



تحلیل حساسیت و آزمون نقاط دورافتاده در مدل های آمیخته نیمه پارامتری

هادی امامی^۱، پروانه منصوری^۲

دانشگاه زنجان

چکیده: مدل های آمیخته نیمه پارامتری کاربرد گسترده ای در مباحث اقتصاد سنجی و زیست سنجی به ویژه داده های طولی دارند. این مقاله تاثیر تحلیل حساسیت روی برآوردگرهای ماکسیمم درستنمایی تاوانیده را در مدل های حذف موردی و حذف آزمودنی بررسی می کند. هم چنین معادل بودن مدل حذف موردی و حذف آزمودنی با مدل انتقال میانگین نقاط دورافتاده از مباحث دیگر این مقاله می باشد. در ادامه با استفاده از یک مثال شبیه سازی شده کاربرد تحلیل حساسیت ارزیابی می شود. واژه های کلیدی: فاصله کوک، تابع درستنمایی تاوانیده، رگرسیون نیمه پارامتری، اسپلاین هموارسازی کد موضوع بندی ریاضی (۲۰۱۰): ۶۲۲۰، ۶۲۲۰۲.

۱ مقدمه

تحلیل حساسیت یکی از بخش های مهم تحلیل آماری به شمار می رود. آماره کوک (۱۹۷۷) یک میحث مهم تحلیل حساسیت می باشد. امروزه آماره کوک کاربرد گسترده ای در مدل های خطی پیدا کرده است و محاسبه آن با نرم افزارهای معروف آماری مانند SPSS و SAS کمک زیادی به محققان می کند. در مدل حذف موردی، آماره کوک بدون نیاز به برآورد دوباره مدل به ازای مشاهدات حذف شده محاسبه می شود. زارع و راسخ (۲۰۱۱) تحلیل حساسیت در مدل های آمیخته خطی با خطا در اندازه گیری را بررسی کرده اند. در این مقاله تاثیر تحلیل حساسیت برای برآوردگرهای ماکسیمم درستنمایی تاوانیده (MPLE) بر پایه حذف موردی و حذف آزمودنی بررسی می شود. هدف ما محاسبه مستقیم معیار حساسیت برای اثرات ثابت پارامتری و برآورد تابع ناپارامتری است که مشابه با آماره کوک و آماره DFFITS بلسلی و همکاران (۱۹۸۰) در مدل رگرسیون خطی است. در این مقاله جهت تمایز بین مفاهیم مدل رگرسیون خطی با مدل آمیخته نیمه پارامتری از DFIT به جای DFFITS استفاده می شود. آماره کوک (۱۹۷۷) به عنوان معیار حساسیت در برآورد پارامتر خطی و DFIT برای اندازه گیری تغییرات یک برازش ناپارامتری کاربرد دارد. هم چنین می توان از آن دو برای شناسایی موارد یا آزمودنی

^۱ پروانه منصوری: p.mansoori@znu.ac.ir

های موثر برای MPLE در مدل های آمیخته نیمه پارامتری استفاده کرد.

۲ مدل و برآورد پارامترها

آزمایشی با m آزمودنی و n_i مشاهده طوری که $i = 1, \dots, m$ را در نظر بگیرید. Y_{ij} متغیر پاسخ i امین آزمودنی در زمان t_{ij} به صورت زیر مدل بندی می شود:

$$Y_{ij} = \mathbf{X}_{ij}^T \beta + f(t_{ij}) + \mathbf{Z}_{ij}^T \mathbf{b}_i + U_i(t_{ij}) + \varepsilon_{ij} \quad (1.2)$$

