

تحلیل نیم پارامتری مدل های رگرسیونی برای پاسخ های سری توانی آماسیده صفر با متغیرهای تبیینی گم شده

نفیسه خجسته و احسان بهرامی

گروه آمار، دانشکده علوم ریاضی، دانشگاه شهید بهشتی

تاریخ دریافت: ۱۳۹۷/۰۱/۲۰ تاریخ پذیرش: ۱۳۹۸/۱۷/۱۳

چکیده: در این مقاله پاسخ های شمارشی با تعداد صفر زیاد، که داده های آماسیده صفر نامیده می شوند، مورد تحلیل قرار گرفته اند. فرض می شود پاسخ ها از سری توانی آماسیده صفر پیروی می کنند. همچنین به دلیل وجود گم شدگی از نوع تصادفی در متغیرهای تبیینی برخی از داده های کاربردی، انواع روش های برآوردیابی پارامترهای مدل براساس تابع امتیاز با و بدون در نظر گرفتن گم شدگی برای مدل رگرسیونی ارایه شده است. در این میان، معلوم یا نامعلوم بودن احتمال انتخاب متغیر تبیینی گم شده منجر به ارایه روش نیم پارامتری برای برآورد پارامترها در مدل رگرسیونی سری توانی آماسیده صفر می شود. به منظور تشریح روش پیشنهادی، مطالعه ای شبیه سازی برای مدل رگرسیونی دوجمله ای منفی آماسیده صفر با متغیرهای تبیینی گم شده به عنوان یک مدل رگرسیونی سری توانی انجام می شود و سپس مثالی از داده های واقعی ارایه می شود. در انتها، عملکرد روش نیم پارامتری در مقایسه با روش ماکسیمم درستنمایی، مورد-کامل، احتمال وارون وزنی مورد بررسی و ارزیابی قرار گرفته است.

واژه های کلیدی: آماسیده صفر، داده های گم شده، تابع امتیاز، گم شدگی تصادفی، تحلیل نیم پارامتری، احتمال انتخاب.

۱ مقدمه

در بسیاری از رشته‌ها مانند زیست‌شناسی، پزشکی، جرم‌شناسی، اقتصاد، محیط زیست، علوم سیاسی، جامعه‌شناسی و غیره داده‌های شمارشی رخ می‌دهند، به طوری که اغلب این داده‌ها از توزیعی از خانواده سری‌های توانی مانند توزیع پواسون، دوجمله‌ای منفی و مانند آن تبعیت می‌کنند. بیشتر محققان برای تحلیل این داده‌ها از مدل‌های خطی تعمیم‌یافته مانند مدل رگرسیونی پواسون یا مدل رگرسیونی دوجمله‌ای منفی استفاده می‌نمایند. همچنین در تحلیل این داده‌ها با مسأله تعداد صفر زیاد، سروکار خواهیم داشت. به عنوان مثال می‌توان تحلیل داده‌های بیمه شخص ثالث، تعداد مفصل‌های آسیب دیده یک فرد، تعداد دفعاتی که یک دانشجو غیبت غیر موجه دارد و تعداد مقالات منتشر شده در هر سال توسط یک محقق اشاره کرد. به دلیل وجود تعداد صفر زیاد در این داده‌ها، برای مدل‌سازی نمی‌توان از مدل رگرسیونی پواسون یا دوجمله‌ای منفی استفاده نمود. بنابراین برای رفع این مشکل از مدل رگرسیونی پواسون آماسیده صفر^۱ یا مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر استفاده می‌شود. مدل رگرسیونی پواسون آماسیده صفر برای مدل‌بندی داده‌های شمارشی با تعداد صفر زیاد در حالت تک متغیره توسط لمبرت (۱۹۹۲) مورد مطالعه قرار گرفته است. سپس مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر توسط گرین و همکاران (۱۹۹۴) بیان شد. علاوه بر این، آزمون امتیاز برای مدل رگرسیونی پواسون آماسیده صفر در مقابل مدل رگرسیونی پواسون توسط جانساگل و هیند (۲۰۰۲) مطرح شد.

از طرفی موضوع گم‌شدگی داده‌ها، امری است که نمی‌توان از آن چشم‌پوشی کرد. ساختار گم‌شدگی نخستین بار توسط روبین (۱۹۷۶) مطرح شد. وی معتقد بود داده‌های گم‌شده دارای سه نوع ساختار گم‌شدگی هستند. ساختار اول: گم‌شدگی کاملاً تصادفی^۲ ($MCAR$)، ساختار دوم: گم‌شدگی تصادفی^۳ (MAR) و ساختار سوم: گم‌شدگی غیرتصادفی^۴ ($NMAR$)، در همین راستا لیتل و روبین (۲۰۱۴) و ابراهیم و همکاران (۱۹۷۶) استنباط‌هایی را برای مدل‌های خطی تعمیم‌یافته با حضور گم‌شدگی انجام دادند. همچنین روش انتخاب مدل، زمانی که متغیرهای تبیینی به طور تصادفی گم‌شده هستند، توسط چن و فو ارائه شد. از طرفی برآورد پارامترهای مدل رگرسیونی سری توانی، زمانی که متغیرهای تبیینی دارای گم‌شدگی هستند، مسأله مهمی است. در همین راستا برآورد پارامترها در این نوع از مدل‌های رگرسیونی به روش بیزی توسط ماسون و همکاران (۲۰۱۰) مورد پژوهش قرار گرفت. همچنین تحلیل نیم‌پارامتری به عنوان یک روش

^۱Zero Inflated

^۲Missing Completely At Random

^۳Missing At Random

^۴Not Missing At Random

بسیار مهم و کارساز در برآورد پارامترهای مدل مورد استفاده قرار می‌گیرد، از سوی دیگر تحلیل نیم پارامتری مربوط به پارامترهای مدل رگرسیونی اولین بار توسط وانگ و همکاران (۱۹۹۷) مورد بررسی و تحقیق قرار گرفت. از سویی دیگر، تحلیل نیم پارامتری مربوط به پارامترهای مدل‌های خطی تعمیم یافته توسط لیتون و نیلسن (۱۹۹۵) معرفی شدند. تحلیل نیم پارامتری روی مدل‌های رگرسیونی سری توانی آماسیده صفر تاکنون مورد بررسی قرار نگرفته است و تنها می‌توان به تحلیل نیم پارامتری مدل‌های رگرسیونی پواسون آماسیده صفر با امکان گم‌شدگی متغیرهای تبیینی گم‌شده که توسط لوکسا و همکاران (۲۰۱۶) اشاره نمود. نکته قابل توجه آن است لوکسا و همکاران (۲۰۱۶) داده‌هایی که مورد بررسی قرار داده‌اند از توزیع پواسون می‌باشند که با مسأله بیش پراکنش مواجه نیستند. زمانی مسأله بیش پراکنش رخ می‌دهد که واریانس متغیر شمارشی از مقدار مورد انتظار آن متغیر بیشتر شود. بنابراین نمی‌توان از مدل رگرسیونی پواسون آماسیده صفر استفاده نمود. از این رو کار انجام شده توسط لوکسا و همکاران (۲۰۱۶)، ناکارآمد می‌شود. بنابراین با ارایه مدل‌های رگرسیونی دوجمله ای منفی آماسیده صفر که نوعی از مدل‌های رگرسیونی سری توانی آماسیده صفر هستند و دارای بیش پراکنش ذاتی‌اند، می‌توان مدل‌های مناسبی را ارایه نمود. در این مقاله با تعمیم کار لوکسا و همکاران (۲۰۱۶)، تحلیل نیم پارامتری مدل‌های رگرسیونی سری توانی آماسیده صفر با امکان گم‌شدگی تصادفی در متغیرهای تبیینی انجام می‌شود.

