

مباحث تشخیصی در مدل‌های خطی آمیخته نیمه پارامتری با خطای اندازه‌گیری

هادی امامی، پروانه منصوری القناب

گروه آمار، دانشگاه زنجان

تاریخ دریافت: ۱۳۹۵/۱۱/۱۲ تاریخ آخرین بازنگری: ۱۳۹۶/۱۱/۳۰

چکیده: تمام مشاهدات نقش یکسان در مدل‌های آماری ندارند. گاهی برخی از مشاهدات اثرات نامناسبی روی نتایج تحلیل رگرسیونی دارند. بنابراین شناسایی چنین مشاهداتی در تحلیل داده‌ها از اهمیت ویژه‌ای برخوردار است. برای شناسایی چنین مشاهداتی از روش‌های تشخیصی استفاده می‌شود. در مقاله حاضر با استفاده از روش حذف موردی و مدل انتقال میانگین نقاط دورافتاده، مباحث تشخیصی در مدل خطی آمیخته نیمه پارامتری با خطا در اندازه‌گیری مورد بررسی قرار گرفته است. علاوه بر مباحث حذف موردی، مباحث حذف آزمودنی نیز ارائه شده است. در پایان عملکرد مباحث تشخیصی با استفاده از مجموعه داده‌های واقعی و یک مثال شبیه‌سازی نشان داده شده است. **واژه‌های کلیدی:** مدل خطی آمیخته، خطا در اندازه‌گیری، فاصله کوک، روش‌های امتیاز تصحیح شده، حذف موردی.

۱ مقدمه

مباحث تشخیصی در یک تحلیل آماری نقش بسزایی دارند. با حذف یک مشاهده موثر از مجموعه داده‌ها تفاوت معنی‌داری در تحلیل داده‌ها به وجود می‌آید. حذف موردی مشاهدات پایه‌ای برای ساختن آماره‌های تشخیصی است. روش‌های تشخیصی در بسیاری از مدل‌های مختلف مطرح شده است از جمله و مدل‌های خطی آمیخته تعمیم داده شده است. بکمن و همکاران (۱۹۸۷)، لیسافر و وریک (۱۹۹۸) اندازه‌های تاثیر موضعی را در این مدل‌ها مطرح کردند. کریستنسن و همکاران (۱۹۹۲) و بنرجی و فریز (۱۹۹۷) به

ترتیب تحلیل حذف موردی^۱ و حذف آزمودنی^۲ را مورد مطالعه قرار دادند. اخیراً روش حذف موردی توسط راسخ و امامی (۲۰۱۱) و امامی و امامی (۲۰۱۶) در مدل‌های خطی مقید با خطاهای همبسته تعمیم داده شده است. همچنین امامی و امامی (۲۰۱۴) روش اندازه تأثیر موضعی را در چنین مدل‌هایی مورد بررسی قرار داده‌اند. در مدل‌های رگرسیون ناپارامتری مباحث تشخیصی کمتر مورد توجه قرار گرفته‌اند. در میان مطالعات انجام شده، کیم (۱۹۹۶)، کیم و همکاران (۲۰۰۲) باقی‌مانده‌ها، داده‌های نافذ و آماره کوک را در مدل‌های رگرسیونی ناپارامتری و نیمه‌پارامتری مورد کنکاش قرار داده‌اند. مشابه پژوهش کیم و همکاران (۲۰۰۲)، فانگ و همکاران (۲۰۰۲) چنین مباحثی را در مدل‌های خطی آمیخته نیمه‌پارامتری تعمیم داده‌اند. ژانگ و همکاران (۱۹۹۸) این موضوع را مطرح کردند که مدل‌های آمیخته نیمه‌پارامتری اغلب در مطالعات طولی به دلیل ترکیب تأثیر زمان ناپارامتری علاوه بر رگرسیون خطی روی متغیرهای کمکی مفید هستند. از طرفی با توجه به پژوهش‌های دیویدیان و گیلتینان (۱۹۹۵) و لیانگ و لی (۲۰۰۹) متغیرهای مستقل در مدل‌های خطی پارامتری و نیمه‌پارامتری ممکن است با یک خطای غیر قابل اغماضی اندازه‌گیری شده باشند. عدم توجه به این مساله باعث بروز نتایج نادرست در یک تحقیق خواهد شد. بنابراین در چنین مواقعی در نظر گرفتن مدل‌های خطی آمیخته نیمه‌پارامتری با خطا در اندازه‌گیری می‌تواند مطلوب واقع شود. یکی از روش‌های متداول در برآورد پارامترهای مدل خطی با خطا در اندازه‌گیری استفاده از روش امتیازی تصحیح شده ناکامورا (۱۹۹۰) است. اخیراً زارع و همکاران (۲۰۱۱)، قپانی و همکاران (۲۰۱۶) و امامی (۲۰۱۵) چنین روشی را به ترتیب در مدل‌های آمیخته خطی و مدل‌های رگرسیون ریدج با خطا در اندازه‌گیری تعمیم داده‌اند. با این وجود در زمینه مباحث تشخیصی در مدل‌های خطی آمیخته نیمه‌پارامتری با خطا در اندازه‌گیری مطالعات چندانی در دسترس نیست. لذا در این مقاله با توجه به پژوهش‌های صورت گرفته توسط زارع و راسخ (۲۰۱۱)، مدل حذف موردی و مدل انتقال میانگین نقاط دورافتاده^۳ در مدل‌های خطی آمیخته نیمه‌پارامتری با خطای اندازه‌گیری با استفاده از تابع امتیازی تصحیح شده ناکامورا^۴ مورد بررسی و کنکاش قرار می‌گیرد. برای کشف نقاط دورافتاده از آماره آزمون امتیازی اصلاح شده بر پایه مدل‌های انتقال میانگین استفاده خواهد شد. علاوه بر این چندین تحلیل حذف موردی و حذف آزمودنی به عنوان ابزاری برای تحلیل مباحث تشخیصی صورت خواهد گرفت.

^۱Case deletion^۲Subject deletion^۳Mean shift outlier model^۴Nakamura corrected score function

۲ تعریف مدل

شکل ماتریسی مدل خطی آمیخته نیمه پارامتری با خطای اندازه گیری به صورت

$$Y = X_*\beta + Zb + f(t) + \varepsilon, \quad X = X_* + \Delta \quad (۱)$$

تعریف می‌شود، که در آن $Y = (Y_1, \dots, Y_m)^T$ طوری Y_i بردار $n \times 1$ از مشاهدات در i مین آزمودنی (خوشه) است. بردار β یک بردار پارامتری $p \times 1$ از ضرایب رگرسیونی، $X_* = Z$ و $Z = [Z_1^T, \dots, Z_m^T]^T$ ماتریس‌های پیش‌بین به ترتیب با ابعاد $n \times p$ و $n \times q$ هستند. Z_i ها ماتریس طرح معلوم و $n \times q_i$ از اثرات تصادفی هستند. $b^T = (b_1^T, \dots, b_m^T)$ و هرکدام از b_i ها بردار تصادفی نامشخص از اثرات تصادفی با توزیع $N(0, \sigma_i^2 I)$ که $i = 1, \dots, m$ است. $f(\cdot)$ یک تابع نامعلوم هموار و دوبار مشتق پذیر روی بازه متناهی، $t = (t_1, \dots, t_n)$ که t_i ها اسکالر و $a \leq t_1, \dots, t_n \leq b$ هموار و دوبار مشتق پذیر روی بازه متناهی، $\varepsilon = (\varepsilon_1^T, \dots, \varepsilon_m^T)$ بردار $n \times 1$ از خطاهای تصادفی با توزیع $N(0, \sigma^2 I)$ و X مقدار مشاهده شده X_* با خطای اندازه‌گیری Δ است، Δ ماتریس تصادفی با توزیع $N(0, I \otimes \Lambda)$ و Λ نیز یک ماتریس معین مثبت است. در مدل (۱) فرض می‌شود b_i ، ε و Δ دوجه دو از هم مستقلند و بردار تصادفی b دارای توزیع $N(0, \sigma^2 D)$ است. D ماتریس قطریست که i امین بلوک قطری آن $\gamma_i I$ با $\gamma_i = \sigma_i^2 / \sigma^2$ مشخص می‌شود. بنابراین توزیع Y نرمال $N(X_*\beta + f(t), \sigma^2 V)$ با $V = I + \sum_{i=1}^m \gamma_i Z_i Z_i^T$ است. فرض کنید مقادیر گسسته مرتب شده $t_1 \dots t_n$ را با $s_1 \dots s_q$ نشان داده شود. ارتباط بین $t_1 \dots t_n$ با $s_1 \dots s_q$ را ماتریس مجاورت N که $n \times q$ هست مشخص می‌کند؛ طوری که $N_{ij} = 1$ اگر $s_j = t_i$ ؛ در غیر این صورت برابر صفر است. فرض کنید f برداری است که مقدار $a_i = f(s_j)$ را شامل شود. بنابراین مدل (۱) را می‌توان به صورت