که β یک بردار $1 \times p$ از ضرایب رگرسیونی، \mathbf{X}_{ij} ماتریس متغیرهای توصیفی، $f(t)$ تابع هموار، b_i بردار مستقل $1 \times q_i$ اثرات تصادفی با ماتریس متغیرهای تصادفی \mathbf{Z}_{ij} ، $\mathbf{U}_i(t)$ فرایند تصادفی مستقل و ε_{ij} بردار خطاهای تصادفی می باشد. بردارهای ویژه آزمودنی را به این صورت نامگذاری می کنیم: $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T$ ، به همین ترتیب \mathbf{X}_i ، \mathbf{U}_i ، \mathbf{Z}_i و ε_i نیز به صورت مشابه تعریف می شوند. فرض کنید $t^\circ = (t_1^\circ, \dots, t_r^\circ)^T$ بردار مقادیر گسسته مرتب شده نقاط زمانی t_{ij} ، $(i = 1, \dots, m, j = 1, \dots, n_i)$ ، \mathbf{N}_i ماتریس رخداد $n_i \times r$ آزمودنی i ام باشد، طوری که $t_{ij} = t_l^\circ$ اگر عنصر (j, l) ام آن برابر ۱ باشد. متغیرهای ε, Y, X, N, U بردارهایی هستند که از جمع کردن m بردار ویژه آزمودنی روی هم به دست می آیند. در تمام مقاله ماتریس ها و بردارها به صورت پیرنگ نمایش داده می شوند. بنابراین مدل ۱.۲ را به صورت ماتریس و بردار به صورت زیر بازنویسی می کنیم:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{Nf} + \mathbf{Zb} + \mathbf{U} + \varepsilon \quad (2.2)$$

خطاها (ε) دارای توزیع نرمال $(0, \sigma^2)$ ، اثر تصادفی $\mathbf{b} \sim N(0, \mathbf{D})$ ، $\mathbf{D} = \text{diag}(\mathbf{D}_1, \dots, \mathbf{D}_m)$ و \mathbf{U} دارای توزیع $(0, \Gamma)$ که $\Gamma = \text{diag}(\Gamma_1, \dots, \Gamma_m)$ است. ماتریس کواریانس \mathbf{Y} برابر است با: $\mathbf{V} = \text{diag}(\mathbf{V}_1, \dots, \mathbf{V}_m)$ $\text{cov}(\mathbf{y}) = \mathbf{V}$ طوری که

$$\mathbf{V}_i = \mathbf{Z}_i \mathbf{D}_i \mathbf{Z}_i^T + \sigma^2 \mathbf{I}_i + \Gamma_i$$

تابع درستنمایی (β, f) را در زیر در نظر بگیرید:

$$l(\beta, \mathbf{f}; \mathbf{Y}) = -\frac{1}{\nu} \log |\mathbf{V}| - \frac{1}{\nu} (\mathbf{Y} - \mathbf{X}\beta - \mathbf{Nf})^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta - \mathbf{Nf})$$

با توجه به تابع فوق، با ماکسیم کردن تابع زیر برآوردگرهای ماکسیم درستنمایی تاوانیده حاصل می شوند:

$$\mathbf{L}(\beta, \mathbf{f}) = l(\beta, \mathbf{f}; \mathbf{Y}) - \frac{\lambda}{\nu} \int f''(t) \mathbf{y} dt = l(\beta, \mathbf{f}; \mathbf{Y}) - \frac{\lambda}{\nu} \mathbf{f}^T \mathbf{K} \mathbf{f} \quad (3.2)$$

λ پارامتر هموارسازی و \mathbf{K} ماتریس هموارساز منتهای نامنفی می باشد. در این مقاله فرض می شود \mathbf{V} و λ پارامترهای معلوم هستند؛ در این صورت رابطه ۳.۲ یک تابع درجه دوم است و MPLE از رابطه خطی زیر به دست می آید:

$$\mathbf{C} \begin{pmatrix} \beta \\ \mathbf{f} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \mathbf{V}^{-1} \\ \mathbf{N}^T \mathbf{V}^{-1} \end{pmatrix} \mathbf{Y} \quad (4.2)$$

$$\mathbf{C} = \begin{pmatrix} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} & \mathbf{X}^T \mathbf{V}^{-1} \mathbf{N} \\ \mathbf{N}^T \mathbf{V}^{-1} \mathbf{X} & \mathbf{N}^T \mathbf{V}^{-1} \mathbf{N} + \lambda \mathbf{K} \end{pmatrix}$$