در بخش ۲ مدل رگرسیونی سری توانی آماسیده صفر معرفی می‌شود. سپس در بخش ۳ انواع شیوه‌های برآورد در حالت با و بدون گم‌شدگی در متغیرهای تبیینی مطرح می‌شود. سپس با معرفی روش برآورد نیم پارامتری وزنی برای مدل‌های رگرسیونی سری توانی آماسیده صفر، جبرئیات این روش معرفی می‌شود. در بخش ۴ به مطالعه شبیه سازی روش‌های پیشنهادی بیان شده در بخش ۳ بر روی مدل رگرسیونی دوجمله ای منفی آماسیده صفر با متغیرهای تبیینی گم‌شده پرداخته می‌شود. در نهایت در بخش ۵، روش‌های پیشنهادی در مورد داده‌های واقعی پیاده سازی می‌شوند.

۲ مدل رگرسیونی سری توانی آماسیده صفر

متغیر تصادفی Y دارای توزیع سری توانی^۵ است، هرگاه دارای تابع جرم احتمال به صورت

$$f(y|\theta) = \frac{a(y)\theta^y}{g(\theta)}, \quad y = 0, 1, 2, \dots$$

^۵Power Series

باشد، که در آن $a(y)$ تابعی نامنفی ($a(y) > 0$) و $g(\theta) = \sum_{y=0}^{\infty} a(y)\theta^y$ تابعی مثبت، متناهی و مشتق‌پذیر از θ است و به عنوان ثابت نرمال‌کننده در این توزیع نقش ایفا می‌کند. واضح است با در نظر گرفتن $a(y) = \binom{m+y-1}{m-1}$ و $g(\theta) = (\frac{1}{1-\theta})^m = (\frac{1}{p})^m$ به طوری که $\theta = 1-p$ ، توزیع سری توانی به توزیع دوجمله‌ای منفی با پارامترهای (m, p) تبدیل می‌شود. همچنین با در نظر گرفتن $a(y) = \frac{1}{y!}$ و $g(\theta) = e^\theta$ ، این توزیع به توزیع پواسون با پارامتر θ تبدیل می‌شود.

عموماً توزیع داده‌های آماسیده صفر به صورت آمیخته از توزیع تباهیده با جرم ۱ در صفر و توزیع سری‌های توانی (PS) مانند پواسون و دوجمله‌ای منفی تعریف می‌شوند و با نماد $ZIPS(\theta, \omega)$ قابل نمایش هستند. تابع جرم احتمال توزیع سری‌های توانی آماسیده صفر به صورت

$$P(Y = y | \omega, \theta) = \begin{cases} \omega + (1 - \omega)f(0 | \theta) & y = 0 \\ (1 - \omega)f(y | \theta) & y = 1, 2, \dots \end{cases}$$

در نظر گرفته می‌شود، که در آن $f(0 | \theta)$ تابع جرم احتمال توزیع سری‌های توانی و ω پارامتر آمیختگی است که در بازه $0 < \omega < 1$ تغییر می‌کند. به عبارت دیگر این پارامتر، احتمال صفر بودن y را که نمی‌تواند به طور کامل توسط فرض مدل برآورد شود، معرفی می‌نماید و برای مدل‌هایی با تعداد صفرهای زیاد به کار می‌رود. احتمال تعداد صفر برابر با $\omega + (1 - \omega)f(0 | \theta)$ است در حالی که احتمال برای y هایی که مقدار مثبت را اختیار می‌کنند، برابر با $(1 - \omega)f(y | \theta)$ خواهد بود.

فرض کنید Y متغیر پاسخ از توزیع سری توانی آماسیده صفر با پارامترهای ω و θ باشد، همچنین با در نظر گرفتن X و Z به عنوان ماتریس طرح مربوط به متغیرهای تبیینی، مدل رگرسیونی سری توانی آماسیده صفر به صورت

$$Y \sim ZIPS(\omega, \theta), \quad \log(\theta) = \beta^\top \chi, \quad \text{logit}(\omega) = \gamma^\top \chi,$$

است، که در آن $\chi = (1^\top, X^\top, Z^\top)^\top$ ماتریس طرح مربوط به متغیرهای تبیینی و $\eta = (\beta^\top, \gamma^\top)^\top$ بردار پارامترهای مدل هستند. فرض کنید

$$\omega = H(\gamma^\top \chi), \quad \theta = \exp(\beta^\top \chi),$$

که در آن $H(u) = (1 + \exp(-u))^{-1}$ تابع پیوند لوژستیک و $\text{logit}(u) = H^{-1}(u)$ است. در

این صورت مدل رگرسیونی سری توانی آماسیده صفر معرفی می‌شود. تابع جرم احتمال Y در این مدل به صورت

$$P(Y = y|X, Z) = H(\gamma^\top \chi) I(y = \circ) + [1 - H(\gamma^\top \chi)] \frac{a(y)[\exp(\beta^\top \chi)]^y}{g(\exp(\beta^\top \chi))},$$

به دست می‌آید. از سوی دیگر

$$\begin{aligned} P(Y = \circ|X, Z) &= H(\gamma^\top \chi) + [1 - H(\gamma^\top \chi)] \frac{a(\circ)}{g(\exp(\beta^\top \chi))} \\ &= H(\gamma^\top \chi) \left[1 - \frac{a(y)}{g(\exp(\beta^\top \chi))}\right] + \frac{a(y)}{g(\exp(\beta^\top \chi))}, \end{aligned}$$

و تابع لگاریتم درستنمایی برای نمونه تصادفی $\{(Y_i, \chi_i) : i = 1, \dots, n\}$ به صورت

$$\begin{aligned} \ell(\eta) &= \sum_{i=1}^n I(Y_i = \circ) \left\{ \log[H(\gamma^\top \chi_i) \left(1 - \frac{a(\circ)}{g(\exp(\beta^\top \chi))}\right) + \frac{a(\circ)}{g(\exp(\beta^\top \chi))}] \right\} \\ &\quad + I(Y_i > \circ) \left\{ \log[1 - H(\gamma^\top \chi_i)] + \log(a(Y_i)) + Y_i \beta^\top \chi - \log(g(\exp(\beta^\top \chi))) \right\}, \end{aligned}$$

محاسبه می‌شود، که در آن $I(Y_i = \circ)$ و $I(Y_i > \circ)$ توابع نشانگر هستند. در ادامه شیوه‌های مختلف برآورد پارامترهای این مدل در حالت با و بدون گم‌شدگی در متغیرهای تبیینی مورد بررسی قرار می‌گیرند. نکته قابل تأمل آن است که با توجه به نوع مدل رگرسیونی و نوع روش برآورد پارامترها، برآوردهای پارامتری، ناپارامتری یا نیم‌پارامتری حاصل می‌شوند. یادآوری می‌شود که مدل رگرسیونی از نوع پارامتری است، یعنی رابطه بین متغیرهای پاسخ و متغیرهای تبیینی معلوم است اما یک جزء مدل به روش ناپارامتری برآورد می‌شود، این موضوع به تحلیل و برآورد نیم‌پارامتری می‌انجامد، درحالی که مدل نیم‌پارامتری فرض نشده است.