$$Y = X\beta + Zb + Nf + \varepsilon \quad (۲)$$

بازنویسی کرد. تابع لگاریتم درستنمایی تاوانیده شده برای مدل (۲) به صورت

$$\begin{aligned} \ell(\beta, b, f; X_*, Y) = & C_{\sigma^2} - \frac{1}{2\sigma^2} (Y - X_*\beta - Zb - Nf)^T (Y - X_*\beta - Zb - Nf) \\ & - \frac{1}{2\sigma^2} b^T D^{-1} b - \frac{\lambda}{2\sigma^2} \int f''(t)^2 dt, \end{aligned} \quad (۳)$$

۲۲۲ مباحث تشخیصی در مدل‌های خطی آمیخته نیمه پارامتری

بدست می‌آید، که در آن $C_{\sigma^2} = -\frac{1}{2} \log(2\pi\sigma^2)^{n+q} - \frac{1}{2} \log |D|$ طبق روش هارویل (۱۹۷۷) با حل معادله $\frac{\partial \ell(\beta, b, f; X_*, Y)}{\partial b} = 0$ برآورد b به صورت

$$\tilde{b}_{\beta, f}(X_*) = (Z^T Z + D^{-1})^{-1} Z^T (Y - X_* \beta - N f)$$

بدست می‌آید. اکنون با جایگذاری برآورد اخیر در معادله (۳) داریم:

$$\ell_p(\beta, f; X_*, Y) = C_{\sigma^2} - \frac{1}{2\sigma^2} (Y - X_* \beta - N f)^T V^{-1} (Y - X_* \beta - N f) - \frac{\lambda}{2\sigma^2} f^T K f,$$

که در آن λ پارامتر هموارسازی است و با مینیم کردن معیار اعتبارسنجی متقابل تعمیم‌یافته تعیین می‌شود و K ماتریس هموارساز متناهی نامنفی است.

۳ برآورد پارامترها

همان‌طور که بیان شد X_* با خطا اندازه‌گیری شده است. با توجه به ناکامورا (۱۹۹۹) اگر X_* با X جایگزین شود برآوردهایی که از توابع امتیاز به‌دست می‌آید سازگار نیستند. لذا باید خطای اندازه‌گیری در مدل تصحیح شود. از جمله روش‌های معمول جهت تصحیح خطا در مدل‌های با خطا در اندازه‌گیری استفاده از روش ناکامورا (۱۹۹۹) است. در این روش از تابع امتیاز تصحیح شده‌ای استفاده می‌شود که امیدریاضی آن نسبت به توزیع با خطا در اندازه‌گیری با تابع امتیاز بر حسب متغیر مستقل یعنی X_* برابر باشد. باید روابط

$$\begin{aligned} E^*[\partial \ell^*(\beta, b, f; X, y)/\partial b] &= \partial \ell(\beta, f; X_*, y)/\partial b. \\ E^*[\partial \ell_p^*(\beta, f; X, y)/\partial \beta] &= \partial \ell(\beta, f; X_*, y)/\partial \beta. \end{aligned} \quad (۴)$$

برای تابع درستنمایی توانیده تصحیح شده $\ell^*(\beta, b, f; X)$ برقرار باشند. با Λ معلوم $\ell^*(\beta, b, f; X, Y)$ به صورت

$$\ell^*(\beta, b, f; X, Y) = C_{\sigma^2} - \frac{1}{2\sigma^2} \{(Y - X\beta - Zb - Nf)^T (Y - X\beta - Zb - Nf)\}$$

$$-tr(V^{-1})\beta^T \Lambda \beta\} - \frac{1}{2\sigma^2} b^T D^{-1} b - \frac{\lambda}{2\sigma^2} f^T K f$$

تغییر می‌کند. بنابراین برآورد b با حل معادله $\frac{\partial \ell^*(\beta, b, f; X, Y)}{\partial b} = 0$ به صورت

$$\begin{aligned} \tilde{b}_{\beta, f}(X) &= (Z^T Z + D^{-1})^{-1} Z^T (Y - X\beta - Nf) \\ &= DZ^T V^{-1} (Y - X\beta - Nf) \end{aligned}$$

به دست می‌آید. با وارد کردن $\tilde{b}_{\beta, f}(X)$ در تابع درستنمایی ℓ^* تابع درستنمایی

$$\begin{aligned} \ell_p^*(\beta, f; X, Y) &= \ell^*(\beta, \tilde{b}_{\beta, f}(X), f; X, Y) \\ &= C_{\sigma^2} - \frac{1}{2\sigma^2} (Y - X\beta - Nf)^T V^{-1} (Y - X\beta - Nf) \\ &\quad + \frac{1}{2\sigma^2} tr(V^{-1})\beta^T \Lambda \beta - \frac{\lambda}{2\sigma^2} f^T K f. \end{aligned} \quad (5)$$

حاصل می‌شود. اکنون با مشتق‌گیری از (۵) نسبت به β و f و σ^2 برآوردهای ماکسیمم درستنمایی توانیده تصحیح شده^۵ (MCPL) پارامترها به صورت

$$\hat{\beta} = (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W Y$$

$$\hat{f} = (N^T W_* N + \lambda K)^{-1} N^T W_* Y$$

$$\hat{\sigma}^2 = \frac{1}{n+q} [(Y - X\hat{\beta} - N\hat{f})^T V^{-1} (Y - X\hat{\beta} - N\hat{f}) - tr(V^{-1})\hat{\beta}^T \Lambda \hat{\beta} - \lambda \hat{f}^T K \hat{f}]$$

به دست می‌آیند، که در آن

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

$$W_* = V^{-1} - V^{-1} X (X^T V^{-1} X - tr(V^{-1})\Lambda)^{-1} X^T V^{-1}.$$

⁵Maximum corrected penalized likelihood estimators

در ادامه برآورد پارامتر مرتبط با اثرات تصادفی مدل به صورت

$$\begin{aligned}\hat{b} &= (Z^T Z + D^{-1})^{-1} Z^T (Y - X\hat{\beta} - N\hat{f}) \\ &= DZ^T V^{-1} (Y - X\hat{\beta} - N\hat{f})\end{aligned}\quad (۶)$$

حاصل می شود. مقادیر برازش شده مدل را نیز می توان به صورت

$$\hat{Y} = X\hat{\beta} + N\hat{f} + Z\hat{b} = HY$$

نوشت، که در آن $H = I - V^{-1} + V^{-1}\bar{H}$ و

$$\bar{H} = \begin{bmatrix} X & N \end{bmatrix} \begin{bmatrix} X^T V^{-1} X - tr(V^{-1})\Lambda & X^T V^{-1} N \\ N^T V^{-1} X & N^T V^{-1} N + \lambda K \end{bmatrix}^{-1} \begin{bmatrix} X^T \\ N^T \end{bmatrix} V^{-1}.\quad (۷)$$

هم چنین بردار باقی مانده مدل از $\tilde{e} = Y - X\hat{\beta} - N\hat{f} = (I - \bar{H})Y$ بدست می آید.