در تمام مقاله فرض می شود C پرتبه $p + r$ است. این مطلب این اطمینان را ایجاد می کند که ماتریس (X, NT) هم پرتبه است طوری که $T = (1, t^0)$ و 1 یک بردار یکه $1 \times r$ می باشد. طبق رابطه ۴.۲ برآورد پارامترها به صورت زیر می باشد:

$$\hat{\beta} = (X^T W_x X)^{-1} X^T W_x Y \quad (5.2)$$

$$\hat{f} = (N^T W_f N + \lambda K)^{-1} N^T W_f Y \quad (6.2)$$

$$W_x = V^{-1} - V^{-1} N (N^T V^{-1} N + \lambda K)^{-1} N^T V^{-1}$$

$$W_f = V^{-1} - V^{-1} X (X^T V^{-1} X)^{-1} X^T V^{-1}$$

هم چنین برآورد اثرات ثابت با امید ریاضی شرطی به شکل زیر بدست می آید:

$$\hat{b}_i = D_i Z_i^T V_i^{-1} (Y_i - X_i \hat{\beta} - N_i \hat{f}) \quad (7.2)$$

$$\hat{b} = (\hat{b}_1^T, \dots, \hat{b}_m^T)^T$$

با کمی محاسبه از دو رابطه ۵.۲ و ۷.۲ مقادیر برازش شده برابر $\hat{Y} = X\hat{\beta} + N\hat{f} + Z\hat{b} = HY$ با $H = I - \Sigma V^{-1} + \Sigma V^{-1} \bar{H}$ و \bar{H} برابر با:

$$\bar{H} = \begin{pmatrix} X & N \end{pmatrix} C^{-1} \begin{pmatrix} X^T \\ N^T \end{pmatrix} V^{-1} \quad (8.2)$$

می باشند. I یک ماتریس همانی است و $\Sigma = \sigma^2 I + \Gamma$. ماتریس H نقش ماتریس برازشی در رگرسیون خطی را ایفا می کند؛ اما در اینجا داده های نافذ اثرات ثابت وابسته به ماتریس \bar{H} می باشند، این وابستگی در مورد باقی مانده ها هم صدق می کند:

$$\bar{e} = Y - X\hat{\beta} - N\hat{f} = (I - \bar{H})Y \quad (9.2)$$

۳ تحلیل حذف موردی

حذف موردی مشاهدات پایه ای برای ساختن آماره های تحلیل حساسیت می باشد. در مدل های رگرسیونی خطی معمولاً از روش حذف موردی به منظور بررسی تأثیر تک مشاهدات بر روی برآورد پارامترها استفاده می شود. در این بخش تأثیر حذف یک تک مشاهده روی برآورد (β, f) بررسی می شود.

۱.۳ برآوردها تحت حذف موردی

فرض کنید $\hat{\theta}_{(ij)} = (\hat{\beta}_{(ij)}^T, \hat{f}_{(ij)}^T)^T$ برآورد θ بدون (i, j) امین مشاهده باشد. محاسبه $\hat{\theta}_{(ij)}$ به ازای تمام (i, j) ها و مقایسه آن با $\hat{\theta}$ مخصوصاً زمانی که اندازه نمونه بزرگ باشد زمان زیادی را می طلبد. با استفاده از قضیه ۱.۳ می توان برآورد پارامترها را تحت حذف مشاهده بدون بررسی کل نمونه بدست آورد. برای سادگی مشاهدات شماره گذاری می شوند؛ فرض کنید (i, j) امین مشاهده شماره $c = n_1 + \dots + n_{i-1} + j$ را به خود می گیرد. هم چنین فرض کنید d_c یک بردار $1 \times n$ می باشد طوری که در مرتبه c مقدار

یک و در بقیه جاها مقدار صفر را خواهد داشت.