۳ برآوردیابی پارامترهای مدل با و بدون گم‌شدگی در متغیرهای تبیینی

وقتی متغیرهای تبیینی بدون گم‌شدگی هستند، یک روش برآورد بردار پارامترهای مدل رگرسیونی سری توانی آماسیده صفر، بر اساس تابع امتیاز

$$U_{F,n}(\eta) = \frac{1}{\sqrt{n}} \frac{\partial \ell(\eta)}{\partial \eta} = \frac{1}{\sqrt{n}} \left(\frac{\partial \ell(\eta)}{\partial \gamma} \right) = \frac{1}{\sqrt{n}} \sum_{i=1}^n S_i(\eta),$$

۱۰۰ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

است، که در آن

$$S_i(\eta) = \frac{\partial \ell_i(\eta)}{\partial \eta} = \left(S_{i1}^\top(\eta), S_{i2}^\top(\eta) \right)^\top, \quad i = 1, \dots, n,$$

$$S_{i1}(\eta) = \frac{\partial \ell_i(\eta)}{\partial \gamma}, \quad (1)$$

$$S_{i2}(\eta) = \frac{\partial \ell_i(\eta)}{\partial \beta}. \quad (2)$$

از طرفی $E(U_{F,n}(\eta)) = 0$ نشان می‌دهد $U_{F,n}(\eta)$ برآوردگر نااریب است. همچنین با قرار دادن $U_{F,n}(\eta) = 0$ برآوردگر ماکسیمم درست‌نمایی η که با $\hat{\eta}_F$ نشان داده می‌شود، به دست می‌آید. فرض کنید X_i متغیر تبیینی برای فرد i ام با گم‌شدگی از نوع MAR و

$$\delta_i = \begin{cases} 1 & X_i \text{ مشاهده شود} \\ 0 & X_i \text{ مشاهده نشود} \end{cases}$$

متغیری نشانگر باشد. همچنین W_i را یک متغیر جانشین برای X_i در نظر بگیرید، به طوری که Z_i و W_i همواره مشاهده می‌شوند. تحت ساختار MAR احتمال انتخاب به صورت

$$P(\delta_i = 1 | Y_i, X_i, V_i) = \pi(Y_i, V_i),$$

است (روبین، ۱۹۷۶)، که در آن $V = (Z^\top, W^\top)^\top$ ، در این حالت روش‌های برآورد پارامترهای مدل رگرسیونی سری توانی آماسیده صفر، وقتی متغیرهای تبیینی با گم‌شدگی هستند، به شرح زیر خواهند بود.

الف) روش برآورد مورد کامل:

برآورد مورد کامل^۶ یکی از روش‌های برخورد با داده‌های گم‌شدگی است و در بسیاری از نرم افزارهای آماری به صورت پیش فرض وجود دارد. در این روش داده‌های گم‌شده حذف می‌شوند. با توجه به آنکه ساختار گم‌شدگی داده‌ها MAR است، تابع برآورد CC به صورت

$$U_{CC,n}(\eta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \delta_i S_i(\eta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \delta_i \begin{pmatrix} S_{i1}(\eta) \\ S_{i2}(\eta) \end{pmatrix},$$

^۶Complete Case

به دست می‌آید، که در آن $S_{i1}(\eta)$ و $S_{i2}(\eta)$ به ترتیب در روابط (۱) و (۲) تعریف شده‌اند. با حل معادله $U_{CC,n}(\eta) = 0$ ، برآوردگر CC به دست می‌آید، که با $\hat{\eta}_{CC}$ نشان داده می‌شود. نکته قابل توجه آن است که $E(U_{CC,n}(\eta)) \neq 0$ (وانگ و همکاران، ۱۹۹۷). بنابراین برآوردگر $U_{CC,n}(\eta)$ ، اریب است. اگرچه تحلیل CC ساده است اما دو معضل به همراه دارد: ۱- از دست دادن میزان کارایی به علت حذف نمونه‌های ناقص، ۲- عامل بالقوه برای برآوردهای ناسازگار وقتی مجموعه داده‌های باقیمانده حاصل از حذف داده‌های گم‌شده از مجموعه داده‌های کامل، زیرنمونه‌ای تصادفی از داده‌های اصلی نیست.

(ب) روش برآورد احتمال معکوس وزنی:

ایده اصلی این روش از برآوردگر هورویتز-تامپسون (۱۹۵۲) است که توسط ژائو و همکاران (۱۹۹۲) توسعه یافت و منجر به معرفی برآوردگر شد. روش احتمال معکوس وزنی^۷ وقتی متغیرهای تبیینی دارای گم‌شدگی هستند برای داده‌های مشاهده شده وزن‌های وارون در نظر می‌گیرد و باعث می‌شود که اریبی حاصل از نمونه‌های ناقص حذف شده، کاهش یابد. از سوی دیگر اگر X_i مشاهده شود $\delta_i = 1$ و اگر X_i مشاهده نشود، $\delta_i = 0$. همچنین تحت ساختار MAR، احتمال انتخاب یا مشاهده X_i به صورت

$$P(\delta_i = 1 | Y_i, X_i, V_i) = \pi(Y_i, V_i),$$

در نظر گرفته می‌شود. بنابراین، تابع برآورد وزنی به صورت

$$U_{W,n}(\eta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\delta_i}{\pi(Y_i, V_i)} S_i(\eta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\delta_i}{\pi(Y_i, V_i)} (S_{i1}(\eta)),$$

به دست می‌آید، که در آن $S_{i1}(\eta)$ و $S_{i2}(\eta)$ به ترتیب در روابط (۱) و (۲) تعریف شده‌اند. با توجه به شرایط $\pi(Y_i, V_i)$ ، تابع برآورد متفاوت و در نتیجه برآوردهای متفاوتی به شرح زیر حاصل خواهند شد.

الف) $\pi(Y_i, V_i)$ معلوم باشد:

در بسیاری از مطالعات $\pi(Y_i, V_i)$ در مرحله طراحی مشخص می‌شود (برسلو و همکاران، ۱۹۸۸). در این حالت اگر $\pi(Y_i, V_i)$ از ابتدا به طور درست مشخص شده باشد، تابع برآورد وزنی بدون تغییر خواهد شد. در این صورت با حل معادله $U_{W,n}(\eta) = 0$ ، برآوردگر وزنی برای η ، که با $\hat{\eta}_W$ نشان داده می‌شود، به دست می‌آید. از طرفی دیگر $E(U_{W,n}(\eta)) = 0$ (وانگ و همکاران، ۱۹۹۷)، بنابراین $U_{W,n}(\eta)$ برآوردگر نااریب خواهد شد.

