر در نقطه ۰، با قبول اندکی انحراف به میزان ϵ ، تقریب LQA را بهبود دادند. سپس برآوردگر را براساس الگوریتم LQA با حل مسئله بهینه‌سازی به صورت $\hat{\beta} = \arg \max_{\beta} \{Q_i(\Theta | \hat{\Theta}^{(r)})\}$ برآوردگر و همچنین سرعت همگرایی تقریب، ضروری است. فن ولی (۲۰۰۱) نشان دادند که الگوریتم LQA پس از تعداد کمی تکرار به همگرایی می‌رسد. در هر تکرار به محض اینکه برخی از مؤلفه‌های $|\beta_k|$ کوچکتر از $\epsilon = 10^{-6}$ شود، متغیرهای مربوط به آن‌ها به عنوان متغیر بی‌تاثیر شناخته و مقدار صفر برای آن‌ها در نظر گرفته می‌شود.

¹SCAD

²Local quadratic approximation (LQA)

گام M: به هنگام کردن $\hat{\beta}^{(r)}$ ، $\hat{D}^{(r)}$ و $\hat{\Sigma}^{(r)}$ از طریق ماکسیم کردن $Q(\Theta|\hat{\Theta}^{(r)})$ که به صورت

$$\hat{\theta}^{(r+1)} = \left(\sum_{i=1}^n \tilde{X}_i^T \hat{\Sigma}_i^{(r)-1} \tilde{X}_i + nE_n(\hat{\theta}^{(r)})^{-1} \right)^{-1} \sum_{i=1}^n \tilde{X}_i^T \hat{\Sigma}_i^{(r)-1} (y_i - z_i \hat{b}_i^{(r)}),$$

$$\hat{D}^{(r+1)} = n^{-1} \sum_{i=1}^n \hat{B}_i^{(r)}, \quad \hat{\Sigma}^{(r+1)} = n^{-1} \sum_{i=1}^n \hat{\Psi}_i^{(r)}(\beta),$$

حاصل می‌شوند، که در آن $E_n(\hat{\theta}^{(r)}) = \text{diag}\left\{\frac{q_{\lambda_n}(|\beta_1|)}{\epsilon+|\beta_1|}, \dots, \frac{q_{\lambda_n}(|\beta_p|)}{\epsilon+|\beta_p|}, \mathbf{0}_{h_n}\right\}$ و $\epsilon = 10^{-6}$ برداری h_n بعدی با عناصر صفر است.

گام‌های E و M تا زمان همگرایی به مقداری ثابت تکرار می‌شود تا برآوردگر ماکسیم درست‌نمایی توانیده برای پارامترهای $(\hat{\theta}, \hat{D}, \hat{\Sigma}) = \hat{\Theta}$ ، بدست آید.

۳.۱ انتخاب پارامترهای تاوان

در مطالعات عددی مقاله، از اسپلاین مکعبی ($d = 3$) استفاده می‌شود. این نکته که تعداد گره‌ها با افزایش حجم نمونه افزایش یابند اهمیت داشته و از سوی دیگر تعداد زیاد گره‌ها ممکن است واریانس برآوردگرها را افزایش دهد. بنابراین، تعداد گره‌ها باید به درستی انتخاب شود تا بین اریبی و واریانس برآوردگرها تعادل ایجاد گردد. برای سهولت در محاسبات معمولاً گره‌هایی با فواصل برابر و تعداد گره‌ها به صورت $L_n \approx n^{1/(2r+1)}$ انتخاب می‌شود. این استراتژی در هی و همکاران (۲۰۰۲)، چن و چو (۲۰۰۷) و سینها و ستار (۲۰۱۵) اجرا شده است. واضح است که $r > 0$ زیرا $r = 0$ تعداد n گره را در نظر می‌گیرد که انتخاب درستی نیست. معمولاً r را برابر ۱ در نظر می‌گیرند؛ به طور مثال برای نمونه‌هایی به حجم ۵۰، ۱۰۰ و ۲۰۰، تعداد گره‌ها به ترتیب برابر ۴، ۵ و ۶ خواهد بود. بنابراین اگر اسپلاین مکعبی و تعداد گره‌ها به صورت $L_n \approx n^{1/(2r+1)}$ در نظر گرفته شود، تعداد توابع پایه به صورت $h_n \approx n^{1/(2r+1)} + 4$ است؛ که برای حجم نمونه‌های ۵۰، ۱۰۰ و ۲۰۰ به ترتیب ۸، ۹ و ۱۰ تابع پایه وجود خواهد داشت. فن و لی (۲۰۰۱) از دیدگاه بیزی مقدار $a = 3.7$ را به عنوان یک مقدار مناسب برای مسائل مختلف پیشنهاد دادند و برای انتخاب مقدار مناسب پارامتر تنظیم λ_n ، از معیار اعتبار سنجی متقابل تعمیم یافته به صورت

$$GCV(\lambda_n) = \frac{RSS(\lambda_n)/n}{(1 - \text{tr}(d(\lambda_n))/n)^2}$$

استفاده کردند، که در آن $RSS(\lambda_n) = (\mathbf{y} - \tilde{\mathbf{X}}^\top \hat{\boldsymbol{\theta}} - \mathbf{Z}^\top \hat{\mathbf{b}})^\top \hat{\mathbf{V}}^{-1} (\mathbf{y} - \tilde{\mathbf{X}}^\top \hat{\boldsymbol{\theta}} - \mathbf{Z}^\top \hat{\mathbf{b}})$ و

$$d(\lambda_n) = \mathbf{X}(\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{X} + \lambda_n \mathbf{I})^{-1} \mathbf{X}^\top \hat{\mathbf{V}}^{-1} + \mathbf{Z} \hat{\mathbf{D}} \mathbf{Z}^\top (\mathbf{I} - \mathbf{X}[(\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \hat{\mathbf{V}}^{-1}]).$$

مقدار λ_n از حل مساله بهینه‌سازی $GCV(\lambda_n)$ برآورد می‌شود.
 $\hat{\lambda}_n = \arg \min_{\lambda_n \in \mathbb{R}^+} GCV(\lambda_n)$