قضیه ۱.۳: با استفاده از نمادگذاری های بالا نتیجه می شود که:

$$\hat{\beta}_{(ij)} = \hat{\beta} - \frac{(\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{d}_c \mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{e}}}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c} \quad (1.3)$$

$$\hat{\mathbf{f}}_{(ij)} = \hat{\mathbf{f}} - \frac{(\mathbf{N}^T \mathbf{W}_f \mathbf{N} + \lambda \mathbf{K})^{-1} \mathbf{N}^T \mathbf{W}_f \mathbf{d}_c \mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{e}}}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c} \quad (2.3)$$

برای داده های مستقل با $\mathbf{V} = \sigma^2 \mathbf{I}$ رابطه ۱.۳ به شکل زیر تغییر می کند:

$$\hat{\beta}_{(c)} = \hat{\beta} - \frac{(\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{d}_c \bar{\mathbf{e}}_c}{1 - \bar{h}_{cc}} \quad (3.3)$$

$\bar{\mathbf{e}}_c$ مولفه c ام \mathbf{e} در ۹.۲ و \bar{h}_{cc} ، c مین عنصر قطری \bar{H} در ۸.۲ می باشد. اگر جزء ناپارامتری یعنی f در مدل نباشد ۳.۳ همان فرمول کوک (۱۹۷۷) در رگرسیون خطی می باشد.

۲.۳ آماره کوک و DFIT

آماره کوک تعمیم یافته به صورت مربع آماره مقیاس شده $\hat{\theta}$ از θ به شکل زیر تعریف می شود:

$$CD_{ij}(\beta, \mathbf{f}) = (\hat{\theta} - \hat{\theta}_{(ij)})^T \mathbf{C} (\hat{\theta} - \hat{\theta}_{(ij)})$$

واضح است که هرچقدر این فاصله بزرگتر باشد احتمال اینکه مشاهده (i, j) ام یک مشاهده موثر باشد بیشتر است. این آماره را می توان با جایگذاری از قضیه ۱.۳ به صورت مستقیم به دست آورد:

$$CD_{ij}(\beta, \mathbf{f}) = \mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{H}} \mathbf{d}_c \left\{ \frac{\mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{e}}}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c} \right\}^2 = \frac{\mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{H}} \mathbf{d}_c}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c} t_c^2$$

$$t_c = \frac{\mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{e}}}{\sqrt{\{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c\}}}$$

t_c باقی مانده استیوندت شده مورد c ام می باشد. در مدل های نیمه پارامتری معمولاً تاثیر جزئی روی برآورد f ، β را به صورت مجزا بررسی می کنند. بنابراین آماره کوک برای جزء پارامتری مدل یعنی β به شکل زیر به دست خواهد آمد:

$$CD_{ij}(\beta) = (\hat{\beta} - \hat{\beta}_{(ij)})^T \{(\mathbf{I}_p, \circ) \mathbf{C}^{-1} (\mathbf{I}_p, \circ)^T\}^{-1} (\hat{\beta} - \hat{\beta}_{(ij)})$$

به بیانی ساده تر و با استفاده از قضیه ۱.۳ نتیجه می شود که:

$$CD_{ij}(\beta) = \frac{\mathbf{d}_c^T \mathbf{W}_x \mathbf{X} (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{d}_c}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c} t_c^2$$

همان طور که مشاهده می شود تاثیر جزئی مورد c ام روی $\hat{\beta}$ به مقادیر نافذ یا به اندازه t_c بستگی دارد. مقادیر نافذ جزء ناپارامتری، عناصر قطری $\mathbf{H}_\beta = \mathbf{W}_x \mathbf{X} (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x$ می باشد. پیرو کوک (۱۹۷۷) جهت معناداری مشاهدات موثر این معیار را با مقادیر توزیع $\chi_{p,\alpha}^2$ می سنجند.