⁷Inverse Probability Weighted

۱۰۲ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

(ب) $\pi(Y_i, V_i)$ نامعلوم باشد:

برآورد $\pi(Y_i, V_i)$ به روش ناپارامتری یا روش پارامتری امکان‌پذیر است. فرض کنید v_1, \dots, v_m مقادیری متمایز از V_i را وقتی $V_i = (Z_i^T, W_i^T)$ گسسته در نظر گرفته شود را نشان دهند، آنگاه برآورد ناپارامتری $\pi(y, v)$ به صورت

$$\hat{\pi}(y_i, v_i) = \frac{\sum_{k=1}^n \delta_k I(Y_k = y, V_k = v)}{\sum_{i=1}^n I(Y_i = y, V_i = v)},$$

به دست می‌آید، که در آن مقادیر Y به صورت $0, 1, \dots$ و $y = 0$ و مقادیر V نیز به صورت $\{v_1, \dots, v_m\}$ در نظر گرفته می‌شوند. در این مقاله فرض می‌شود که V_i گسسته است. (در حالتی که V_i ها پیوسته باشند برای برآورد $\pi(Y_i, V_i)$ ، به وانگ و همکاران (۱۹۹۷) مراجعه شود). بنابراین با جایگذاری $\hat{\pi}(Y_i, V_i)$ در تابع برآورد وزنی، تابع برآورد نیم‌پارامتری وزنی به صورت

$$U_{Ws,n}(\eta, \hat{\pi}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\delta_i}{\hat{\pi}(Y_i, V_i)} S_i(\Theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\delta_i}{\hat{\pi}(Y_i, V_i)} (S_{i1}(\eta), S_{i2}(\eta)),$$

به دست می‌آید، که در آن $S_{i1}(\eta)$ و $S_{i2}(\eta)$ به ترتیب در روابط (۱) و (۲) تعریف شده‌اند. با حل معادله $U_{Ws,n}(\eta) = 0$ ، برآوردگر نیم‌پارامتری وزنی برای η ، که با $\hat{\eta}_{Ws}$ نشان داده می‌شود، به دست می‌آید. از سوی دیگر رابطه

$$U_{Ws,n}(\eta, \hat{\pi}) - U_{Ws,n}(\eta, \pi) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[\frac{\delta_i - \pi(Y_i, V_i)}{\pi(Y_i, V_i)} \right] S_i^*(\eta) + o_p(1),$$

برقرار است، که در آن $i = 1, \dots, n$ ، $S_i^*(\eta) = E(S_i(\eta) | Y_i, V_i)$ (کوچک) برای تقریب فرمول‌ها به کار می‌رود به این مفهوم که تقریب‌های نامتناهی و تقریب‌های چند جمله‌ای نامتناهی را وقتی متغیر به سمت صفر میل می‌کند، مشخص می‌کند. ($o_p(1)$ مجموعه‌ای از عبارات هستند که بزرگترین توان آنها، توان یک است. وقتی متغیر به سمت صفر میل می‌کند، این عبارت برابر با صفر است.) در این صورت مشخص می‌شود که $U_{Ws,n}(\eta, \hat{\pi})$ ، تابع برآورد تقریباً نااریب است.

۴ مطالعه شبیه‌سازی

تولید داده‌ها در یک فرایند دو مرحله‌ای شبیه‌سازی می‌شوند. ابتدا داده‌های کامل تولید می‌شود. برای این منظور ابتدا متغیرهای تبیینی و متغیر پاسخ بر اساس مدل رگرسیون دوجمله‌ای منفی آماسیده صفر تولید

می‌شود. سپس با ایجاد گم‌شدگی از نوع تصادفی در متغیر تبیینی، تولید داده‌ها صورت می‌پذیرد. آنگاه نحوه برآورد پارامترهای مدل بیان می‌شود. ابتدا متغیرهای تبیینی X و Z به ترتیب از توزیع یکنواخت روی بازه $(-۱, ۲)$ و توزیع برنولی با احتمال موفقیت $۰/۴$ به تصادف تولید می‌شوند. در مرحله دوم در متغیرهای تبیینی X ، گم‌شدگی از نوع تصادفی ایجاد شده است. پس از تولید متغیرهای تبیینی، ابتدا متغیر پاسخ Y_i^{NB} از روی مدل رگرسیونی دوجمله‌ای منفی به صورت

$$Y_i^{NB} \sim NB(\Theta_i, m), \quad m = ۱۰,$$

$$\log \theta_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i,$$

تولید می‌شود. حال برای تولید Y_i از مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر

$$Y_i \sim ZINB(\theta, m, \omega), \quad m = ۱۰,$$

$$\text{logit}(\omega_i) = \gamma_0 + \gamma_1 X_i + \gamma_2 Z_i,$$

$$\log \theta_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i,$$

استفاده می‌شود، که در آن پارامتر آمیختگی از روی مدل به صورت

$$\omega_i = \frac{\exp(\gamma_0 + \gamma_1 X_i + \gamma_2 Z_i)}{1 + \exp(\gamma_0 + \gamma_1 X_i + \gamma_2 Z_i)}$$

به دست می‌آید. با در نظر گرفتن متغیر تصادفی $U_i^{(۱)}$ که از توزیع یکنواخت روی بازه $(۰, ۱)$ تولید می‌شود. می‌توان $U_i^{(۱)}$ را با پارامتر محاسبه شده ω_i مورد مقایسه قرار داد و نتایج زیر را به دست آورد. اگر $U_i^{(۱)} \geq \omega_i$ ، آنگاه مقدار $Y_i = Y_i^{NB}$ در نظر گرفته می‌شود. به عبارت دیگر مقدار Y_i همان نمونه تولید شده از توزیع دوجمله‌ای منفی است.

اگر $U_i^{(۱)} < \omega_i$ ، آنگاه مقدار Y_i صفر در نظر گرفته می‌شود. شایان یادآوری است مقادیر واقعی بردار پارامترهای مدل رگرسیونی دو جمله‌ای منفی آماسیده صفر به صورت

$$\eta = (\gamma, \beta)^T = (\gamma_0, \gamma_1, \gamma_2, \beta_0, \beta_1, \beta_2)^T,$$

$$\eta_1 = (-۱, -۱, ۰/۵, ۱, ۰/۷, ۱)^T,$$

$$\eta_2 = (-۰/۱, -۱, ۰/۵, ۱, ۰/۷, ۱)^T,$$

است. بنابراین به این شیوه یک نمونه تصادفی از مدل رگرسیون دوجمله‌ای منفی آماسیده صفر تولید شد. در این بخش با استفاده از مکانیسم گم‌شدگی تصادفی، احتمال انتخاب متغیر تبیینی X_i به صورت

۱۰۴ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

$$\begin{aligned}\pi(Y_i, V_i) &= P(\delta_i = 1 | Y_i, X_i, V_i) \\ &= H(\mu_0 + \mu_1 Y_i I(Y_i < \epsilon) + \mu_2 I(Y_i \geq \epsilon) + \mu_3 Z_i + \mu_4 W_i),\end{aligned}$$

در نظر گرفته می‌شود، که در آن $V_i = (Z_i, W_i)$ و W_i متغیر تصادفی دودویی به صورت

$$W_i = \begin{cases} 1 & X_i \leq 0 \\ 0 & X_i > 0 \end{cases}$$

و H تابع پیوند لوژستیک است. همانند بخش ۱، یک متغیر تصادفی مانند $U_i^{(2)}$ در نظر گرفته می‌شود که از توزیع یکنواخت روی بازه $(0, 1)$ تولید می‌شود. سپس با محاسبه $\pi(Y_i, V_i)$ و مقایسه آن با $U_i^{(2)}$ نتایج زیر در مورد تولید داده‌های گم‌شده به دست می‌آید:

- اگر $U_i^{(2)} < \pi(y_i, v_i)$ ، آن‌گاه مقدار X_i گم‌شده در نظر گرفته می‌شود.