۴ بررسی رفتار مجانبی برآوردگرها

این بخش به بیان توزیع مجانبی برآوردهای $\hat{\boldsymbol{\beta}}_n$ و \hat{g} پرداخته است. در اکثر متون آماری در حیطه انتخاب متغیر به منظور سادگی در محاسبات، پارامترهای درست $\boldsymbol{\beta}_{n^0}$ به صورت $\boldsymbol{\beta}_{n^0} = (\boldsymbol{\beta}_{n^0,1}^\top, \boldsymbol{\beta}_{n^0,2}^\top)^\top$ و ماتریس طرح مدل به صورت $\mathbf{X}_i = (\mathbf{X}_{i(1)}, \mathbf{X}_{i(2)})$ تفکیک می‌شوند. بدون کاستن از کلیت مساله فرض می‌شود $\boldsymbol{\beta}_{n^0,2} = \mathbf{0}$ و تمام اعضای $\boldsymbol{\beta}_{n^0,1}$ که برداری s_n^* بعدی است غیر صفر باشد. در مقاله حاضر، ضرایب رگرسیونی درست به صورت $\boldsymbol{\theta}_{n^0} = (\boldsymbol{\beta}_{n^0,1}^\top, \boldsymbol{\beta}_{n^0,2}^\top, \boldsymbol{\alpha}_{n^0}^\top)^\top$ است، که در آن $\boldsymbol{\alpha}_{n^0}$ برداری h_n بعدی و وابسته به تابع ناپارامتری g_0 است. به منظور جداسازی فضای پارامتر به فضای صفر و غیر صفر، $\boldsymbol{\theta}_{n^0}$ را به صورت $(\boldsymbol{\theta}_{n^0,1}^\top, \boldsymbol{\theta}_{n^0,2}^\top)^\top$ بازتفکیک کرده به طوری که $\boldsymbol{\theta}_{n^0,1} = (\boldsymbol{\beta}_{n^0,1}^\top, \boldsymbol{\alpha}_{n^0}^\top)^\top$ یک بردار $(s_n = s_n^* + h_n)$ بعدی است که تمام اعضای آن غیر صفراند و $\boldsymbol{\theta}_{n^0,2} = \boldsymbol{\beta}_{n^0,2} = \mathbf{0}$. متعاقباً، مقادیر برآورد و ماتریس طرح به ترتیب به صورت $\hat{\boldsymbol{\theta}}_n = (\hat{\boldsymbol{\theta}}_{n1}^\top, \hat{\boldsymbol{\theta}}_{n2}^\top)^\top$ و $\tilde{\mathbf{X}}_i = (\tilde{\mathbf{X}}_{i(1)}^\top, \tilde{\mathbf{X}}_{i(2)}^\top)^\top$ بازتفکیک می‌شوند، که در آن $\hat{\boldsymbol{\theta}}_{n1} = (\hat{\boldsymbol{\beta}}_{n1}^\top, \hat{\boldsymbol{\alpha}}_n^\top)^\top$ ، $\tilde{\mathbf{X}}_{i(1)} = (\mathbf{X}_{i(1)}^\top, \mathbf{B}_i(t_i)^\top)^\top$ ، $\hat{\boldsymbol{\theta}}_{n2} = \hat{\boldsymbol{\beta}}_{n2}$ و $\tilde{\mathbf{X}}_{i(2)} = \mathbf{X}_{i(2)}$ فرض کنید $\{p'_{\lambda_n} | \beta_{n^0,k}, \beta_{n^0,k} \neq 0\}$ و $a_n = \max_{1 \leq k \leq p_n} \{p'_{\lambda_n} | \beta_{n^0,k}, \beta_{n^0,k} \neq 0\}$ و $b_n = \max_{1 \leq k \leq p_n} \{p''_{\lambda_n} | \beta_{n^0,k}, \beta_{n^0,k} \neq 0\}$ مولفه k ام $\boldsymbol{\beta}_{n^0}$ است.

اگر گام M در الگوریتم EM چند جواب داشته باشد، تنها دنباله‌ای از برآوردگر سازگار $\hat{\boldsymbol{\theta}}_n$ در نظر گرفته می‌شود، یعنی وقتی $n \rightarrow \infty$ آنگاه $\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_{n^0} \rightarrow \mathbf{0}$ و $\sup_t |\mathbf{B}^\top(t) \hat{\boldsymbol{\alpha}}_n - g_0(t)| \xrightarrow{p} 0$. علاوه بر این، فرض کنید به ازای هر i ، $n_i = m < \infty$ باشد؛ مجموعه سه تایی $(\mathbf{y}_i, \mathbf{X}_i, t_i)$ ، $i = 1, \dots, n$ متغیرهایی مستقل و هم توزیع باشند و مشتق k ام تابعی مانند $d(\cdot)$ به صورت $d^{(k)}(\cdot)$ تعریف شود. برای مطالعه رفتار حدی برآوردگرها، شرایط نظم زیر را در نظر بگیرید.

(۱.۰A) به ازای هر $k \geq 2$ ، مشتق k ام تابع $g_0(t_i)$ کراندار است.

(۲.۰A) ماتریس کوواریانس برآورد شده $\widehat{\Sigma}_i$ در شرط $\|\widehat{\Sigma}_i - \Sigma_i\| = O_p(n^{-1/2})$ صدق می‌کند.

(۳.۰A) $b_1 \leq \lambda_{\min}(n^{-1} \sum_{i=1}^n \tilde{X}_i^T \Sigma_i^{-1} \tilde{X}_i) \leq \lambda_{\max}(n^{-1} \sum_{i=1}^n \tilde{X}_i^T \Sigma_i^{-1} \tilde{X}_i) \leq b_2$ که در آن b_1 و b_2 مقادیر ثابت و مثبت و λ_{\min} و λ_{\max} کوچکترین و بزرگترین مقدار ویژه هستند.

(۴.۰A) وقتی $n \rightarrow \infty$ ، آنگاه $\min_{1 \leq k \leq s_n} |\theta_{n \cdot k}| / \lambda_n \rightarrow \infty$ ، $s_n^2 n^{-1} = o(1)$ ، $\lambda_n \rightarrow \circ$ ، $p_n s_n^4 (\log n)^6 = o(n^2 \lambda_n^2)$ ، $\log(p_n) = o(n \lambda_n^2 / (\log n))$ ، $s_n^2 (\log n)^4 = o(n \lambda_n^2)^2$ و $p_n s_n^2 (\log n)^4 = o(n^2 \lambda_n^4)$.

(۵.۰A) تابع تاوان در شرایط $\liminf_{n \rightarrow \infty} \liminf_{\theta \rightarrow \theta_0^+} p'_{\lambda_n}(\theta) / \lambda_n > \circ$ ، $a_n = O(n^{-1/2})$ ، $b_n \rightarrow \infty$ صدق می‌کند و همچنین مقادیر ثابت C و D وجود دارند به طوری که به ازای هر $\theta_1, \theta_2 > \lambda_n C$ خواهیم داشت $|p''_{\lambda_n}(\theta_1) - p''_{\lambda_n}(\theta_2)| \leq D|\theta_1 - \theta_2|$.