$$(\hat{\beta} - \beta)^T \{(\mathbf{I}_p, \circ) \mathbf{C}^{-1} (\mathbf{I}_p, \circ)^T\}^{-1} (\hat{\beta} - \beta) \leq \chi_{p,\alpha}^2$$

که یک بیضی اطمینان در سطح $1 - \alpha$ است. از کوک (۱۹۷۷) علاقه مند به $\hat{\beta}_{(ij)}$ های هستیم که در ناحیه اطمینان 100% درصد باشد. از آن جایی که ممکن است اثر مشاهدات به دلیل هم پوشانی در یک همسایگی موضعی بر روی برآوردهای اسپلین، معیار حساسیت را کاهش دهد؛ می توان از اثر جزئی با استفاده از رابطه ۴.۳ استفاده کرد.

$$DFIT_{ij} = \frac{|d_c^T N(\hat{f}_{(ij)} - \hat{f})|}{s_{ij}} \quad (4.3)$$

طوری که s_{ij}^2 عنصر قطری c ام $N(0, I_r)C^{-1}(0, I_r)^T N^T$ می باشد.

۳.۳ مدل انتقال میانگین نقاط دورافتاده و آزمون نقاط دورافتاده

از مدل انتقال میانگین نقاط دورافتاده اغلب برای آزمون اینکه آیا مشاهده مورد بررسی یک نقطه دورافتاده برای مدل محسوب می شود استفاده می شود. به ازای مورد (i, j) ام این مدل به فرم زیر می باشد:

$$Y_{kl} = X_{kl}^T \beta + f(t_{kl}) + Z_{kl}^T b_i + U_k(t_{kl}) + \varepsilon_{kl} \quad k \neq i, l \neq j,$$

$$Y_{ij} = X_{ij}^T \beta + f(t_{ij}) + Z_{ij}^T b_i + U_i(t_{ij}) + \gamma + \varepsilon_{ij} \quad (5.3)$$

که γ یک پارامتر اضافی برای نشان دادن وجود نقطه دورافتاده می باشد. بنابراین آزمون نقاط دورافتاده را می توان به صورت آزمون فرض $\gamma = 0$ بیان کرد. فرض کنید $\hat{\beta}_{ij}$ و \hat{f}_{ij} و $\hat{\gamma}_{ij}$ ؛ MPLE مدل ۵.۳ باشد. قضیه زیر برقرار است:

قضیه ۳.۳: برآوردهای ماکسیم درستنمایی تحت مدل ۲.۲ برابر برآوردهای ماکسیم درستنمایی مدل ۵.۳ با مجموعه تمام مشاهدات می باشد. یعنی:

$$\hat{f}_{(ij)} = \hat{f}_{ij} \quad \hat{\beta}_{(ij)} = \hat{\beta}_{ij}$$

معادل بودن مدل حذف موردی و مدل نقاط دورافتاده در رگرسیون خطی را در کوک و ویزبرگ (۱۹۸۲) می توان دید. این مساله توسط وی (۱۹۹۸) در کلاس گسترده تری از مدل های پارامتری بسط داده شده است. قضیه ۳.۳ برابری این دو مدل را در مدل های آمیخته نیمه پارامتری تأیید می کند. هم چنین MPLE برای γ برابر است با:

$$\hat{\gamma}_{ij} = (d_c^T V^{-1} d_c)^{-1} d_c^T V^{-1} (Y - X \hat{\beta}_{(ij)} - N \hat{f}_{ij}) = \frac{d_c^T V^{-1} \bar{e}}{d_c^T V^{-1} (I - \bar{H}) d_c}$$