- اگر $U_i^{(2)} \geq \pi(y_i, v_i)$ ، آن‌گاه مقدار X_i تغییر نمی‌کند.

حال با در نظر گرفتن مقادیر واقعی برای بردار پارامترهای احتمال انتخاب $\mu = (\mu_0, \dots, \mu_4)^T$ احتمال انتخاب متغیر X_i محاسبه می‌شود. برای این منظور، می‌توان دو حالت در نظر گرفت:

الف- نسبت گم‌شدگی متغیر تبیینی X_i کمتر از 50% درصد باشد.

ب- نسبت گم‌شدگی متغیر تبیینی X_i بیش‌تر از 50% درصد باشد.

برای اینکه حالت (الف) اتفاق بیافتد لازم است مقادیر واقعی پارامترهای مربوط به احتمال انتخاب X_i به صورت

$$\mu_1 = (\mu_0, \dots, \mu_4)^T = (-1, 0/3, 10, 0/5, 1)^T,$$

در نظر گرفته شوند. در این حالت درصد داده‌های گم‌شده در متغیر تبیینی X کمتر از 50% درصد خواهد شد. همچنین برای این‌که حالت (ب) رخ دهد مقادیر واقعی بردار پارامترهای مربوط به احتمال انتخاب X_i به صورت

$$\mu_2 = (\mu_0, \dots, \mu_4)^T = (-1/8, 0/3, 10, 0/5, 1)^T,$$

در نظر گرفته می‌شود. حال احتمال انتخاب متغیر تبیینی X_i برای دو حالت (الف) و (ب) به طور جداگانه محاسبه و با متغیر $U_i^{(2)}$ مقایسه و دو سری نمونه تصادفی برای متغیر X_i (حالت الف و ب) تولید می‌شود.

به طوری که متغیر X_i در کنار متغیر تبیینی Z_i روی متغیر پاسخ Y_i که در بخش اول تولید شده‌اند برازش داده می‌شوند.

حال براساس مراحل ب و مدل مورد نظر، چهار روش برآوردیابی برای برآورد بردار پارامترهای Θ مورد استفاده قرار می‌گیرد و چهار برآوردگر به صورت زیر حاصل می‌شوند.

- $\hat{\eta}_F$: برآوردگر ماکسیم درستنمایی براساس تابع امتیاز $U_{F,n}(\eta)$ برای داده‌های کامل.

- $\hat{\eta}_{CC}$: برآوردگر CC براساس تابع برآورد $U_{CC,n}(\eta)$.

- $\hat{\eta}_W$: برآوردگر IPW وزنی براساس تابع برآورد $U_{W,n}(\eta, \pi)$.

- $\hat{\eta}_{Ws}$: برآوردگر IPW نیم‌پارامتری براساس تابع برآورد نیم پارامتری وزنی $U_{Ws}(\eta, \hat{\pi})$.

در جدول ۱ نتایج شبیه‌سازی برای حالت (الف) ارائه شده است، که در آن میزان اریبی ($Bias$)، خطای استاندارد (SE)، احتمال پوشش (CP) هریک از پارامترهای مدل برای حالت با و بدون گم‌شدگی در متغیر تبیینی X برای حجم نمونه‌های $n = 500$ و $n = 1000$ محاسبه شده است. در این حالت درصد گم‌شدگی کمتر از ۵۰ درصد X است. در جدول ۲ نتایج شبیه‌سازی برای حالت (ب) بیان شده است. این جدول نیز میزان اریبی، خطای استاندارد و احتمال پوشش هر یک از پارامترهای مدل برای حالت با و بدون گم‌شدگی در متغیر تبیینی برای حجم نمونه‌های $n = 1000$ و $n = 2000$ محاسبه شده است. همچنین درصد گم‌شدگی بالاتر از ۵۰ درصد X_i ، فرض شده است. شایان یادآوری است که برآوردگر ماکسیم درستنمایی براساس داده‌های بدون گم‌شدگی ($\hat{\eta}_F$) به عنوان یک معیار مقایسه در مقابل برآوردهای دارای گم‌شدگی در مطالعات شبیه‌سازی به کار برده می‌شود. در جدول ۱ میزان گم‌شدگی در متغیر تبیینی X_i برای حجم‌های $n = 500$ و $n = 1000$ به ترتیب ۲۶/۰ و ۳۳/۰ است. از سوی دیگر درصد صفرهای تولید شده ($Y = 0$) از مجموعه داده‌های کامل (بدون گم‌شدگی) برای حجم‌های $n = 500$ و $n = 1000$ به ترتیب ۴۲/۰ و ۴۷/۰ است. همچنین درصد صفرهای تولید شده ($Y = 0$) از مجموعه داده‌های ناکامل (با گم‌شدگی) برای حجم‌های $n = 500$ و $n = 1000$ به ترتیب ۱۷/۰ و ۲۳/۰ است. جدول ۱ نشان می‌دهد که میزان اریبی برای هر کدام از برآوردهای $\hat{\eta}_F$ ، $\hat{\eta}_{CC}$ ، $\hat{\eta}_W$ و $\hat{\eta}_{Ws}$ وجود دارد. به طور خاص، برآوردگر CC برای پارامترهای رگرسیون لوژیستیک $(\lambda_0, \lambda_1, \lambda_2)^T$ دارای میزان اریبی بالایی است. همچنین مقادیر SE با افزایش تعداد نمونه، کاهش یافته‌اند. همچنین مقادیر SE تحت برآوردگر $\hat{\eta}_{Ws}$ از مقادیر SE تحت برآوردگر $\hat{\eta}_{CC}$ کمتر هستند. به عبارت دیگر می‌توان نشان داد که $SE(\hat{\eta}_F) \leq SE(\hat{\eta}_{Ws}) \leq SE(\hat{\eta}_W)$. احتمال پوشش براساس برآوردگر CC بدترین وضعیت ممکن را دارا است. از آنجا که احتمال پوشش ۹۵/۰ فرض می‌شود، مقادیر احتمال پوشش CC اختلاف

۱۰۶ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

بسیاری با احتمال پوشش ۰/۹۵ دارند. احتمال پوشش با روش ماکسیمم درستنمایی از روش‌های دیگر بهتر است، اما این روش به مجموعه داده‌های کامل نیاز دارد. از طرفی احتمال پوشش $\hat{\eta}_W$ از احتمال پوشش $\hat{\eta}_{W_s}$ بهتر بوده است اما این روش برآورد IPW وزنی، ارزش عملی محدودی دارد به این خاطر که استفاده از آن به یک دانش خاصی برای تعیین مدل احتمال انتخاب نیازمند است. روش نیم‌پارامتری در مقابل روش IPW دارای مزیت بوده و در این روش شکل پارامتری برای احتمال انتخاب در نظر گرفته نمی‌شود.