ملاحظه ۱. در مدل اولیه فرض می‌شود که $g(\cdot)$ تابعی هموار است، یعنی کلیه مشتقات آن وجود داشته باشد، لذا شرط (۱.۰A) برقرار است. شرط (۲.۰A) مشابه شرط استفاده شده در لیانگ و زیگر (۱۹۸۶) است که فرض می‌کند پارامتر همبستگی $\widehat{\theta}$ در شرط $\sqrt{n}(\widehat{\theta} - \theta_0) = O_p(1)$ صدق می‌کند. با توجه به این که دنباله برآوردگرهای سازگار در الگوریتم EM در نظر گرفته می‌شود، این شرط نیز برقرار است. شرط (۳.۰A) یکی از شرایط معمول در رگرسیون برای داده‌های مستقل است. شرایط (۴.۰A) و (۵.۰A) از شرایط معمول در متون انتخاب متغیر است که توسط فن ولی (۲۰۰۱) اثبات شده است.

قضیه ۱. با در نظر گرفتن شرایط (۱.۰A)-(۵.۰A)، وقتی $n \rightarrow \infty$ و $n^{-1} p_n^2 = o(1)$ ، آنگاه $\widehat{\theta}_n$ یک جواب الگوریتم EM است، به طوری که

$$i. \|\widehat{\theta}_n - \theta_{n \cdot}\| = O_p(\sqrt{p_n}(n^{-1/2} + a_n)),$$

$$ii. \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} (\widehat{g}(t_{ij}) - g_0(t_{ij}))^2 = O_p(n^{-2m/(2m+1)}).$$

برهان: برای اثبات قسمت (i)، فرض کنید $\alpha_n = \sqrt{p_n}(n^{-1/2} + a_n)$ و $\|U\| = C$ ، به طوری که یک مقدار ثابت و به اندازه کافی بزرگ است. نشان داده می‌شود که به ازای هر $\xi > \circ$ یک مقدار ثابت C وجود دارد که $\mathcal{P}\{\inf_{\|U\|=C} \ell_c^{Pen}(\theta_{n \cdot} + \alpha_n U) \geq \ell_c^{Pen}(\theta_{n \cdot})\} \geq 1 - \xi$ به عبارت دیگر با احتمال نزدیک به ۱، در مجموعه $\{\theta_{n \cdot} + \alpha_n U : \|U\| \leq C\}$ یک جواب محلی برای $\ell_c^{Pen}(\theta_{n \cdot})$

وجود دارد. با استفاده از $p_{\lambda_n}(\circ) = \circ$ داریم:

$$\begin{aligned} D_n(U) &= \ell_c^{Pen}(\boldsymbol{\theta}_{n\circ} + \alpha_n U) - \ell_c^{Pen}(\boldsymbol{\theta}_{n\circ}) \\ &\geq \ell_c(\boldsymbol{\theta}_{n\circ} + \alpha_n U) - \ell_c(\boldsymbol{\theta}_{n\circ}) + n \sum_{i=1}^n \{p_{\lambda_n}(|\theta_{n\circ,k} + \alpha_n U_k|) - p_{\lambda_n}(|\theta_{n\circ,k}|)\} \\ &\equiv I + II \end{aligned}$$

با انجام محاسبات جبری داریم

$$I = \frac{n\alpha_n^2}{\gamma} U^\top \hat{\Gamma}_n U - \alpha_n U^\top \sum_{i=1}^n \tilde{X}_i \Sigma_i^{-1} (Y_i - \tilde{X}_i \boldsymbol{\theta} - Z_i b_i) \equiv I_1 - I_2,$$

که در آن $\hat{\Gamma}_n = n^{-1} \sum_{i=1}^n \tilde{X}_i^\top \hat{\Sigma}_i^{-1} \tilde{X}_i$ و $\Gamma_n = n^{-1} \sum_{i=1}^n \tilde{X}_i^\top \Sigma_i^{-1} \tilde{X}_i$ عبارت‌های I_1 و I_2 را به صورت

$$\begin{aligned} I_1 &= \frac{n\alpha_n^2}{\gamma} U^\top (\Gamma_n - \hat{\Gamma}_n + \hat{\Gamma}_n) U = \frac{n\alpha_n^2}{\gamma} U^\top \hat{\Gamma}_n U + O_p(n^{-1/2}) n\alpha_n^2 \|U\|^2 \\ I_2 &= \alpha_n U^\top \sum_{i=1}^n \tilde{X}_i \Sigma_i^{-1} \epsilon_i \end{aligned}$$

خواهیم داشت. با استفاده از شرط (۳.A) می‌توان نشان داد

$$\begin{aligned} \left| \alpha_n U^\top \sum_{i=1}^n \tilde{X}_i \Sigma_i^{-1} \epsilon_i \right| &\leq \alpha_n \left\| \sum_{i=1}^n \tilde{X}_i \Sigma_i^{-1} \epsilon_i \right\| \|U\| \\ &= O_p(\alpha_n \sqrt{np_n}) \|U\| \\ &= O_p(\alpha_n^2 n) \|U\|. \end{aligned}$$

بنابراین عبارت I_2 به صورت $I_2 = O_p(\alpha_n^2 n) \|U\|$ بدست می‌آید. مشابه اثبات فن و لی (۲۰۰۱) می‌توان نشان داد

$$II = O_p(\sqrt{s_n n \alpha_n a_n}) \|U\| + O_p(n \alpha_n^2 b_n) \|U\|.$$

برای اثبات کافی است نشان داده شود عبارت

$$\frac{\partial \ell_c^{Pen}(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}} = \begin{cases} < \circ & -\xi_n < \theta_{n,k} < \circ \\ > \circ & \circ < \theta_{n,k} < \xi_n \end{cases}$$

به ازای هر $\theta_{n,k} \neq \circ$ و $k = s_n + 1, \dots, p_n$ برقرار است، که در آن $\xi_n = C\sqrt{p_n/n}$. واضح است که $\frac{\partial \ell_c^{Pen}(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}} = \frac{\partial \ell_c(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}} + np'_{\lambda_n}(|\theta_{n,k}|)\text{sgn}(\theta_{n,k})$ و مشابه فن ولی (۲۰۰۱) می‌توان نشان $\frac{\partial \ell_c(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}} = O_p(\sqrt{np_n})$ بنا براین

$$\frac{\partial \ell_c^{Pen}(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}} = n\lambda_n \left\{ \frac{p'_{\lambda_n}(|\theta_{n,k}|)}{\lambda_n} \text{sgn}(\theta_{n,k}) + O_p\left(\frac{\sqrt{p_n/n}}{\lambda_n}\right) \right\}.$$

از آنجایی که $\frac{\sqrt{p_n/n}}{\lambda_n} \rightarrow \circ$ و $\liminf_{\theta \rightarrow \circ^+} p'_{\lambda_n}(\theta)/\lambda_n > \circ$ ، علامت $\theta_{n,k}$ تعیین کننده علامت $\frac{\partial \ell_c^{Pen}(\boldsymbol{\theta}_n)}{\partial \theta_{n,k}}$ است. بنا براین اثبات کامل می‌شود.