و واریانس آن تحت فرض صفر برابر است با:

$$var(\hat{\gamma}_{ij}) = \frac{d_c^T V^{-1} (I - \bar{H}) V (I - \bar{H}) V^{-1} d_c}{\{d_c^T V^{-1} (I - \bar{H}) d_c\}^2}$$

که مجانباً برابر با واریانس بیزی آن یعنی:

$$var_B(\hat{\gamma}_{ij}) = \frac{1}{d_c^T V^{-1} (I - \bar{H}) d_c}$$

می باشد. مقدار استاندارد شده $\hat{\gamma}_{ij}$ با واریانس بیزی آن برابر با باقی مانده های استیودنت شده می باشد و از آن می توان برای بررسی این که آیا مشاهده (i, j) ام در مدل یک نقطه پرت محسوب می شود استفاده کرد. همانطور که در مدل رگرسیون خطی ساده از باقی مانده های مقیاس شده برای شناسایی نقاط دورافتاده استفاده می شود؛ در این مدل از نمودار t_c که نمودار باقی مانده های استیودنت شده می باشد برای غربالگری نقاط دور افتاده می توان استفاده نمود.

۴ تحلیل حذف آزمودنی

در مطالعات طولی مشاهدات آزمودنی های مشابه اغلب مقادیر مشابهی در متغیرهای کمکی می گیرند. بنابراین شناسایی آزمودنی های موثر در مقابل مشاهدات موثر می تواند مفید واقع شود. در این بخش تاثیر حذف آزمودنی i ام روی $\hat{\beta}$ و \hat{f} در مدل بررسی می شود. فرض کنید $\hat{\theta}_{[i]} = (\hat{\beta}_{[i]}^T, \hat{f}_{[i]}^T)^T$ برآورد β و f با حذف آزمودنی i ام باشد. قضیه زیر فرمولی برای برآورد ارائه می کند.

قضیه ۴: فرض کنید $E_i = (\circ, \dots, \circ, I_{n_i}, \circ, \dots, \circ)^T$ یک ماتریس $n \times n_i$ باشد. بنابراین:

$$\hat{\theta}_{[i]} = \hat{\theta} - C^{-1} \begin{pmatrix} X^T \\ N^T \end{pmatrix} V^{-1} E_i (I_{n_i} - \bar{H}_i)^{-1} E_i^T \bar{e}$$

طوری که

$$\bar{H}_i = \begin{pmatrix} X_i & N_i \end{pmatrix} C^{-1} \begin{pmatrix} X_i^T \\ N_i^T \end{pmatrix} V_i^{-1}$$

آماره کوک تعمیم یافته تحت مدل حذف آزمودنی به شکل زیر تعریف می شود:

$$CD_{[i]}(\beta, f) = (\hat{\theta}_{[i]} - \theta)^T C (\hat{\theta}_{[i]} - \theta)$$

با استفاده از قضیه ۴ فرمول را می توان بازنویسی کرد:

$$CD_{[i]}(\beta, f) = \bar{e}^T E_i (I_{n_i} - \bar{H}_i^T)^{-1} V_i^{-1} \bar{H}_i (I_{n_i} - \bar{H}_i)^{-1} E_i^T \bar{e}$$

که $E_i^T \bar{e}$ بردار باقی مانده $n_i \times 1$ متناظر با آزمودنی i ام می باشد. همانند بخش قبل برای شناسایی تاثیر آزمودنی i ام روی برآورد β باید $CD_{[i]}(\beta)$ را با توزیع $\chi_{p,\alpha}^2$ ارزیابی کرد. از آنجایی که حذف آزمودنی i ام شامل n_i نقطه زمانی می باشد؛ به محاسبه $DFIT$ در n_i نقطه زمانی برای شناسایی تاثیر موضعی روی منحنی برازش شده f نیاز هست. هم چنین از مدل انتقال میانگین نقاط دورافتاده برای i امین مورد می توان استفاده کرد:

$$Y_i = X_i \beta + N_i f + Z_i b_i + \Delta + U_i + \varepsilon_i$$

$$Y_j = X_j \beta + N_j f + Z_j b_j + U_j + \varepsilon_j \quad j \neq i \quad (1.4)$$