در جدول ۲ میزان گم‌شدگی در متغیر تبیینی X ، برای حجم‌های $n = 1000$ و $n = 2000$ به ترتیب ۰/۶۰ و ۰/۶۷ است. از سوی دیگر درصد صفرهای تولید شده ($Y = 0$) از مجموعه داده‌های کامل برای حجم‌های $n = 1000$ و $n = 2000$ به ترتیب ۰/۵۳ و ۰/۵۷ است. همچنین درصد صفرهای تولید شده از مجموعه داده‌های ناکامل (با گم‌شدگی) برای حجم‌های $n = 1000$ و $n = 2000$ به ترتیب ۰/۳۲ و ۰/۳۶ است. با این وجود با در نظر گرفتن فرض‌ها و نتایج عددی متفاوت برای دو جدول، نتایج تحلیلی یکسان هستند. برای مثال، میزان بالای گم‌شدگی در متغیر تبیینی X_i در جدول ۲ باعث می‌شود احتمال پوشش تحت برآوردگر $\hat{\eta}_{CC}$ بسیار کم باشد. درحالی که در جدول ۱ نیز احتمال پوشش کم است اما در جدول ۲ این مقدار بیشتر است. هدف از در نظر گرفتن درصدهای متفاوت گم‌شدگی برای متغیر تبیینی بررسی عملکرد برآوردگرها درحالت‌های متفاوت بوده است. نتایج نشان از عملکرد یکسان در حالت‌های متفاوت دارند.

۵ مثال کاربردی

داده‌ها مربوط به تعداد ماهی‌هایی است ($count$) که در یک پارک آبی توسط بازدیدکنندگان گرفته شده‌اند. در این مطالعه از ۲۵۰ گروه از بازدیدکنندگان تشکیل شده است. عوامل موثر بر تعداد ماهی‌های گرفته شده، تعداد کودکان در گروه ($C.child$)، تعداد بزرگسالان در گروه ($C.persons$) و وجود یا عدم وجود راهنما ($leder$) هستند. تعداد ماهی‌ها به عنوان متغیر پاسخ شمارشی و تعداد کودکان، تعداد بزرگسالان و وجود یا عدم وجود راهنما به عنوان متغیر تبیینی گسسته در نظر گرفته شده‌اند. برخی شاخص‌های مربوط به متغیرهای مورد علاقه به طور خلاصه در جدول ۳ گردآوری شده است.

همان‌طور که ملاحظه می‌شود، تعداد صفر زیادی در متغیر پاسخ تعداد ماهی‌ها ملاحظه می‌شود. همچنین در متغیر تبیینی تعداد بزرگسالان درصد گم‌شدگی ۰/۲۹ درصد وجود داشته است. با وجود تعداد زیاد صفر در داده‌ها برای اطمینان از این که کدام مدل برای برازش روی داده‌های مربوطه مناسب است

جدول ۱: نتایج شبیه‌سازی برای حالت (الف) و $m = ۱۰$

$n = ۱۰۰۰$				$n = ۵۰۰$				حجم نمونه
$\hat{\eta}_{Ws}$	$\hat{\eta}_W$	$\hat{\eta}_{CC}$	$\hat{\eta}_F$	$\hat{\eta}_{Ws}$	$\hat{\eta}_W$	$\hat{\eta}_{CC}$	$\hat{\eta}_F$	پارامتر
—/۰/۱۲	—/۰/۴۰	—/۰/۶۳	—/۰/۳۰	—/۰/۳۳	—/۰/۴۲	—/۰/۷۱۰	—/۰/۳۱	γ_0
۰/۳۰۶	۰/۴۱۳	۰/۴۰۷	۰/۲۹۱	۰/۳۴۴	۰/۴۷۳	۰/۴۶۸	۰/۳۱۲	<i>Bias</i>
۰/۹۴۰	۰/۹۴۱	۰/۵۰۲	۰/۹۴۴	۰/۹۳۴	۰/۹۳۷	۰/۴۰۸	۰/۹۴۸	<i>S.E.</i>
								<i>CP</i>
—/۰/۳۲	—/۰/۳۵	—/۰/۶۴۴	—/۰/۲۰	—/۰/۳۹	—/۰/۴۵	—/۰/۸۱۲	—/۰/۲۴	γ_1
۰/۲۹۳	۰/۴۱۲	۰/۴۰۶	۰/۲۷۶	۰/۳۲۲	۰/۴۳۳	۰/۴۲۷	۰/۳۰۵	<i>Bias</i>
۰/۹۴۵	۰/۹۵۰	۰/۴۰۶	۰/۹۵۱	۰/۹۳۷	۰/۹۴۸	۰/۳۰۵	۰/۹۵۰	<i>S.E.</i>
								<i>CP</i>
۰/۰/۲۹	—/۰/۰/۳۸	—/۰/۵۶۴	۰/۰/۳۱	۰/۰/۳۲	—/۰/۰/۴۰	—/۰/۷۲۰	۰/۰/۴۱	γ_2
۰/۲۶۳	۰/۴۴۲	۰/۴۳۴	۰/۲۵۴	۰/۳۱۱	۰/۴۶۷	۰/۴۴۱	۰/۲۹۸	<i>Bias</i>
۰/۹۴۶	۰/۹۵۰	۰/۵۰۹	۰/۹۴۷	۰/۹۳۹	۰/۹۴۷	۰/۴۰۷	۰/۹۵۰	<i>S.E.</i>
								<i>CP</i>
۰/۰/۰۰۹	—/۰/۰/۱۰	۰/۰/۵۹	۰/۰/۱۲	۰/۰/۱۰	—/۰/۰/۱۱	۰/۰/۶۰	—/۰/۰/۲۰	β_0
۰/۰/۵۳۱	۰/۰/۵۹	۰/۰/۵۸	۰/۰/۵۲	۰/۰/۵۹	۰/۰/۶۶	۰/۰/۶۱	۰/۰/۵۴	<i>Bias</i>
۰/۹۳۶	۰/۹۴۶	۰/۳۵۴	۰/۹۵۴	۰/۹۲۵	۰/۹۴۱	۰/۳۲۴	۰/۹۴۹	<i>S.E.</i>
								<i>CP</i>
—/۰/۰/۱۵	۰/۰/۱۰	—/۰/۰/۴۹	۰/۰/۱۰	—/۰/۰/۱۶	۰/۰/۱۲	—/۰/۰/۵۹	۰/۰/۱۱	β_1
۰/۰/۴۵	۰/۰/۵۲	۰/۰/۵۱	۰/۰/۴۳	۰/۰/۴۸	۰/۰/۵۳	۰/۰/۵۴	۰/۰/۴۴	<i>Bias</i>
۰/۹۳۹	۰/۹۴۵	۰/۶۵۰	۰/۹۴۶	۰/۹۳۵	۰/۹۳۴	۰/۴۵۱	۰/۹۵۲	<i>S.E.</i>
								<i>CP</i>
۰/۰/۱۴	—/۰/۰/۱۲	—/۰/۰/۵۰	۰/۰/۲۹	۰/۰/۱۵	—/۰/۰/۱۲	—/۰/۰/۵۲	۰/۰/۳۰	β_2
۰/۰/۴۱	۰/۰/۴۹	۰/۰/۵۲	۰/۰/۳۹	۰/۰/۴۵	۰/۰/۵۲	۰/۰/۵۲	۰/۰/۴۳	<i>Bias</i>
۰/۹۳۶	۰/۹۳۹	۰/۶۴۸	۰/۹۵۴	۰/۹۳۹	۰/۹۴۳	۰/۵۰۵	۰/۹۴۹	<i>S.E.</i>
								<i>CP</i>

چهار مدل را بر روی داده‌ها برازش داده می‌شود. با برازش چهار مدل رگرسیونی پواسون، دوجمله‌ای منفی، پواسون آماسیده صفر و دوجمله‌ای منفی آماسیده صفر با گم‌شدگی در متغیرهای تبیینی بر روی داده‌های مثال واقعی که در متغیر پاسخ آن تعداد زیادی صفر مشاهده شده است به کمک معیار AIC استدلال می‌شود.