بخش همگرایی در احتمال و همگرایی در توزیع مستقیماً از طریق قانون ضعیف اعداد بزرگ، قضیه حد مرکزی و قضیه اسلاتسکی اثبات می‌شود.

۵ مطالعه شبیه‌سازی

در این بخش، فرایند برازش مدل، برآورد پارامترهای مدل و انتخاب متغیر به صورت عددی ارزیابی می‌شود. ابتدا مجموعه‌ای از داده‌های شبیه‌سازی شده توسط روش پیشنهادی مقاله یعنی برآوردگر ماکسیمم درستنمایی تاوانیده در مدل خطی-جزئی با اثرات آمیخته (P-PLMM) مورد تحلیل قرار می‌گیرد. سپس به منظور بررسی عملکرد روش P-PLMM در مسئله برآورد و انتخاب متغیر همزمان، نتایج حاصل از تحلیل داده‌ها بدون رهیافت تابع تاوان یعنی روش برآوردگر ماکسیمم درستنمایی در مدل خطی-جزئی با اثرات آمیخته (PLMM) نیز گزارش می‌شود. در نظر داشته باشید که اگر داده‌ها ذاتاً دارای جزء ناپارامتری باشند، مدل خطی-جزئی بهتر از مدل خطی عمل می‌کند. برای نشان دادن این موضوع روش P-PLMM با روش برآوردگر ماکسیمم درستنمایی تاوانیده در مدل خطی با اثرات آمیخته (P-LMM) مقایسه می‌شود.

متغیرهای پاسخ از مدل زیر شبیه‌سازی می‌شوند

$$y_i(t_{ij}) = x_{\lambda,i}(t_{ij})\beta_{\lambda} + x_{\nu,i}(t_{ij})\beta_{\nu} + g(t_{ij}) + b_i + \varepsilon_i(t_{ij}),$$

که در آن مولفه‌ها به صورت $n_i = 5, j = 1, \dots, n_i, n = 50, 100, 200, i = 1, \dots, n$ ، $\mathbf{X}_{ij}^T = (x_{ij}^{(1)}, \dots, x_{ij}^{(p_n)})$ ، $\beta^T = (2, 3, 15, 2, 0, \dots, 0)$ ، $\sigma_b^2 = 0.25$ ، $b_i \sim N(0, \sigma_b^2)$ ، $t_{ij} \sim U(0, 1)$ و $x_{ij}^{(k)} \sim U(-1, 1)$ نویسندگان پیشنهاد‌های متفاوتی از جمله $p_n = \lfloor \frac{n}{4} \rfloor$ ، $p_n = \lfloor 4.5n^{1/4} \rfloor$ و $p_n = \lfloor \frac{n}{b \log(n)} \rfloor$ داشته‌اند که در آن‌ها $b > 1$ و $[a]$ جزء صحیح مقدار a است. موارد بیان شده برای حالت $p_n < n$ است. برای زمانی که $p_n > n$ نیز نرخ‌هایی مانند $\log(p_n) = o_p(n^b)$ پیشنهاد شده است که در آن $0 < b < 1$. انتخاب مقادیر p_n بسیار بزرگتر از n چالش‌هایی را به همراه خواهد داشت، به طوری که در این حالت اگر چه احتمال انتخاب متغیرهای موثر در مدل افزایش می‌یابد در عین حال متغیرهای بی‌تاثیر بیشتری در مدل وارد خواهند شد که این موضوع بر عملکرد روش‌های برآورد و انتخاب متغیر تاثیر منفی خواهد گذاشت. به منظور همپوشانی تمامی نرخ‌های پیشنهادی، موارد $(200, 16)$ ، $(100, 14)$ ، $(50, 11)$ برای حالت $p_n < n$ و موارد $(200, 2000)$ ، $(100, 500)$ ، $(30, 100)$ برای حالت $p_n > n$ در نظر گرفته شده است.

دقت برآورد هر یک از روش‌ها توسط میانگین توان دوم خطاها (MSE) براساس ۱۰۰ تکرار از شبیه‌سازی، مورد سنجش قرار گرفت که به صورت $MSE = \sum_{s=1}^{100} \|\hat{\beta}_n^s - \beta_{n^0}\|^2 / 100$ تعریف می‌شود که در آن برآوردگر $\hat{\beta}_n^s$ در s -امین مجموعه داده تولید شده است. ارزیابی عملکرد روش‌ها در مساله انتخاب متغیر توسط معیارهای (C, I) سنجیده می‌شود به طوری که C میانگین ضرایب صفر که به درستی صفر تخمین زده شده‌اند و I میانگین ضرایب غیر صفر است که به اشتباه صفر شده‌اند. برای ارائه توصیف جامع‌تر از عملکرد مساله انتخاب متغیر روش‌ها از معیارهای دیگری نیز استفاده شده است. UF بیانگر نسبت تکرارهایی است که ضرایب تاثیرگذار را به اشتباه صفر در نظر گرفته است، CF بیانگر نسبت تکرارهایی است که مدل را به درستی تعیین کرده‌اند و نسبت تکرارهایی که علاوه بر متغیرهای مهم، متغیرهای بی‌تاثیر را نیز در مدل وارد کرده‌اند با OF نشان داده شده است.

جدول ۱، نتایج مربوط به عملکرد روش‌های PLMM و P-LMM، P-PLMM در برآورد و انتخاب مدل برای مقادیر مختلف (n, p_n) است. از نظر دقت برآورد روش PLMM، MSE بیشتری نسبت به دو روش دیگر دارد، در حالیکه دو روش P-LMM و P-PLMM بسیار نزدیک به هم عمل می‌کنند. با