که در آن Δ یک بردار $n_i \times 1$ و نشان دهنده آزمودنی های دورافتاده می باشد. فرض کنید $\bar{\beta}_{[i]}$ ، $\bar{f}_{[i]}$ و $\bar{\Delta}_{[i]}$ نمایانگر MPLE برای β ، f و Δ مدل ۱.۴ باشد. مشابه قضیه ۳.۳ می توانیم نشان دهیم:

$$\bar{\beta}_{[i]} = \hat{\beta}_{[i]} \quad \bar{f}_{[i]} = \hat{f}_{[i]} \quad \bar{\Delta}_{[i]} = R_i$$

۵ تحلیل داده شبیه سازی شده

برای ارزیابی معیار های تحلیل حساسیت در این مقاله داده شبیه سازی شده با دو نقطه دور افتاده را بررسی می کنیم. متغیر کمکی X_i اندازه های غیر نرمال BMI از ۲۰ آزمودنی اول داده های پروژسترون سوئز و همکاران (۱۹۹۸) می باشد. مقادیر آن از ۲۱,۵۳۲۸ تا

۳۸,۰۱۶۵ مرتب شده اند. مدل به صورت زیر است:

$$Y_{ij} = 0.5(X_i - 26) + f(t_{ij}) + b_i + U_i(t_{ij}) + \varepsilon_{ij} \quad i = 1, \dots, 20 \quad j = 1, \dots, 4$$

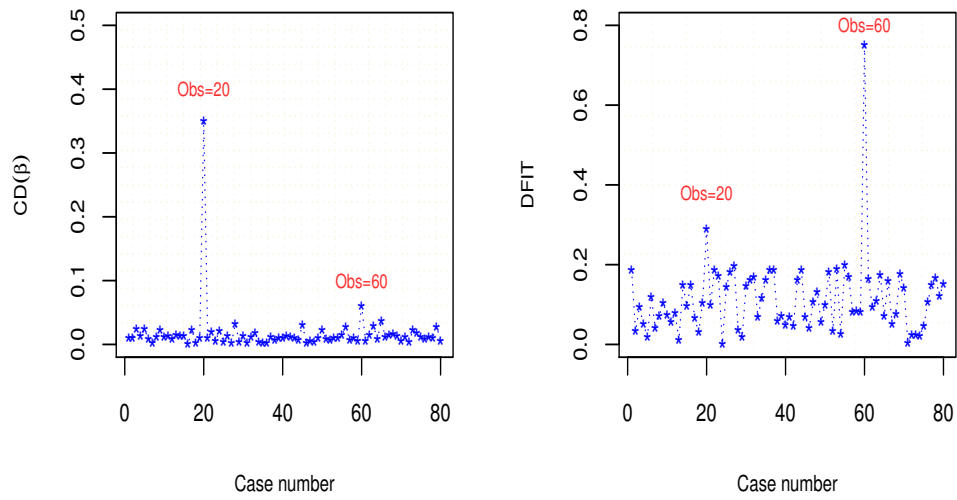
که در آن $U_i(t)$, $b_i \sim N(0, 0.12)$, $t_{i,j+1} = t_{ij} + 0.25$ ($j = 2, 3, 4$), $t_{i1} \sim U(0, 0.25)$, $f(t) = (t - 0.5)^2$ فرآیند اتورگرسیو مرتبه اول با واریانس 0.22 و خودهمبستگی مرتبه اول 0.3 و $\varepsilon_{ij} \sim N(0, 0.1)$ می باشد. آزمودنی ها از هم مستقل هستند. متغیر پاسخ Y_{ij} از مقدار -2.8 تا 6.9 مرتب شده است. الگوریتم MPLE این نتایج را به دنبال دارد: $\hat{\beta} = 0.452$ و $\lambda = 0$ شکل ۱، $CD(\beta)$ و DFIT برای برآورد f را نشان می دهد. مورد 2 اثر زیادی روی برآورد β و مورد 6 اثر زیادی روی برآورد f دارد. از آنجایی که مورد 6 نقطه نفوذ ضعیفی محسوب می شود اثر آن روی برآورد β خیلی بزرگ نیست. با توجه به بحث حذف مشاهدات، همان طوری که در شکل ۲ مشاهده می کنیم آزمودنی پنجم تاثیر زیادی روی برآورد β دارد. هم چنین با توجه به شکل ۲ آزمودنی ۵ و ۱۵ به عنوان نقاط دورافتاده شناسایی شده اند. در این مثال دو نقطه دورافتاده با اندازه های مباحث تشخیصی که بحث شد شناسایی شدند. این نتایج کاملا رضایت بخش می باشد.