همان‌طور که در جدول ۴ ملاحظه می‌شود، مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر که نسبت به مدل‌های دیگر دارای کمترین مقدار AIC است، مدل برتر محسوب می‌شود. این مدل به صورت

$$Y_i \sim ZINB(\theta_i, m, \omega_i),$$

$$\text{logit}(\omega_i) = \gamma_0 + \gamma_1 C.child_i + \gamma_2 C.persons_i + \gamma_3 lder_i,$$

$$\log \theta_i = \beta_0 + \beta_1 C.child_i + \beta_2 C.persons_i + \beta_3 lder_i,$$

در نظر گرفته می‌شود، که در آن متغیر نشانگر δ_i و متغیر W_i به عنوان متغیر جانشین برای $C.persons_i$

۱۰۸ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

جدول ۲: نتایج شبیه‌سازی برای حالت (ب) و $m = ۱۰$

حجم نمونه پارامتر	$n = ۲۰۰۰$				$n = ۱۰۰۰$			
	$\hat{\eta}_{Ws}$	$\hat{\eta}_W$	$\hat{\eta}_{CC}$	$\hat{\eta}_F$	$\hat{\eta}_{Ws}$	$\hat{\eta}_W$	$\hat{\eta}_{CC}$	$\hat{\eta}_F$
γ_0								
Bias	۰/۰۳۰	۰/۰۴۲	-۰/۸۵۴	-۰/۰۱۱	۰/۰۴۵	۰/۰۴۵	-۰/۹۴۴	۰/۰۱۱
S.E.	۰/۲۰۰	۰/۳۱۵	۰/۳۱۴	۰/۲۰۰	۰/۳۰۲	۰/۳۳۲	۰/۳۱۹	۰/۲۰۱
CP	۰/۹۴۰	۰/۹۴۲	۰/۰۰۲	۰/۹۴۹	۰/۹۳۸	۰/۹۴۵	۰/۵۵۰	۰/۹۵۵
γ_1								
Bias	۰/۰۷۲	۰/۰۷۰	-۰/۸۰۰	-۰/۰۴۰	۰/۰۷۵	-۰/۰۷۲	-۰/۸۰۳	-۰/۰۴۰
S.E.	۰/۳۱۵	۰/۳۳۲	۰/۳۳۲	۰/۳۱۶	۰/۳۱۵	۰/۳۳۳	۰/۳۳۲	۰/۳۱۴۳
CP	۰/۹۱۶	۰/۹۲۱	۰/۰۰۲	۰/۹۲۶	۰/۸۹۵	۰/۹۳۰	۰/۰۰۸	۰/۹۳۹
γ_2								
Bias	۰/۰۲۲	-۰/۰۱۱	۰/۹۹۱	۰/۰۰۳	-۰/۰۳۷	-۰/۰۲۱	-۱/۰۰۳	-۰/۰۰۹
S.E.	۰/۳۵۵	۰/۳۸۹	۰/۳۸۷	۰/۳۴۴	۰/۳۵۸	۰/۳۹۶	۰/۳۹۵	۰/۳۵۹
CP	۰/۹۳۹	۰/۹۴۵	۰/۰۰۲	۰/۹۳۷	۰/۹۳۹	۰/۹۴۰	۰/۳۸۱	۰/۹۵۳
β_0								
Bias	۰/۰۲۸	۰/۰۲۰	۰/۴۰۸	-۰/۰۱۱	۰/۰۳۰	-۰/۰۲۰	۰/۴۰۹	-۰/۰۱۳
S.E.	۰/۰۴۰	۰/۰۵۱	۰/۰۵۶	۰/۰۴۶	۰/۰۴۵	۰/۰۵۹	۰/۰۵۸	۰/۰۴۸
CP	۰/۹۰۱	۰/۹۴۰	۰/۰۰۱	۰/۹۳۴	۰/۸۷۶	۰/۹۳۷	۰/۰۳۰	۰/۹۵۰
β_1								
Bias	۰/۰۲۵	۰/۰۲۰	۰/۰۱۹۸	۰/۰۱۰	۰/۰۲۸	-۰/۰۲۱	-۰/۲۰۱	۰/۰۱۱
S.E.	۰/۰۳۴	۰/۰۴۱	۰/۰۴۰	۰/۰۳۳	۰/۰۳۷	۰/۰۴۶	۰/۰۴۸	۰/۰۳۴
CP	۰/۹۱۳	۰/۹۴۱	۰/۰۱۵	۰/۹۳۳	۰/۸۶۶	۰/۹۴۴	۰/۵۰۰	۰/۹۴۲
β_2								
Bias	۰/۰۱۲	۰/۰۱۰	-۰/۲۰۴	۰/۰۱۱	۰/۰۱۶	۰/۰۱۰	-۰/۲۱۵	۰/۰۱۰
S.E.	۰/۰۵۰	۰/۰۵۵	۰/۰۵۴	۰/۰۴۲	۰/۰۴۳	۰/۰۵۶	۰/۰۵۵	۰/۰۴۳
CP	۰/۹۲۵	۰/۹۳۰	۰/۰۰۱	۰/۹۴۴	۰/۹۲۲	۰/۹۳۵	۰/۰۵۰	۰/۹۴۰

جدول ۳: شاخص‌های مربوط به متغیر پاسخ

متغیر	مقادیر	فراوانی	درصد فراوانی
تعداد ماهی‌ها	۰	۱۱۸	۰/۴۷
	۱	۴۶	۰/۱۸
	≥ ۲	۸۶	۰/۳۵

به صورت

$$\delta_i = \begin{cases} ۱ & C.persons_i \text{ مشاهده شود} \\ ۰ & \text{در غیر این صورت} \end{cases} \quad W_i = \begin{cases} ۱ & C.persons_i \geq ۲/۵ \\ ۰ & C.persons_i < ۲/۵ \end{cases}$$

جدول ۴: شاخص‌های مربوط به متغیرهای تبیینی

متغیر	درصد گم‌شدگی	بیشترین مقدار	کمترین مقدار	میانگین	انحراف معیار
تعداد بزرگسالان	۰/۲۹	۴	۱	۲/۵۱	۰/۱۴۷
تعداد کودکان	۰	۳	۰	۰/۶۸	۰/۰۲۴
تعداد راهنما	۰	۱	۰	۰/۲۴	۰/۰۰۳

جدول ۵: معیار AIC برای مدل‌های مختلف با گم‌شدگی در متغیرهای تبیینی منفی آماسیده صفر

مدل رگرسیونی	پواسون	دوجمله‌ای منفی	پواسون آماسیده صفر	دوجمله‌ای منفی آماسیده صفر
AIC	۶۵۴/۲	۷۱۱/۸۱	۵۴۲/۶	۴۳۲/۱

تعریف می‌شوند. در نهایت احتمال انتخاب تحت ساختار MAR به صورت

$$\begin{aligned} \text{logit}(\pi(Y_i, V_i)) &= \text{logit}P(\delta_i = 1 | Y_i, V_i) \\ &= \eta_0 + \eta_1 \text{Count}_i + \eta_2 W_i + \eta_3 C.\text{child}_i + \eta_4 \text{slider}_i. \end{aligned}$$

فرض می‌شود. در داده‌های واقعی نسبت تعداد صفرهای مربوط به متغیر پاسخ در $n = ۲۵۰$ از داده‌های کامل برابر با $۰/۵۶$ درصد و نسبت تعداد صفرهای مربوط به متغیر پاسخ در $n = ۲۵۰$ از مجموعه داده‌های ناکامل $۰/۱۵$ درصد است. همچنین میزان گم‌شدگی متغیر تعداد بزرگسالان $۰/۲۹$ درصد است. نتایج جدول ۵ نشان می‌دهد که میزان ارزیابی برای هر کدام از برآوردکننده $\hat{\theta}_W$ ، $\hat{\theta}_{CC}$ ، $\hat{\theta}_F$ و $\hat{\theta}_{WS}$ وجود دارد.