۴ مباحث تشخیصی حذف موردی

یک مشاهده که به طور قابل توجهی متفاوت تر از سایر مشاهدات است می تواند تفاوت زیادی در نتایج تحلیل رگرسیونی ایجاد کند. مباحث تشخیصی تکنیکی برای تشخیص این مشاهدات است. در مباحث تشخیصی از دو مدل متداول حذف موردی و مدل انتقال میانگین نقاط دورافتاده استفاده می شود.

۱.۴ مدل حذف موردی

حذف موردی مشاهدات پایه ای برای ساختن آماره های تشخیصی است. در مدل های رگرسیونی خطی معمولاً از روش حذف موردی به منظور بررسی تأثیر تک مشاهدات روی برآورد پارامترها استفاده می شود. قضیه زیر روشی برای برآورد پارامترها در مدل حذف موردی ارائه می کند.

قضیه ۱: برای مدل (۲) روابط

$$\hat{\beta}_{(ij)} = \hat{\beta} - \frac{(X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W d_c d_c^T V^{-1} \tilde{e}}{d_c^T V^{-1} (I - \bar{H}) d_c}$$

$$\hat{f}_{(ij)} = \hat{f} - \frac{(N^T W_* N + \lambda K)^{-1} N^T W_* d_c d_c^T V^{-1} \tilde{e}}{d_c^T V^{-1} (I - \bar{H}) d_c}$$

برقرار است، که در آن (i, j) امین مشاهده شماره $j + n_{i-1} + \dots + n_1 + c$ را شامل می شود. d_c یک بردار $1 \times n$ است که مولفه c ام آن برابر یک و بقیه مولفه ها صفر هستند.

برهان: فرض کنید $c = 1$ ، ماتریس افراز شده

$$V = \begin{bmatrix} v_{cc} & \mathbf{v}_c^T \\ \mathbf{v}_c & V_{(c)} \end{bmatrix}$$

را در نظر بگیرد و فرض کنید Y_c ، N_c و X_c ردیف c ام ماتریس های Y ، N و X باشند و $Y_{(c)}$ ، $N_{(c)}$ و $X_{(c)}$ ماتریس هایی باشند که ردیف c ام آن ها حذف شده است. اکنون با تعریف روابط

$$\begin{aligned} Y_i^* &= (Y_{i1}, \dots, Y_{ij-1}, Y_{ij}^*, Y_{ij+1}, \dots, Y_{in_i})^T \\ Y^* &= (Y_1^T, \dots, Y_{i-1}^T, Y_i^{*T}, Y_{i+1}^T, \dots, Y_m^T)^T \\ Y_{ij}^* &= X_{ij}^T + \hat{f}_{(ij)}(t_{ij}) + \mathbf{v}_c^T V_{(c)}^{-1} (Y_{(c)} - X_{(c)} \hat{\beta}_{(ij)} - N_{(c)} \hat{f}_{(ij)}) \end{aligned}$$

با توجه به فانگ و همکاران (۲۰۰۲) برای پارامتر β و منحنی هموارساز f خواهیم داشت:

$$\begin{aligned} \hat{\beta}_{(ij)} &= (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W Y^* \\ &= \hat{\beta} - (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W d_c (Y_{ij} - Y_{ij}^*), \\ \hat{f}_{(ij)} &= (N^T W_* N + \lambda K)^{-1} N^T W_* Y^* \\ &= \hat{f} - (N^T W_* N + \lambda K)^{-1} N^T W_* d_c (Y_{ij} - Y_{ij}^*). \end{aligned}$$

از طرفی با توجه به تعاریف

$$\mathbf{Y}_{(c)} = Y_{ij}, \quad \mathbf{N}_{(c)} \mathbf{f} = f(t_{ij}), \quad \mathbf{d}_c^T \mathbf{V}^{-1} = s_c^{-1}(\mathbf{1}, -\mathbf{v}_c^T \mathbf{V}_{(c)}) \quad (۸)$$

و $\mathbf{X}\hat{\boldsymbol{\beta}}_{(ij)} + \mathbf{N}\hat{\mathbf{f}}_{(ij)} = \bar{\mathbf{H}}\mathbf{Y}^*$ داریم

$$\begin{aligned} Y_{ij} - Y_{ij}^* &= s_c \mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{(ij)} - \mathbf{N}\hat{\mathbf{f}}_{(ij)}) \\ &= s_c \mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{Y} - \bar{\mathbf{H}}\mathbf{Y}^*) \\ &= s_c \mathbf{d}_c^T \mathbf{V}^{-1} \{ \mathbf{Y} - \bar{\mathbf{H}}\mathbf{Y} + \bar{\mathbf{H}}(\mathbf{Y} - \mathbf{Y}^*) \} \\ &= s_c \mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{Y} - \bar{\mathbf{H}}\mathbf{Y}) + s_c \mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{H}} \mathbf{d}_c (Y_{ij} - Y_{ij}^*). \end{aligned}$$

که در آن

$$Y_{ij} - Y_{ij}^* = \frac{s_c \mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{Y} - \bar{\mathbf{H}}\mathbf{Y})}{\mathbf{1} - s_c \mathbf{d}_c^T \mathbf{V}^{-1} \bar{\mathbf{H}} \mathbf{d}_c} = \frac{\mathbf{d}_c^T \mathbf{V}^{-1} \tilde{\mathbf{e}}}{\mathbf{d}_c^T \mathbf{V}^{-1} (\mathbf{I} - \bar{\mathbf{H}}) \mathbf{d}_c}.$$

۲.۴ مدل انتقال میانگین نقاط دورافتاده

با توجه به این‌که داده‌های دورافتاده در تمام مراحل مربوط به تحلیل و تفسیر اطلاعات تاثیرگذار هستند در این بخش مدل تشخیصی انتقال میانگین نقاط دورافتاده مورد بررسی قرار می‌گیرد. مدل انتقال میانگین نقاط دورافتاده برای مدل (۲) به صورت

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \mathbf{N}\mathbf{f} + \boldsymbol{\phi} + \boldsymbol{\varepsilon}_{ij} \quad (۹)$$

تعریف می‌شود، که در آن $\boldsymbol{\phi}$ نمادی برای نشان دادن نقطه دورافتاده در مدل است. از آزمون فرض $\boldsymbol{\phi} = \mathbf{0}$ می‌توان به عنوان آزمون نقاط دورافتاده استفاده کرد. برای مدل (۹) برآوردهای درست‌نمایی تاوانیده تصحیح شده پارامترهای $\boldsymbol{\beta}$ ، $\boldsymbol{\phi}$ و \mathbf{f} را با علائم $\hat{\boldsymbol{\beta}}_{mij}$ ، $\hat{\boldsymbol{\phi}}_{mij}$ و $\hat{\mathbf{f}}_{mij}$ نشان داده می‌شود.