جدول ۱. مقایسه روش‌های ۱: PLMM، ۲: P-LMM و ۳: P-PLMM

حالت $p_n > n$						حالت $p_n < n$						روش
$(n, p) = (50, 100)$						$(n, p) = (50, 11)$						
OF	CF	UF	$I(\circ)$	$C(\gamma)$	MSE	OF	CF	UF	$I(\circ)$	$C(\gamma)$	MSE	
۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۰/۱۰	۸/۰۵	۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۰/۴۰	۱
۰/۷۷	۰/۲۳	۰/۰۰	۰/۰۰	۹۴/۱۹	۰/۴۶	۰/۵۰	۰/۵۰	۰/۰۰	۰/۰۰	۶/۲۸	۰/۱۸	۲
۰/۷۳	۰/۲۷	۰/۰۰	۰/۰۰	۹۴/۶۵	۰/۳۹	۰/۴۶	۰/۵۴	۰/۰۰	۰/۰۰	۶/۴۸	۰/۱۹	۳
$(n, p) = (100, 500)$						$(n, p) = (100, 14)$						
۱/۰۰	۰/۰۰	۰/۰۳	۰/۰۳	۴۷/۰۲	۱۴۹۹/۱۴	۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۲	۰/۲۷	۱
۰/۷۰	۰/۳۰	۰/۰۰	۰/۰۰	۴۹۵/۷۲	۰/۰۶	۰/۳۸	۰/۶۲	۰/۰۰	۰/۰۰	۹/۴۷	۰/۰۹	۲
۰/۴۶	۰/۵۴	۰/۰۰	۰/۰۰	۴۹۶/۲۵	۰/۰۴	۰/۳۱	۰/۶۹	۰/۰۰	۰/۰۰	۹/۵۶	۰/۰۹	۳
$(n, p) = (200, 2000)$						$(n, p) = (200, 16)$						
۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۱۱۳۷/۶۲	۱۲۵/۴۱	۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۹	۰/۱۴	۱
۰/۱۱	۰/۸۹	۰/۰۰	۰/۰۰	۱۹۹۶/۸۹	۰/۰۲	۰/۰۷	۰/۹۳	۰/۰۰	۰/۰۰	۱۱/۸۳	۰/۰۴	۲
۰/۰۸	۰/۹۲	۰/۰۰	۰/۰۰	۱۹۹۶/۹۳	۰/۰۱۸	۰/۰۴	۰/۹۶	۰/۰۰	۰/۰۰	۱۱/۸۶	۰/۰۴	۳

جدول ۲. برآوردهای P-PLMM در حالت‌های $p_n > n$ و $p_n < n$

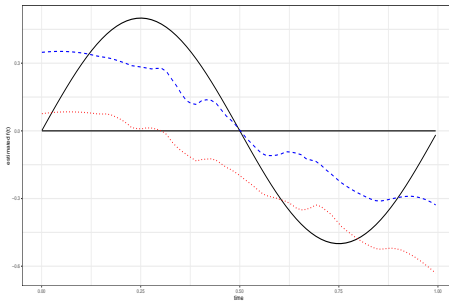
حالت $p_n > n$					حالت $p_n < n$					
β_4	β_3	β_2	β_1	(n, p_n)	β_4	β_3	β_2	β_1	معیار	(n, p_n)
۰/۰۶	۰/۰۵	۰/۱۱	۰/۰۸	$(50, 100)$	۰/۰۱	۰/۰۰	۰/۰۱	۰/۰۴	Bias	$(50, 11)$
۰/۰۸	۰/۰۸	۰/۰۹	۰/۰۹		۰/۱۰	۰/۱۱	۰/۱۱	۰/۱۰	SD۱	
۰/۰۷	۰/۰۷	۰/۰۸	۰/۰۹		۰/۰۸	۰/۰۸	۰/۰۸	۰/۰۷	SD۲	
۰/۹۳	۰/۹۹	۰/۹۶	۰/۹۶		۰/۹۱	۰/۹۵	۰/۹۷	۰/۹۵	CP	
۰/۰۳	۰/۰۲	۰/۰۳	۰/۰۴	$(100, 500)$	۰/۰۲	۰/۰۱	۰/۰۵	۰/۰۱	Bias	$(100, 14)$
۰/۰۸	۰/۰۷	۰/۰۸	۰/۰۷		۰/۰۸	۰/۰۷	۰/۰۸	۰/۰۷	SD۱	
۰/۰۵	۰/۰۵	۰/۰۵	۰/۰۵		۰/۰۵	۰/۰۵	۰/۰۵	۰/۰۵	SD۲	
۰/۹۴	۰/۹۵	۰/۹۳	۰/۹۳		۰/۹۴	۰/۹۵	۰/۹۴	۰/۹۱	CP	
۰/۰۱	۰/۰۱	۰/۰۴	۰/۰۲	$(200, 2000)$	۰/۰۱	۰/۰۱	۰/۰۰	۰/۰۴	Bias	$(200, 16)$
۰/۰۸	۰/۰۷	۰/۰۸	۰/۰۸		۰/۰۶	۰/۰۶	۰/۰۶	۰/۰۶	SD۱	
۰/۰۶	۰/۰۵	۰/۰۵	۰/۰۵		۰/۰۴	۰/۰۴	۰/۰۴	۰/۰۴	SD۲	
۰/۹۴	۰/۹۲	۰/۹۶	۰/۹۳		۰/۹۵	۰/۹۵	۰/۹۵	۰/۹۶	CP	

این وجود MSE روش پیشنهادی همواره کمتر از P-LMM است. از نظر انتخاب مدل مشاهده می‌شود که PLMM انتخاب مدل انجام نمی‌دهد اما روش‌های P-LMM و P-LMM تمامی متغیرهای مهم را انتخاب کرده است زیرا معیار I آن‌ها صفر است. همچنین روش پیشنهادی به دلیل داشتن مقادیر بالا در نرخ‌های C و CF و مقادیر کمتر در OF ، نسبت به رقیب خود یعنی P-LMM عملکرد بهتری داشته است. با توجه به معیارهای CF و OF ، هنگامی که $p_n > n$ نسبت به حالت $p_n < n$ روش‌ها در انتخاب مدل عملکرد پایینی دارند و ضرایب صفر تمایل بیشتری دارند که در مدل گنجانده شوند، اما با