بحث و نتیجه گیری

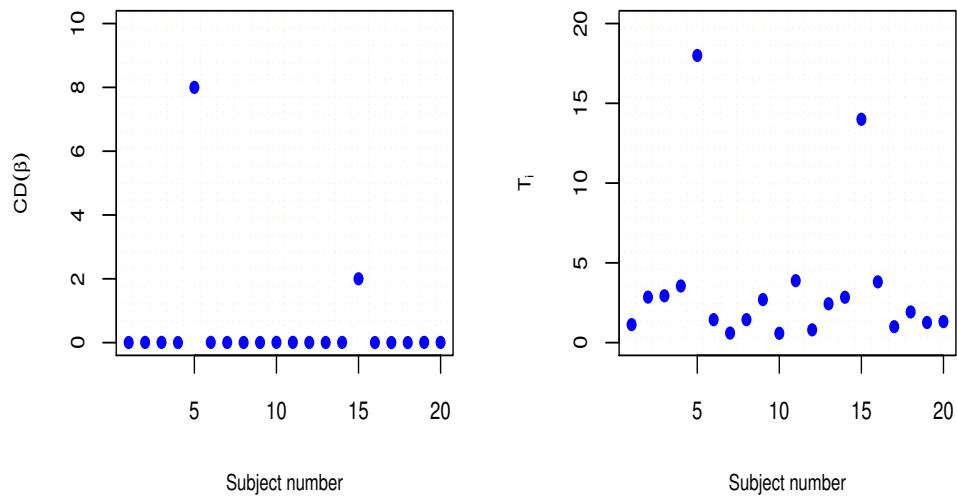
در این مقاله مدل های حذف موردی و حذف آزمودنی را در مدل های آمیخته نیمه پارامتری بررسی کردیم که بسیار با اهمیت هستند و نقش به سزایی در تحلیل داده ها دارند.

مراجع

- Belsley D. A. , Kuh, E. and Welsh, R. E. (1980) *Regression Diagnostics*, Wiley, New York.
- Cook R. D. (1977) Detection of influential observations in linear regression. *Technometrics* 19, 15-18.
- Cook R. D. and Weisberg S. (1982) , *Residual and influence in Regression*, Chapman & Hall, New York.
- Sowers, M. F., Crutchfield, M., Randolph, J. F., Shiapo, B., Zhang, B. , Pietra, M.L., Schork, M.A. (1998), Urinary ovarian and gonadotrophin hormone levels in premenopausal women with low bone mass. *J. Bone Min. statist. Res.*, 13, 1191-1202.
- Wei B. C. (1998) *Exponential Family Nonlinear Models*, Springer, Singapore.
- Zare K. and Rasekh A. (2011) , Diagnostics measure for linear mixed measurement error models. *Sort* 35, 13, 125-144.



شکل ۱: تاثیر حذف موردی روی داده های شبیه سازی شده



شکل ۲: تاثیر حذف آزمودنی روی داده های شبیه سازی شده