به طور کلی تحلیل‌های داده‌های واقعی با داده‌های شبیه‌سازی نتایج یکسان را نشان می‌دهند و این موضوع را مشخص می‌کنند که روش برآورد نیم‌پارامتری برای مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر با متغیرهای تبیینی گم‌شده یک روش مطلوب است و می‌تواند به عنوان یک شیوه بامزیت به کار برده شود.

۶ بحث و نتیجه‌گیری

در این مقاله به بررسی و تحلیل نیم‌پارامتری مدل‌های رگرسیونی سری توانی آماسیده صفر مانند مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر با و بدون داده‌های گم‌شده در متغیرهای تبیینی پرداخته شد. این

۱۱۰ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

جدول ۶: برآورد پارامترهای مدل رگرسیون دوجمله‌ای منفی آماسیده صفر با امکان گم‌شدگی تصادفی در متغیرهای تبیینی

روش‌های برآوردیابی							
$\hat{\Theta}_{Ws}$	$\hat{\Theta}_W$	$\hat{\Theta}_{CC}$	$\hat{\Theta}_F$	$\hat{\Theta}_{Ws}$	$\hat{\Theta}_W$	$\hat{\Theta}_{CC}$	$\hat{\Theta}_F$
$\beta_0 (constant)$				$\gamma_0 (constant)$			
۰/۰۹۳	۰/۱۹۶	۰/۵۶۱	۰/۰۴۹	۰/۱۹۰	۰/۲۳۰	-۰/۹۶۱	۰/۰۷۲ Bias
$\beta_1 (C.child)$				$\gamma_1 (C.child)$			
۰/۱۳۲	۰/۱۸۹	۰/۷۵۲	۰/۰۴۰	۰/۱۶۵	۰/۲۲۰	-۰/۸۶۰	۰/۰۵۱ Bias
$\beta_2 (C.persian)$				$\gamma_2 (C.persian)$			
۰/۱۶۱	۰/۱۵۱	۰/۲۱۹	۰/۰۳۶	۰/۱۹۰	۰/۲۳۵	-۰/۸۳۱	۰/۰۴۷ Bias
$\beta_3 (lider)$				$\gamma_3 (lider)$			
۰/۱۵۳	۰/۲۱۱	۰/۶۰۱	۰/۰۳۲	۰/۱۷۴	۰/۲۳۰	-۰/۸۳۰	۰/۰۴۲ Bias

تحلیل نیم‌پارامتری را می‌توان برای مدل‌های رگرسیونی آماسیده در یک نقطه غیر از صفر نیز مورد بررسی قرار داد. با تعمیم این روش‌ها روی مدل‌های رگرسیونی سری توانی آماسیده در بیش از یک نقطه، تحلیل‌های جدیدی مورد بررسی و به عنوان پژوهش‌های آتی مورد توجه قرار خواهد گرفت. هم‌چنین با در نظر گرفتن نوع مکانیسم گم‌شدگی غیر تصادفی، تحلیل‌های نیم‌پارامتری جدیدی را می‌توان برای مدل‌های بیان شده ارایه داد.

۷ تقدیر و تشکر

از داوران و ویراستاران محترم مجله که با توصیه‌های بسیار مفید سبب ارتقای این مقاله و ارایه بهتر آن شده‌اند کمال تشکر را داریم.

مراجع

اسماعیل‌زاده، م. بهرامی سامانی، ا. (۱۳۹۷)، مدل‌های اثر تصادفی دومتغیره آماسیده برای پاسخ‌های آمیخته سری توانی نرمال، مجله علوم آماری، مقاله آماده انتشار.

Lambert, D. (1992), Zero-Inflated Poisson Regression, with an Application to Defects in Manufacturing, *Technometrics*, **34**, 1-14.

Greene, W. H. (1994), Accounting for Excess Zeros and Sample Selection in Poisson and Negative Binomial Regression Models, *Leonard N, Department of Economics, New York*.

- Jansakul, N., and Hinde, J. P. (2002), Score Tests for Zero-Inflated Poisson Models, *Computational Statistics and Data Analysis*, **40** , 75-96.
- Rubin, D. B. (1976), Inference and Missing Data, *Biometrika*, **63**, 581-592.
- Little, R. J., and Rubin, D. B. (2014), *Statistical Analysis with Missing Data*, John Wiley and Sons, New York.
- Ibrahim, J. G., Chen, M. H., Lipsitz, S. R., and Herring, A. H. (2005), Missing-Data Methods for Generalized Linear Models: A Comparative Review, *Journal of the American Statistical Association*, **100**, 332-346.
- Chen, X. D., and Fu, Y. Z. (2011), Model Selection for Zero-Inflated Regression with Missing Covariates, *Computational Statistics and Data Analysis*, **55**, 765-773.
- Mason, A., Best, N., Richardson, S., and PLEWIS, I. (2010), Strategy for Modelling Non-Random Missing Data Mechanisms in Observational Studies using Bayesian Methods, *Journal of Official Statistics*, **28** , 279-302.
- Wang, C. Y., Wang, S., Zhao, L. P., and Ou, S. T. (1997), Weighted Semi-parametric Estimation in Regression Analysis With Missing Covariate Data, *Journal of the American Statistical Association*, **92**, 512-525.
- Linton, O., and Nielsen, J. P. (1995), A Kernel Method of Estimating Structured Nonparametric Regression Based on Marginal Integration, *Biometrika*, **82**, 93-100.
- Lukusa, T. M., Lee, S. M., and Li, C. S. (2016), Semiparametric Estimation of a Zero-Inflated Poisson Regression Model with Missing Covariates, *Metrika*, **79**, 457-483.
- Horvitz, D. G., and Thompson, D. J. (1952), A Generalization of Sampling without Replacement from a Finite Universe, *Journal of the American Statistical Association*, **47**, 663-685.
- Zhao, L. P., and Lipsitz, S. (1992), Designs and Analysis of Two Stage Studies, *Statistics in Medicine*, **11**, 769-782.
- Breslow, N. E., and Cain, K. C. (1988), Logistic Regression for Two-Stage Case-Control Data, *Biometrika*, **75** , 11-20.

..... ۱۱۲ تحلیل نیم‌پارامتری مدل‌های رگرسیونی برای پاسخ‌های سری توانی آماسیده صفر

Wang, C. Y., Wang, S., Zhao, L. P., and Ou, S. T. (1997), Weighted Semiparametric Estimation in Regression Analysis with Missing Covariate Data, *Journal of the American Statistical Association*, **75**, 11-20.