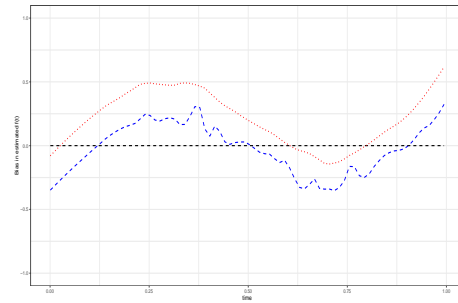
افزایش حجم نمونه هر دو معیار بهبود می‌یابند. برای بررسی بیشتر عملکرد روش پیشنهادی، نتایج مربوط به اریبی، خطای استاندارد مجانبی، خطای استاندارد نمونه‌ای و احتمال پوشش نمونه‌ای برای بازه اطمینان ۹۵٪ در جدول ۲ گزارش شده است. نتایج نشان می‌دهد که خطاهای استاندارد نمونه‌ای و مجانبی نزدیک به هم برآورد شده‌اند و احتمال پوشش در اغلب موارد نزدیک به ۹۵٪ است که اینها بیانگر عملکرد خوب روش است. به عبارت دیگر برآوردها رفتار مجانبی خوبی را نشان داده‌اند که این امر با نتایج قضایای حدی بیان شده مطابقت دارد. همچنین رفتار تابع $g(t)$ از طریق رسم نمودارهایی بررسی شده است. این نمودارها شامل منحنی‌های $\hat{g}(t)$ ، اریبی، انحراف استاندارد و احتمال پوشش هستند. شکل‌های ۱ و ۲ به ترتیب نتایج مربوط به حجم نمونه ۵۰ و ۱۰۰ هستند. در هر دو تصویر نتایج برای حالت‌های $p_n < n$ (منحنی خط‌چین) و $p_n > n$ (منحنی نقطه‌چین) رسم شده‌اند. ملاحظه می‌شود که منحنی‌های برآورد شده تقریباً دارای تابع سینوسی هستند و این امر در حجم نمونه بیشتر یعنی شکل ۲ مشهودتر است. منحنی اریبی برآوردها حول خط صفر در نوسان هستند و برای حجم نمونه ۱۰۰ میزان اریبی کمتر از حجم نمونه ۵۰ است. همچنین برای حجم نمونه ۵۰ میزان اریبی حالت $p_n < n$ کمتر از $p_n > n$ است در حالی که در حجم نمونه بالا میزان اریبی در دو حالت رفتار یکسانی را نشان می‌دهد. در هر دو شکل مقادیر انحراف استاندارد حالت $p_n < n$ کمتر از $p_n > n$ است. همچنین منحنی احتمال پوشش در هر دو شکل حول خط ۹۵٪ در نوسان است.

۶ تحلیل داده‌های HIV

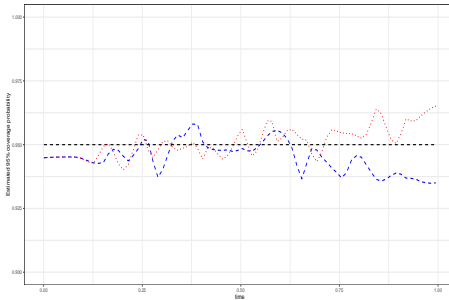
در این قسمت با استفاده از داده‌های مربوط به آزمایشات بیماری ایدز، رفتار برآوردهای پیشنهادی مطالعه می‌شود. این داده‌ها تعداد ۳۶۹ مرد مبتلا به HIV را، مورد بررسی قرار می‌دهد. در این مثال، متغیر پاسخ، تعداد سلول‌های $CD4$ و متغیرهای تبیینی شامل استفاده یا عدم استفاده از مواد مخدر ($DRUG$)، تعداد شرکای جنسی ($SEXP$)، تعداد بسته‌های سیگار مورد استفاده در روز ($SMOKE$)، میزان افسردگی ($CESD$)، سن (AGE) و سال ($YEAR$) است. به منظور اینکه توزیع داده‌ها، به توزیع نرمال نزدیک باشد، تبدیل ریشه دوم متغیر پاسخ در نظر گرفته شده است. زیگر و دیگل (۱۹۹۴) دریافتند که این داده‌ها دارای همبستگی درونی از نوع ساختار همبستگی تبادل پذیر با پارامتر همبستگی $\rho = ۰,۵۰۹$ هستند. برای بهره‌مندی از انعطاف بالای مدل‌های خطی-جزئی متغیر $YEAR$ به عنوان جزء ناپارامتری و سایر متغیرها به صورت خطی وارد مدل شده‌اند. از آنجایی که اثرات متقابل نیز ممکن است در مدل اهمیت داشته باشند، تمامی اثرات متقابل متغیرها نیز در مدل وارد شده و از روش پیشنهادی برای



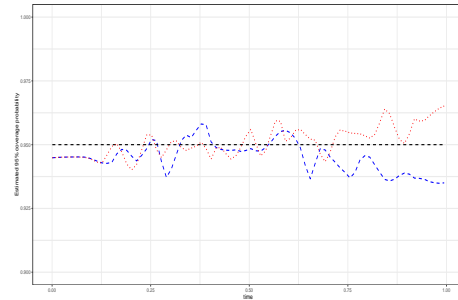
(ب)



(الف)



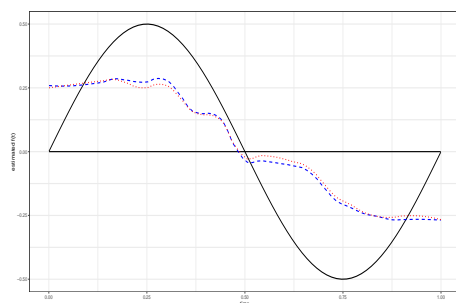
(د)



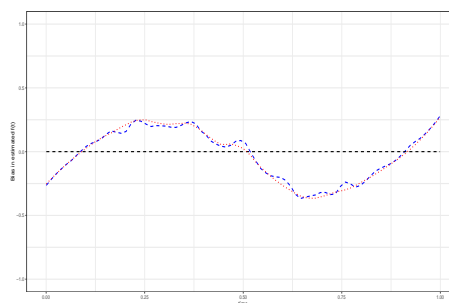
(ج)

شکل ۱. به ازای $n = 50$ ، منحنی‌های واقعی (خط پر)، برای $p = 11$ (خط چین آبی) و برای $p = 100$ (نقطه چین قرمز)، الف: برآورد $g(t)$ ، ب: اریبی، ج: انحراف استاندارد و د: احتمال پوشش.

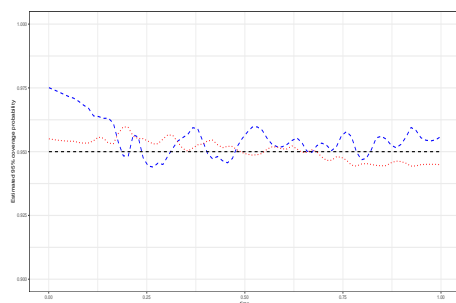
انتخاب متغیرهای مهم استفاده می‌شود. به منظور مقایسه عملکرد روش P-PLMM با دو روش دیگر یعنی P-LMM و PLMM از معیار انحراف استاندارد مجانبی (SD) و کمترین توان های دوم خطای تجربی (MSE) پارامترها استفاده گردید. در اینجا خطای استاندارد مجانبی، جذر واریانس مجانبی پارامترها یعنی جذر عناصر قطر اصلی ماتریس $I_{\beta_{n \times 1}, \beta_{n \times 1}}^{-1}$ در قضیه ۲ بوده و میانگین توان های دوم خطای تجربی (MSE) از طریق میانگین گرفتن توان های دوم خطا از یک نمونه‌گیری بوت‌استرپ که به صورت $MSE(\hat{\beta}_k) = \frac{1}{S} \sum_{s=1}^S (\hat{\beta}_k - \hat{\beta}_k^{(s)})^2$ است، حاصل می‌شود که در آن مقدار برآورد شده پارامتر k ام از کل نمونه و $\hat{\beta}_k^{(s)}$ مقدار برآورد شده پارامتر از نمونه بوت‌استرپ s ام است. همچنین برای شناسایی بهترین مدل پشتیبانی شده توسط داده‌ها، شاخص‌های مجموع توان های دوم رگرسیون (SSR)، مجموع توان های دوم خطا (SSE)، ضریب تعیین (R^2) و میانگین توان های دوم خطای پیشگویی