قضیه ۲: برای مدل انتقال میانگین نقاط دورافتاده داریم:

$$\hat{\beta}_{ij} = \hat{\beta}_{mij} \quad \hat{f}_{ij} = \hat{f}_{mij}$$

برهان: با توجه به فرض های قضیه ۱ می توان گفت $(\hat{\beta}_{(ij)}, \hat{f}_{(ij)})$ تابع درستنمایی تاوانیده تصحیح شده

$$L_{ij}(\beta, f) = (Y_{(c)} - X_{(c)}\beta - N_{(c)}f)^T V_{(c)}^{-1} (Y_{(c)} - X_{(c)}\beta - N_{(c)}f) + \lambda \left(\int f''(t) dt - tr(V^{-1})\beta^T \Lambda \beta \right)$$

را مینیم می سازد و $(\hat{\beta}_{(mij)}, \hat{f}_{(mij)}, \hat{\phi})$ تابع

$$L_{mij}(\beta, f, \phi) = (y - X\beta - Nf - \phi d_c)^T V^{-1} (y - X\beta - Nf - \phi d_c) + \lambda \left(\int f''(t) dt - tr(V^{-1})\beta^T \Lambda \beta \right)$$

را مینیم می کند. با مشتق گیری از تابع درستنمایی $L_{mij}(\beta, f, \phi)$ نسبت به ϕ رابطه

$$\hat{\phi}_{mij}(\beta, f) = (d_c^T V^{-1} d_c)^{-1} d_c^T V^{-1} (Y - X\beta - Nf)$$

می شود که با جایگذاری برآورد ϕ در L_{mij} خواهیم داشت:

$$L_{ij} = (\beta, f) = L_{mij}\{\beta, f, \hat{\phi}_{mij}(\beta, f)\}.$$

قضیه ۳: آماره آزمون امتیازی برای آزمون فرض $\phi = 0$ در مدل انتقال میانگین نقاط دورافتاده عبارت است از:

$$SC_{ij} = \frac{(V^{-1} d_c d_c^T V^{-1}) \tilde{e}}{S_c}$$

۲۲۸ مباحث تشخیصی در مدل‌های خطی آمیخته نیمه پارامتری

برهان: ابتدا برای محاسبه ماتریس اطلاع فیشر تصحیح شده به روش تابع درستنمایی تصحیح شده داریم:

$$L_{mij}(\beta, f, \phi) = (Y - X\beta - Nf - \phi d_c)^T V^{-1} (Y - X\beta - Nf - \phi d_c) + \lambda \left(\int f''(t)^\top dt \right) - \text{tr}(V^{-1}) \beta^T \Lambda \beta$$

بنابراین ماتریس اطلاع فیشر تصحیح شده به صورت

$$J(\beta, \phi) = \begin{bmatrix} X^T V^{-1} X - \text{tr}(V^{-1}) \Lambda & X^T V^{-1} d_c \\ d_c^T V^{-1} X & d_c^T V^{-1} d_c \end{bmatrix}$$

خواهد بود. با توجه به کوک و ویزبرگ (۱۹۸۲) آماره آزمون امتیازی تحت فرض $\phi = 0$ به صورت

$$SC_{ij} = \left[\frac{\partial}{\partial \phi} L_{mij}(\beta, f, \phi) \right]^T J^{\phi\phi} \left[\frac{\partial}{\partial \phi} L_{mij}(\beta, f, \phi) \right] \Big|_{(\hat{\beta}, \hat{f})}$$

محاسبه می‌شود که در آن $J^{\phi\phi}$ گوشه سمت راست پایینی ماتریس $J^{-1}(\beta, \phi)$ است. در نتیجه آماره آزمون امتیازی برای مدل (۹) عبارت است از:

$$SC_{ij} = (Y - X\hat{\beta} - N\hat{f})^T (d_c^T V^{-1} d_c)^{-1} (Y - X\hat{\beta} - N\hat{f})$$

با استفاده از رابطه $S_c = (d_c^T V^{-1} d_c)^{-1}$ و $\tilde{e} = Y - X\hat{\beta} - N\hat{f}$ نتیجه‌ی مورد نظر حاصل می‌شود یعنی:

$$SC_{ij} = \frac{(V^{-1} d_c d_c^T V^{-1}) \tilde{e}^\top}{S_c}.$$

۳.۴ فاصله کوک

روش‌های زیادی برای بررسی تاثیر حذف مشاهده بر جوانب مختلف مدل برازش شده وجود دارد. یکی از این روش‌ها استفاده از فاصله کوک است، که نخستین بار توسط کوک (۱۹۷۷) معرفی شد. فاصله کوک تعمیم یافته با استفاده از توان دوم تغییرات وزنی حاصل از حذف مشاهده (ij) ام در برآورد پارامترها

به صورت

$$CD_{ij}(\beta, f) = \frac{(\hat{\theta} - \hat{\theta}_{(ij)})^T C (\hat{\theta} - \hat{\theta}_{(ij)})}{\hat{\sigma}^2}$$

تعریف می شود، که در آن

$$C = \begin{bmatrix} X^T V^{-1} X - tr(V^{-1})\Lambda & X^T V^{-1} N \\ N^T V^{-1} X & N^T V^{-1} N + \lambda K \end{bmatrix}. \quad (۱۲)$$

اکنون با استفاده از روابط (۷) ، (۸) و قضیه ۱ داریم:

$$CD_{ij}(\beta, f) = \frac{d_c^T V^{-1} \bar{H} d_c}{\hat{\sigma}^2} \left\{ \frac{d_c^T V^{-1} \tilde{e}}{d_c^T V^{-1} (I - \bar{H}) d_c} \right\}^2 = \frac{d_c^T V^{-1} \bar{H} d_c}{\hat{\sigma}^2 \{d_c^T V^{-1} (I - \bar{H}) d_c\}} t_c^2$$

که در آن

$$t_c = \frac{d_c^T V^{-1} \tilde{e}}{\sqrt{\{d_c^T V^{-1} (I - \bar{H}) d_c\}}}$$

باقی مانده استیودنت شده مورد c ام است. فاصله کوک برای جزء ناپارمتری β نیز به صورت

$$CD_{ij}(\beta) = \frac{(\hat{\beta} - \hat{\beta}_{(ij)})^T \{(I_p, \circ) C^{-1} (I_p, \circ)^T\}^{-1} (\hat{\beta} - \hat{\beta}_{(ij)})}{\hat{\sigma}^2}$$

تعریف می شود. از رابطه (۱۲) داریم $(I_p, \circ) C^{-1} (I_p, \circ)^T = X^T W X - tr(V^{-1})\Lambda$ آن گاه می توان نوشت:

$$CD_{ij}(\beta) = \frac{[d_c^T W X (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W d_c] (d_c^T V^{-1} \tilde{e})^2}{\hat{\sigma}^2 \{d_c^T V^{-1} (I - \bar{H}) d_c\}}$$

یا به عبارتی دیگر

$$CD_{ij}(\beta) = \frac{d_c^T W X (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W d_c}{\hat{\sigma}^2 \{d_c^T V^{-1} (I - \bar{H}) d_c\}} t_c^2.$$

مقادیر بزرگ برای $CD_{ij}(\beta)$ نشان دهنده تاثیر مشاهده (i, j) ام روی برآورد β است.