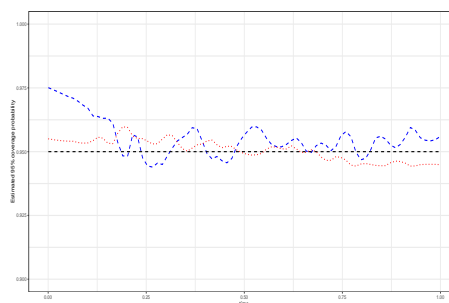
(ب)



(الف)



(د)



(ج)

شکل ۲. به ازای $n = 100$ ، منحنی‌های واقعی (خط پر)، برای $p = 14$ (خط چین آبی) و برای $p = 500$ (نقطه‌چین قرمز)، الف: برآورد $g(t)$ ، ب: اریبی، ج: انحراف استاندارد و د: احتمال پوشش.

$(MSPE)$ محاسبه شدند. نتایج در جدول ۳ نشان می‌دهد که برآوردگر P-PLMM در مقایسه با دو روش دیگر دارای مقادیر SD و MSE کمتری هستند. اما این تفاوت زیاد چشمگیر نیست، با این وجود در روش P-PLMM، شاخص‌های SSR و R^2 دارای بیشترین مقدار بوده و شاخص‌های SSE و $MSPE$ کمترین مقدار را دارند که عملکرد بهتر این روش را تایید می‌کنند. همچنین مشاهده می‌شود که روش PLMM انتخاب متغیر انجام نمی‌دهد و در روش P-LMM به غیر از متغیر $CESD$ تمام متغیرها وارد مدل شده‌اند. تحت مدل P-PLMM، متغیرهای $SMOKE$ ، $SEXP$ ، $CESD$ و اثر متقابل $SMOKE * SEXP$ به عنوان متغیرهای مهم شناسایی شده‌اند.

جدول ۳. برآورد پارامترها به همراه SD و MSE تحت برازش سه مدل برای داده‌های بیماری ایدز.

P-PLMM $\hat{\beta}(SD, MSE)$	P-LMM $\hat{\beta}(SD, MSE)$	PLMM $\hat{\beta}(SD, MSE)$	متغیرها
۰/۰۰ (۰/۰۰, ۰/۰۰)	-۰/۶۶۱ (۰/۰۰۵, ۰/۰۰۱)	۰/۰۸۴ (۰/۰۰۴, ۰/۰۰۱)	AGE
۰/۴۱۷ (۰/۰۱۴, ۰/۰۰۱)	۲/۴۹۸ (۰/۰۳۰, ۰/۰۰۱)	۰/۷۳۰ (۰/۰۲۴, ۰/۱۶۳)	SMOKE
۰ (۰, ۰/۰۰۱)	۶/۶۰۸ (۰/۰۳۸, ۰/۰۰۴)	۰/۸۱۱ (۰/۰۴۷, ۰/۳۳۰)	DRUG
۰/۰۸۴ (۰/۰۰۶, ۰/۰۰۰)	۱/۱۴۵ (۰/۰۰۸, ۰/۰۰۱)	۰/۱۷۰ (۰/۰۰۹, ۰/۰۷۴)	SEXP
-۰/۰۵۴ (۰/۰۰۳, ۰/۰۰۰)	۰ (۰, ۰/۰۰۴)	-۰/۰۶۹ (۰/۰۰۴, ۰/۰۰۱)	CESD
۰ (۰, ۰/۰۰۱)	۰/۰۱۶ (۰/۰۰۱, ۰/۰۰۲)	۰/۰۱۲ (۰/۰۰۱, ۰/۰۰۰)	AGE * SMOKE
۰ (۰, ۰/۰۰۰)	۰/۰۷۰ (۰/۰۰۵, ۰/۰۰۰)	-۰/۰۲۰ (۰/۰۰۴, ۰/۰۰۱)	AGE * DRUG
۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۰ (۰/۰۰, ۰/۰۰)	-۰/۰۱۲ (۰/۰۰۱, ۰/۰۰۱)	AGE * SEXP
۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۰ (۰/۰۰, ۰/۰۰)	-۰/۰۰۵ (۰/۰۰۲, ۰/۰۰۰)	AGE * CESD
۰/۰۰ (۰/۰۰, ۰/۰۰)	-۱/۱۶۶ (۰/۰۳۱, ۰/۰۰۱)	-۰/۲۷۸ (۰/۰۲۲, ۰/۲۹۴)	SMOKE * DRUG
۰/۰۴۲ (۰/۰۰۳, ۰/۰۰۱)	۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۳۸ (۰/۰۰۳, ۰/۰۰۲)	SMOKE * SEXP
۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۰ (۰/۰۰, ۰/۰۰)	-۰/۰۰۵ (۰/۰۰۱, ۰/۰۰۰)	SMOKE * CESD
۰/۰۰ (۰/۰۰, ۰/۰۰)	-۰/۰۵۵۸ (۰/۰۱۰, ۰/۰۰۱)	-۰/۰۷۹ (۰/۰۰۹, ۰/۰۷۴)	DRUG * SEXP
۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۲۰ (۰/۰۰۳, ۰/۰۰۱)	DRUG * CESD
۰/۰۰ (۰/۰۰, ۰/۰۰)	۰/۰۰ (۰/۰۰, ۰/۰۰۱)	-۰/۰۰۱ (۰/۰۰۳, ۰/۰۰۰)	SEXP * CESD
۶۷۵۷۱۸۹۰	۷۳۲۵۳۲۱۰	۶۶۵۶۱۶۸۰	SSR
۳۹۱۲۷۷۶۰	۵۴۳۶۸۸۸۰	۳۹۲۸۷۱۰۰	SSE
۰/۶۳۳	۰/۵۷۴	۰/۶۲۸	R^2
۴۰/۷۵۹	۷۷/۰۲۴	۴۷/۰۱۳	MSPE