۴.۴ فاصله درستنمایی

معیار تشخیص تعمیم یافته دیگری که بر پایه لگاریتم تابع درستنمایی به دست می‌آید فاصله درستنمایی است (کوک و ویزبرگ، ۱۹۸۲). فرض کنید $L(\hat{\beta}, \hat{f}, X, Y)$ تابع لگاریتم درستنمایی تصحیح شده روی مجموعه داده‌ها باشد. در این صورت فاصله درستنمایی به صورت

$$LD_{ij}(\beta) = 2[L(\hat{\beta}, \hat{f}, X, Y) - L(\hat{\beta}_{(ij)}, \hat{f}, X, Y)]$$

تعریف می‌شود. با بسط تیلور $L(\hat{\beta}_{(ij)}, \hat{f}, X, Y)$ حول $\hat{\beta}$ خواهیم داشت:

$$LD_{ij}(\beta) = 2[J^{*T}(\hat{\beta})(\hat{\beta} - \hat{\beta}_{ij}) + \frac{1}{2}(\hat{\beta} - \hat{\beta}_{ij})^T \{-\ddot{J}(\hat{\beta})\}(\hat{\beta} - \hat{\beta}_{ij})]$$

که در آن $\ddot{J} = \frac{\partial^2 l_p^*(\beta, f, X, Y)}{\partial \beta \partial \beta^T} \big|_{\beta=\hat{\beta}, f=\hat{f}}$ و $J^{*T}(\hat{\beta}) = \frac{\partial \ell_p^*(\beta, f, X, Y)}{\partial \beta} \big|_{\beta=\hat{\beta}, f=\hat{f}}$. این نتایج دقیق است زیرا مشتق سوم برابر صفر است. چون $\ddot{J}^*(\hat{\beta}) = 0$ ، فاصله درستنمایی مورد نظر برابر

$$LD_{ij}(\beta) = \frac{d_c^T W X (X^T W X - \text{tr}(V^{-1}) \Lambda^{-1}) X^T W d_c}{\hat{\sigma}^2 \{d_c^T V^{-1} (I - \bar{H}) d_c\}}$$

است. همان‌طور که انتظار می‌رود $LD_{ij}(\beta) = CD_{ij}(\beta)$ است.

۵ مباحث تشخیصی حذف آزمودنی

در این قسمت تاثیر حذف آزمودنی روی برآورد β و f را بررسی می‌شود. فرض کنید $\hat{\theta}_{[i]} = (\hat{\beta}_{[i]}^T, \hat{f}_{[i]}^T)^T$ برآورد β و f با حذف آزمودنی i ام باشد در این صورت داریم:

$$\hat{\theta}_{[i]} = \hat{\theta} - C^{-1} \begin{bmatrix} X^T \\ N^T \end{bmatrix} V^{-1} E_i (I_{n_i} - \bar{H}_i)^{-1} E_i^T \bar{e} \quad (13)$$

که در آن

$$\bar{H}_i = \begin{bmatrix} X_i & N_i \end{bmatrix} \begin{bmatrix} X^T V^{-1} X - tr(V^{-1})\Lambda & X^T V^{-1} N \\ N^T V^{-1} X & N^T V^{-1} N + \lambda K \end{bmatrix}^{-1} \begin{bmatrix} X_i^T \\ N_i^T \end{bmatrix} V_i^{-1}.$$

با استفاده از (۱۳) آماره‌ی کوک تعمیم یافته برای (β, f) به صورت

$$CD_{[i]}(\beta, f) = \frac{\tilde{e}^T E_i (I_{n_i} - \bar{H}_i^T)^{-1} V_i^{-1} \bar{H}_i (I_{n_i} - \bar{H}_i)^{-1} E_i^T \tilde{e}}{\hat{\sigma}^2}$$

تعریف می‌شود. که در آن $\tilde{e} E_i^T$ بردار باقی مانده $1 \times n_i$ متناظر با آزمودنی i ام است. با استفاده از رابطه $\hat{\beta}_{[i]} - \beta = (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W E_i (I_{n_i} - \bar{H}_i)^{-1} E_i^T \tilde{e}$ پارامتر β به صورت

$$CD_{[i]}(\beta) = \frac{R_i^T H_{\beta,i} R_i}{\hat{\sigma}^2}$$

بیان می‌شود، که در آن $\tilde{e} E_i^T$ و $R_i = (I_{n_i} - \bar{H}_i)^{-1} E_i^T \tilde{e}$ را به صورت

$$H_{\beta,i} = E_i^T W X (X^T W X - tr(V^{-1})\Lambda)^{-1} X^T W E_i \quad (14)$$

تعریف می‌شود. از آن‌جا که حذف آزمودنی i ام n_i نقطه زمانی را شامل می‌شود بنابراین به محاسبه فاصله کوک برای تعیین تاثیر موضعی روی منحنی f نیاز است. در نتیجه طبق فرمول فاصله کوک خواهیم داشت:

$$CD_{[i]}(f) = \frac{R_i^T W_* N S^{-1} N_i^T (N_i S^{-1} N_i^T)^{-1} N_i)^{-1} N_i S^{-1} N^T W_* R_i}{\hat{\sigma}^2}$$

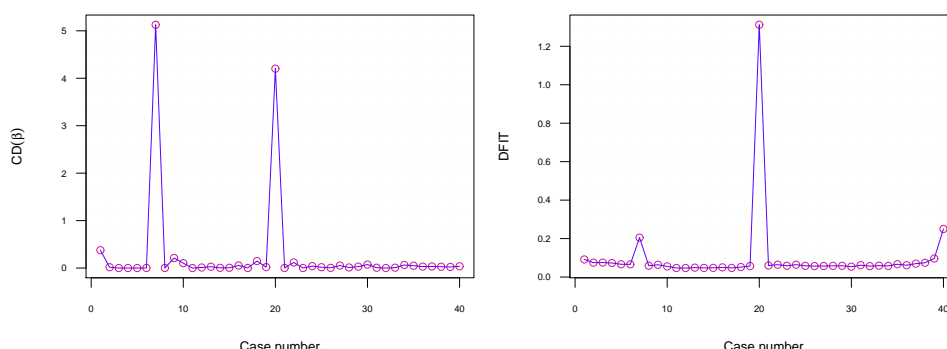
که در آن $S = N^T W_* N + \lambda K$ است.

۶ مطالعه شبیه‌سازی

برای بررسی مباحثی که عنوان شد متغیر پاسخ Y_{ij} از مدل

$$Y_{ij} = x_{ij}^{(1)}\beta_1 + x_{ij}^{(2)}\beta_2 + b_{1j} + f(t_{ij}) + \varepsilon_{ij} \quad i = 1, \dots, m; \quad j = 1, \dots, q.$$

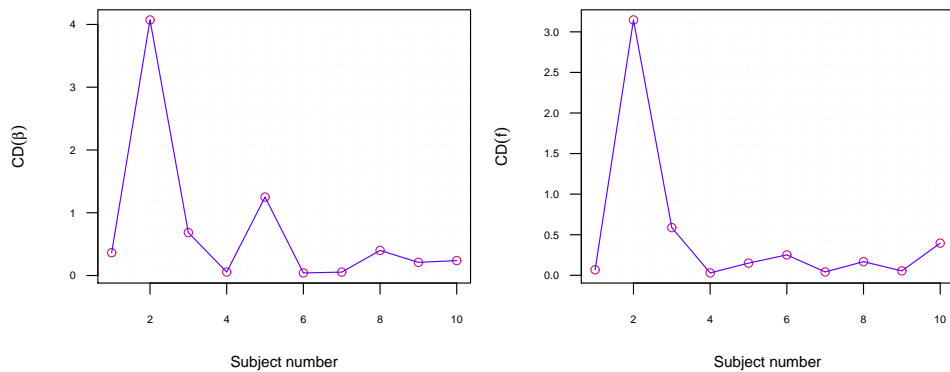
شبیه‌سازی می‌شود، که در آن q تعداد خوشه‌های مستقل و m اندازه خوشه در مطالعات طولی با اندازه نمونه $n = mq$ است (وانگ و همکاران، ۱۹۹۸). فرم ماتریسی این مدل را می‌توان به صورت $r = 1$ ، $b_1 =$ ، $Y = (Y_{11}, \dots, Y_{1q}, Y_{21}, \dots, Y_{2q}, \dots, Y_{m1}, \dots, Y_{mq})^T$ ، $(b_{11}, \dots, b_{1q})^T$ ، X ، f و ε هم ساختاری مانند Y دارند. در این شبیه‌سازی $q = 10$ و $m = 4$ است. همچنین در این شبیه‌سازی $\beta_1 = 1$ ، $\beta_2 = 2$ ، $x_{ij}^{(1)} \sim N(35, 3)$ ، $x_{ij}^{(2)} \sim N(23, 3)$ ، $f(t_{ij}) = (t - 0.5)^2$ و $\varepsilon_{ij} \sim N(0, \sigma^2)$ ، $b_{1j} \sim N(0, \sigma_1^2)$ ، $t_{i,j+1} = t_{ij} + 0.25$ ($j = 2, 3, 4$)، $t_{i1} \sim U(0, 0.25)$ و $\sigma^2 = \sigma_1^2 = 0.6^2 = 0.36$ انتخاب شد و $\Lambda = \text{diag}(0.5^2, 0.5^2)$ است. برای هر ترکیب از پارامترها ۱۰۰۰ بار تکرار انجام شده است. قبل از اجرای شبیه‌سازی تغییراتی در متغیر پاسخ هفتم و بیستم دادیم.