بحث و نتیجه‌گیری

مسئله برآورد و انتخاب متغیر همزمان در مدل‌های نیمه پارامتری با اثرات آمیخته برای داده‌های طولی با بعد بالا در نظر گرفته شده است. مولفه ناپارامتری موجود در مدل با اسپلاین‌های رگرسیونی تقریب زده شده و از طریق بهینه‌سازی تابع هدف مبتنی بر تابع تاوان برآورد و انتخاب متغیر به طور همزمان انجام می‌شود. در ادامه، رفتار حدی برآوردگرهای حاصل در چارچوب داده‌های طولی با بعد بالا که در آن تعداد پارامترها متناسب با افزایش حجم نمونه افزایش می‌یابد، مورد مطالعه قرار گرفت. به منظور پیاده‌سازی روش برآورد پیشنهادی، یک الگوریتم تکراری مناسب برای انتخاب متغیرهای مهم و برآورد ضرایب غیر صفر ارائه گردیده است. در نهایت، عملکرد روش پیشنهادی با مطالعه شبیه‌سازی و تحلیل یک مجموعه داده واقعی مورد ارزیابی قرار گرفته است. نتایج نشان می‌دهند که مدل پیشنهادی در مقایسه با سایر مدل‌های رقیب، در برآورد ضرایب غیر صفر و انتخاب مدل عملکرد بهتری دارد. به عبارت دیگر وقتی داده‌ها دارای روند غیر خطی هستند، مدل‌های خطی-جزئی نسبت به مدل‌های خطی کارآمدتر هستند.

تقدیر و تشکر

نویسندگان از سردبیر، داوران و ویراستار محترم مجله که نظرات ارزشمند ایشان باعث بهبود مطالب ارائه شده در این مقاله گردید، کمال تشکر و قدردانی را دارند.

مراجع

کاظمی، م.، شاهسونی، د. و آرشی، م. (۱۳۹۷)، انتخاب متغیر و تشخیص ساختار در بعد بالا برای مدل‌های جمعی خطی-جزیی، مجله علوم آماری، ۱۲، ۴۸۵-۵۱۲.

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

Cantoni, E., Mills, F. J. and Ronchetti, E. (2005), Variable Selection for Marginal Longitudinal Generalized Linear Models, *Biometrics*, **61**, 507-514.

Fan, J. Q., Huang, T. and Li, R. (2007), Analysis of Longitudinal Data with Semiparametric Estimation of Covariance Function, *Journal of the American Statistical Association*, **102**, 632-641.

Fan, J. Q. and Li, R. (2001), Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties, *Journal of the American Statistical Association*, **96**, 1348-1360.

Fan, J. Q. and Li, R. (2004), New Estimation and Model Selection Procedures for Semiparametric Modeling in Longitudinal Data Analysis, *Journal of the American Statistical Association*, **99**, 710-723.

Fu, W. J. (2003), Penalized Estimating Equations, *Biometrics*, **59**, 126-132.

- He, X. M., Zhu, Z. Y. and Fung, W. K. (2002), Estimation in a Semiparametric Model for Longitudinal Data with Unspecified Dependence Structure, *Biometrika*, **89**, 579-590.
- Hunter, D. and Li, R. (2005), Variable Selection Using MM Algorithm. *Annals of Statistics*, **33**, 1617-1642
- Laird, N. M. and Ware, J. H. (1982), Random Effects Models for Longitudinal Data, *Biometrics*, **38**, 963-974.
- Liang, K. Y. and Zeger, S. L. (1986), Longitudinal Data Analysis Using Generalized Linear Models, *Biometrika*, **73**, 13-22.
- Ma, S., Song, Q. and Wang, L. (2013), Simultaneous Variable Selection and Estimation in Semiparametric Modeling of Longitudinal/Clustered Data, *Bernoulli*, **19**, 252–274.
- Ni, X., Zhang, D. and Zhang, H. H. (2010), Variable Selection for Semiparametric Mixed Models in Longitudinal Studies, *Biometrics*, **66**, 79-88.
- Qin, G. Y. and Zhu, Z. Y. (2007), Robust Estimation in Generalized Semiparametric Mixed Models for Longitudinal Data, *Journal of Multivariate Analysis*, **98**, 1658-1683.
- Schumaker, L. L. (1981), *Spline Functions*, New York: Wiley.
- Sinha, S. K. and Sattar, A. (2015), Inference in Semi-Parametric Spline Mixed Models for Longitudinal Data, *Metron*, **73**, 377-395.
- Taavoni, M. and Arashi, M. (2019), Regularization in Generalized Semiparametric Mixed-Effects Model for Longitudinal Data, Submitted.

Wang, L. and Qu, A. (2009), Consistent Model Selection and Data-Driven Smooth Tests for Longitudinal Data in the Estimating Equations Approach, *Journal of the Royal Statistical Society (Series B)*, **71**, 177-190.

Wang, W. L., Fan, T. H. (2011), Estimation in Multivariate t Linear Mixed Models for Multiple Longitudinal Data, *Statistica Sinica*, **21**, 1857-1880.

Wang, L., Zhou, J. and Qu, A. (2012), Penalized Generalized Estimating Equations for High-Dimensional Longitudinal Data Analysis, *Biometrics*, **68**, 353-360.

Xu, P. R., Fu, W. and Zhu, L. X. (2013), Shrinkage Estimation Analysis of Correlated Binary Data with a Diverging Number of Parameters, *Science China Mathematics*, **56**, 359–377.

Xue, L., Qu, A. and Zhou, J. (2010), Consistent Model Selection for Marginal Generalized Additive Model for Correlated Data, *Journal of the American Statistical Association*, **105**, 1518-1530.

Zeger, S. L. and Diggle, P. J. (1994), Semi-Parametric Models for Longitudinal Data with Application to CD4 Cell Numbers in HIV Seroconverters, *Biometrics*, **50**, 689-699.

Journal of Statistical Sciences, Autumn and Winter, 2020
Vol. 14, No. 2, pp 367-388
DOI: 10.29252/jss.13.2.367

Variable Selection in Semiparametric Mixed Effect Model for High-Dimension Longitudinal Data

Taavoni, M. and Arashi, M.

Department of Statistics, Shahrood University of Technology, Shahrood, Iran.

Abstract: This paper considers the problem of simultaneous variable selection and estimation in a semiparametric mixed-effects model for longitudinal data with normal errors. We approximate the nonparametric function by regression spline and simultaneously estimate and select the variables under the optimization of the penalized objective function. Under some regularity conditions, the asymptotic behaviour of the resulting estimators is established in a high-dimensional framework where the number of parametric covariates increases as the sample size increases. For practical implementation, we use an EM algorithm to select the significant variables and estimates the nonzero coefficient functions. Simulation studies are carried out to assess the performance of our proposed method, and a real data set is analyzed to illustrate the proposed procedure.

Keywords: Longitudinal data, Penalized estimator, Smoothing spline, Semiparametric mixed model, Variable selection, HIV.

Mathematics Subject Classification (2010): 62H20, 62F12.