شکل ۱. مباحث تشخیصی حذف موردی برای داده‌های شبیه‌سازی شده: فاصله کوچک برای β و f

همچنین برای متغیر پیش‌بین X در آزمودنی دوم از شماره پنج تا هشت، هفت واحد به آن‌ها اضافه کردیم. شکل ۱ فاصله کوچک در برابر مشاهدات برای پارامتر β و قسمت ناپارامتری f را نشان می‌دهد. مشاهده هفتم و بیستم تأثیر زیادی روی برآورد β و نیز مشاهده بیستم روی برآورد f هم تأثیر زیادی داشته است. مطابق شکل ۲ نتایج حاصل از شبیه‌سازی نشان می‌دهند که آزمودنی دوم یک مشاهده موثر است. بعد از

آزمودنی دوم، آزمودنی پنج تاثیر بالایی در برآورد فاصله کوک روی پارامتر β دارد.



شکل ۲. مباحث تشخیصی حذف آزمودنی برای داده‌های شبیه‌سازی شده: فاصله کوک برای β و f

۷ تحلیل داده‌های فرسایش شدت باران

خاک‌های مناطق خشک و نیمه‌خشک ایران بشدت تحت تاثیر تجمع نمک‌ها به ویژه گچ و کربنات‌ها هستند. خاک‌های این مناطق ذاتا مواد آلی کم و ساختمان ناپایداری دارند. به علاوه پوشش گیاهی طبیعی سطح خاک‌ها ضعیف و در معرض عوامل فرساینده قرار دارد. اخیرا بلیانی (۲۰۱۷) پژوهشی به منظور بررسی میزان فرسایش پاشماني در خاک‌های با خصوصیات مختلف از استان‌های زنجان و فارس انجام داده است. در پژوهش او ۱۲ نوع بافت خاک مورد مطالعه قرار گرفته است. از بین ۱۲ کلاس بافت خاک ۱۱ بافت خاک از مناطق مختلف استان زنجان و یک بافت خاک سیلتي از استان فارس به تصادف انتخاب شدند. مجموعا تعداد ۱۴۴ جعبه خاک به ابعاد ۲۵ در ۳۵ سانتی متر مورد بررسی قرار گرفته است. در نمونه‌های خاک، درصد آهک (x_1)، درصد شن (x_2) و درصد رس (x_3) به روش هیدرومتری و درصد ماده آلی خاک (x_4) به روش آزمایشگاهی اندازه‌گیری شده است. از آنجایی که درصد گچ موجود در خاک (x_5) با استفاده از روش اندازه‌گیری هدایت الکتریکی است و این روش با خطا در اندازه‌گیری همراه هست میزان این متغیر با ۱۲ تکرار در هربافت معین شده است همچنین میزان سنگ ریزه خاک با جداسازی ذرات اولیه با تعیین نسبت جرمی آن (x_6) مشخص می‌شود. بدون در نظر گرفتن خطای در اندازه‌گیری بلیانی و همکاران رابطه بین میزان شدت فرسایش باران (y) و ویژگی‌های خاک (x_1-x_6)

را با یک مدل رگرسیونی چندگانه مورد مطالعه قرار دادند. در این بخش یک مدل کاملتری همراه با مباحث تشخیصی (تحلیل حساسیت) مربوط به آن ارایه می‌شود. نخست برای تعیین متغیر بخش ناپارامتری با استفاده از آماره‌ی یاتچپو (۲۰۰۳) فرضیه‌ی خطی بودن $f(t)$ بعبارتی فرضیه $f(t) = h(t, \beta)$: H_0 را به ازای یک تابع مشخص $h(\cdot)$ می‌توان آزمود:

$$Z_0 = \frac{n^{1/2}(s_{res}^2 - s_{diff}^2)}{s_{diff}^2} \xrightarrow{D} N(0, 1)$$

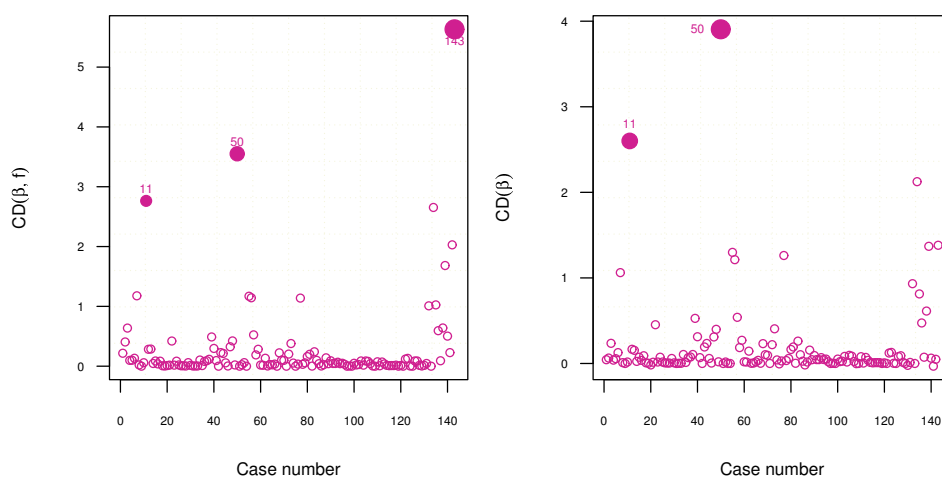
که در آن s_{res}^2 برآورد معمول واریانس باقی‌مانده در رگرسیون خطی و $s_{diff}^2 = \frac{1}{n} \sum_{i=2}^n (y_i - y_{i-1})^2$ است. برای انجام آزمون فوق برای هریک از متغیرهای x_j ، $j = 1, \dots, 6$ ، درایه‌های هر متغیر ابتدا باید به شکل صعودی $x_{1j} < \dots < x_{ij} < \dots < x_{nj}$ مرتب شود طوری که مقدار y_i در s_{res}^2 برابر i امین مشاهده متناظر به x_{ij} مرتب شده است. با انجام این آزمون متغیر x_6 را به عنوان بخش ناپارامتری در نظر گرفته می‌شود زیرا این آماره برای x_6 در بین سایر متغیرهای مستقل معنی‌دار است ($Z_0 = 10/531$). با توجه به اطلاعات فوق مدل

$$y = \sum_{j=1}^4 \beta_j x_j + \beta_5 x_5 + Z_1 b_1 + f(x_6) + \varepsilon, \quad x_5 = x_{*5} + \delta$$

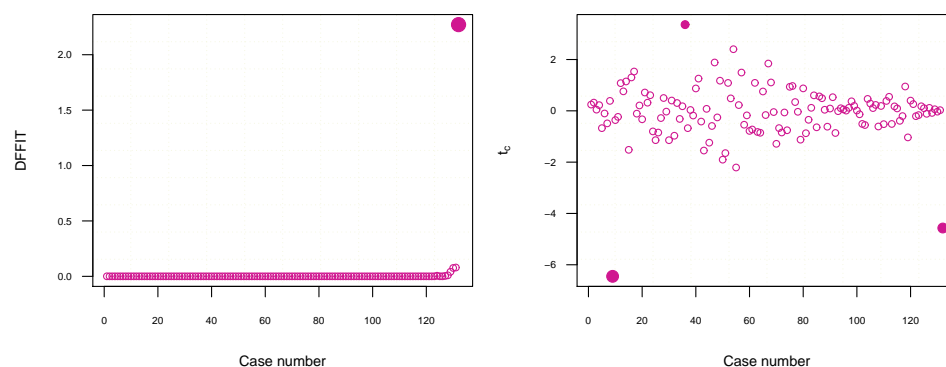
می‌تواند برای برازش داده‌ها مناسب باشد، که در آن Z_1 ماتریس طرح اثر نوع بافت خاک (b_1) با بعد 12×143 و δ بردار تصادفی خطا با توزیع $N(0, \Lambda)$ است. به دلیل گم شدن اندازه‌گیری‌های بعضی متغیرهای مربوط به یک جعبه خاک مدل با $n = 143$ برازش شده است. برای برآورد Λ با در نظر گرفتن \bar{x}_i به عنوان میانگین نمونه‌ای داده‌های تکراری مدل n ، $i = 1, \dots, m_i$ ، $j = 1, \dots, m_i$ ، $x_{ij} = x_{*ij} + \delta_{ij}$ برآورد ناریب Λ با روش گشتاوری به صورت

$$\hat{\Lambda} = \frac{\sum_{i=1}^n \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2}{\sum_{i=1}^n (m_i - 1)}$$

محاسبه می‌شود. فاصله کوک تعمیم‌یافته را برای β و (β, f) تحت مدل حذف موردی محاسبه شده و نمودار آن‌ها در شکل ۳ رسم شده است. همانطور که ملاحظه می‌شود مشاهدات ۵۰ و ۱۱ روی β و همین مشاهدات با مشاهده ۱۴۳ همزمان روی برآورد (β, f) موثر هستند. با توجه به نمودار باقی‌مانده‌های

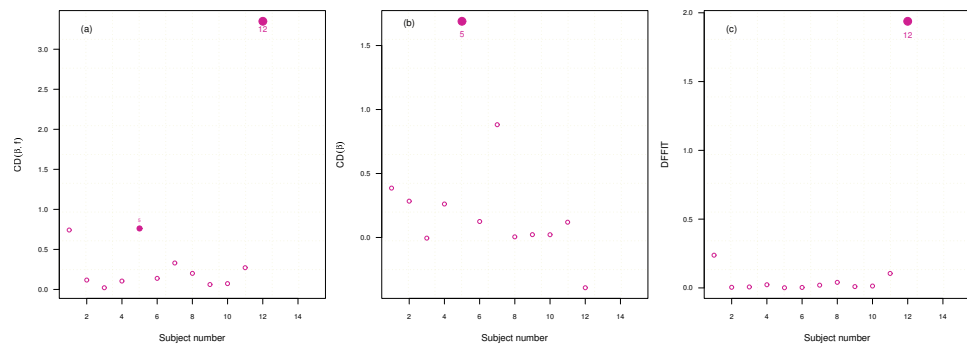


شکل ۳. مباحث تشخیصی حذف موردی برای داده‌های خاکشناسی: فاصله کوک برای β و f

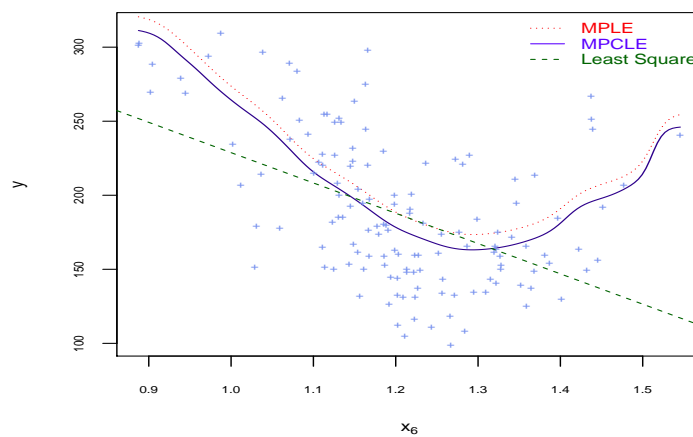


شکل ۴. نمودار باقی‌مانده استیودنت شده و $DFIT$ برای برآورد f

استیودنت شده در شکل ۴ سه مشاهده ۵۰، ۱۱ و ۱۳۴ بعنوان مشاهدات دور افتاده شناسایی می‌شوند. هیچ‌کدام از اندازه‌های $DFIT$ نگرانی جدی ایجاد نمی‌کنند، تنها مشاهده ۱۳۳ نسبت به بقیه مقدار بزرگ‌تری دارد که می‌تواند بعنوان مشاهده موثر روی برآورد f در نظر گرفته شود. برای بررسی اثر هر نوع بافت خاک در برآورد پارامترها نمودار $CD_{[i]}(\beta, f)$ در شکل ۵ رسم شده است. همانطور که ملاحظه می‌شود افت خاک از نوع ۵ روی برآورد β و بافت خاک از نوع ۱۲ روی برآورد f موثر است.



شکل ۵. مباحث تشخیصی حذف آزمودنی: فاصله کوک برای β و (β, f) و $DFIT$ برای f



شکل ۶. برآورد تابع ناپارامتری (MCPLE) منحنی با خط ممتد و MPLE (منحنی نقطه چین)

جدول ۱ برآورد β ، واریانس اثر تصادفی σ^2 ، واریانس خطای مدل σ^2 و واریانس خطای اندازه گیری Λ در صورت حذف و عدم حذف مشاهدات موثر را نشان می‌دهد. مقادیر داخل پرانتز مقادیر انحراف استاندارد برآوردها را نشان می‌دهد. درصد تغییرات مربوط به برآورد ضریب هر متغیر پس از حذف مشاهدات موثر در ستون آخر محاسبه شده است. شکل ۶ برآورد قسمت ناپارامتری مدل یعنی $f(x_6)$ را با استفاده از دو روش برآورد ماکسیمم درستنمایی تاوانیده (MPLE) و برآورد ماکسیمم درستنمایی تاوانیده

تصحیح شده (MCPLE) را نشان می‌دهد. پهنای باند که توسط اعتبارسنجی متقابل انتخاب شده برای هر دو روش برآورد به ترتیب ۰/۰۹۸ و ۰/۲۰۱ است. از شکل مشخص است وقتی از روش تصحیح شده ناکامورا استفاده می‌شود برآورد $f(x_6)$ کمتر از روشی است که خطاها را نادیده در نظر می‌گیرد. این اختلاف ارتفاع در دو منحنی به دلیل تصحیح خطای اندازه‌گیری است.

جدول ۱. برآورد ضرایب پارامتری مدل برازش شده به داده‌های فرسایش شدت باران

ضرایب	کل داده‌ها	حذف مشاهدات ۱۱ و ۵۰	درصد تغییرات
$\hat{\beta}_1$	۰/۸۸۸ (۰/۰۲۱)	۰/۴۷۵ (۰/۰۰۹)	۴۶/۵
$\hat{\beta}_2$	۰/۱۹۲ (۰/۰۰۳)	۰/۱۰۶ (۰/۰۱۱)	۸۱/۱
$\hat{\beta}_3$	۰/۲۲۰ (۰/۰۰۱)	۰/۳۵۱ (۰/۰۰۸)	۵۹/۵
$\hat{\beta}_4$	-۶/۴۴۲ (۱/۷۹۴)	-۵/۲۱۷ (۲/۰۰۳)	۱۸/۹
$\hat{\beta}_5$	-۲/۶۵۷ (۰/۸۴۱)	-۱/۱۰۶ (۰/۴۲۲)	۵۸/۳
$\hat{\sigma}_1^2$	۰/۰۱۶۹	۰/۰۱۰۳	۰/۶۶۰
$\hat{\sigma}^2$	۰/۰۸۵	۰/۱۷۹	۱۱۰/۵۸۸
$\hat{\Lambda}$	۰/۱۰۳	-	-

بحث و نتیجه‌گیری

در این مقاله مباحث تشخیصی در مدل خطی آمیخته نیمه‌پارامتری با خطا در اندازه‌گیری با استفاده از تابع امتیاز تصحیح شده ناکامورا بررسی شده است. مباحث تشخیصی حذف موردی و آزمودنی و نیز چندین ابزار تشخیصی برای شناسایی نقاط موثر پیشنهاد شده است. علاوه بر روش نظری با استفاده از شبیه‌سازی نشان داده شد که روش‌های تشخیصی مورد تعمیم یافته برای شناسایی مشاهدات موثر و دورافتاده مناسب هستند. این روش‌ها به تحلیل‌گر داده‌ها در تحلیل مدل خطی آمیخته نیمه‌پارامتری با خطا در اندازه‌گیری کمک شایانی می‌کند. برای مثال در داده‌های فرسایش شدت باران مشاهدات ۵۰ و ۱۱ از بافت خاک‌های استان زنجان به عنوان مشاهدات موثر روی برآورد ضرایب پارامتری مدل شناسایی شدند طوری که این مشاهدات با توجه به جدول ۱ بیشترین درصد تغییر را روی متغیر درصد شن و واریانس مدل ایجاد می‌کنند.

کنند. همچنین مشاهده ۱۴۳ که مربوط به بافت خاک استان فارس است موثرترین مشاهده روی جزء ناپارامتری مدل شناسایی شد. بعنوان یک نکته پایانی می توان گفت مباحث تشخیصی انجام شده در بخش ۷ را می توان مانند یک نوع تحلیل حساسیت درمباحث خاکشناسی در نظر گرفت که امکان سنجش حساسیت نمونه ها و نوع بافت خاک را در یک مدل دقیق را فراهم می آورد.

تشکر و قدردانی

نویسندگان از داوران، سردبیر و ویراستار محترم نشریه که با رهنمودهای ارزنده خود موجب بهتر شدن مقاله شدند و همچنین از دکتر بلیانی بابت در اختیار گذاشتن داده‌ها کمال تشکر را دارند.

مراجع

- Balyani, A. (2017), Quantification of Rainfall Erosion in Some Soils of the Semi-arid Regions in North West of Iran, *Ph.D. Thesis in University of Zanzan, Zanzan, Iran* .
- Banerjee, M., Frees, L. W. (1997), Influence Diagnostics for Linear Longitudinal Models, *Journal of the American Statistical Association*, **92**, 999–1005.
- Beckman, R. J., Nachtsheim, C. J. and Cook, R. D. (1987), Diagnostics for Mixed-models Analysis of Variance, *Technometrics*, **29**, 413–426.
- Christensen, R., Pearson, L. M. and Johnson, W. (1992), Case Deletion Diagnostics for Mixed Models, *Technometrics*, **34**, 38–45.
- Cook, R. D. (1977), Detection of Influential Observations in Linear Regression, *Technometrics*, **19**, 15–18.
- Cook, R. D. and Weisberg, S. (1982), *Residuals and Influence in Regression*, Chapman and Hall, London.
- Davidian, M, Giltinan, D. M. (1995), *Nonlinear Models for Repeated Measurement Data*, Chapman and Hall, London.
- Emami, H. (2015), Influence Diagnostics in Ridge Semiparametric Regression Models, *Journal of Statistics and Probability Letter*, **105**, 106–115.
- Emami, H. (2015), Influence Measures in Ridge Linear Measurement Error Models, *Journal of Statistical Research of Iran*, **12**, 39-56.

- Emami, H. and Emami, M. (2016), Influence Diagnostics in Constrained General Linear Models. *Communications in Statistics-Theory and Methods*, **45**, 5331-5340.
- Emami, H. and Emami, M. (2014), Local Influence in Constrained General Linear Models. *Journal of Data Science*, **12**, 717-726.
- Emami, H. and Rasekh, A. (2014), Influence Diagnostics on Testing Linear Hypothesis in Linear Models with Correlated Errors, *Communications in Statistics Theory and Methods*, **4**, 1050-1060.
- Fung, W. K., Zhu, Z. Y., Wei, B. C. and He, X. M. (2002), Influence Diagnostics and Outlier Tests for Semiparametric Mixed Models, *Journal of Royal Statitital Society*, **64**, 565- 579.
- Ghapani, F., Rasekh, A., and Babadi, B. (2016), The Weighted Ridge Estimator in Stochastic Restricted Linear Measurement Error Models, *Statistical Paper*, DOI 10.1007/s00362-016-0786-3.
- Harvill, D. A. (1977), Maximum Likelihood Approaches to Variance Component Estimation and Related Problems (with discussion), *Journal of the American Statistical Association*, **72**, 320-340.
- Harrison, D. and Rubinfeld, D. L. (1978), Hedonic Housing Prices And The Demand For Clean Air, *Journal of Environmental Economics and Managemment*, **5**, 81–102.
- Kim, C. (1996), Cook’s Distance In Spline Smoothing, *Statistics and Probability Letters*, **31**, 139–144.
- Kim, C., Kim, W. and Park, B. U. (2002). Influence Diagnostics in Semiparametric Regression Models, *Statistics and Probability Letters*, **60**, 49-58.
- Lesaffre, E. and Verbeke, G. (1998), Local Influence in Linear Mixed Models, *Biometrics*, **54**, 570–582.
- Liang, H. and Li, R. (2009), Variable Selection For Partially Linear Models With Measurement Errors, *Journal of the American Statistical Association*, **485**, 234-248.
- Nakamura, T. (1990), Corrected Score Function For Errors-in-Variables Models: Methodology and Application to Generalized Linear Models, *Biometrics*, **77**, 127–137.
- Wang, N., Lin, X., Gutierrez, R. G. and Carroll, R. J. (1998), Bias Analysis And SIMEX Approach In Generalized Linear Mixed Measurement Error Models, *Journal of American Statistical Association*, **93**, 249–261.

- Yatchew, A. (2003), *Semiparametric Regression for the Applied Econometrician*, Cambridge University press , Cambridge.
- Zare, K. and Rasekh, A. (2011), Diagnostic Measures for Linear Mixed Measurement Error Models, *SORT*, **35**, 125-144.
- Zare, K., Rasekh, A. and Rasekhi, A. (2011), Estimation of Variance Components In Linear Mixed Measurement Error Models, *Statistical Papers*, **53**, 849-863.
- Zhong, D., Lin, X. and Sower, M. (1998), Semiparametric Stochastic Mixed models for Longitudinal Data, *Journal of the American Statistical Association*, **93**, 710-719.
- Zhong, X. P., Fung, W. K. and Wei, B. C. (2002), Estimation in Linear Models with Random Effects and Errors in Variables, *Annals of the Institute of Statistical Mathematics*, **54**, 595-606.