

تأثیر نادیده گرفتن همبستگی در داده های طولی یا زوج شده بر استنباط پارامترها

مجتبی گنجعلی^۱ جواد قاسمیان^۲

چکیده

در مواقعي که اندازه گيري از يك صفت خاص، روی داده های طولی^۳ (Diggle^۴ و Diggran، ۱۹۹۴) یا زوج شده^۵ (Agresti^۶، ۱۹۹۶) در مقاطع زمانی مختلف انجام می دهيم، بحث همبستگي بين داده ها مطرح می شود. در اين مقاله، با استفاده از داده های دودوسي که در آن متغير پاسخ، نشان دهنده داشتن یا نداشت آسم است، نشان می دهيم که درنظر گرفتن یا نگرفتن اين همبستگي، چه تأثيری بر نتایج برآورد پارامترها و انحراف استاندارد برآورد کننده های پارامترهاي منتبه با متغير های تو صify دارد.
واژه های کلیدی: پاسخهای دودویی، مدل پرویت دو متغیر، متغير پنهان، متغير تو صify مانا، آسم.

۱. مقدمه

داده های طولی در آزمایش های رخ می دهند که يك دنباله از اندازه گیریها، بر روی تعدادی آزمودنی در زمانهای مختلف به دست می آید. هر چند که واحدهای مختلف از هم مستقلند، ولی پیشرفت های اخیر در تحلیل این گونه داده ها، تأکید بر استفاده از مدل هایی دارند که وابستگی بین اندازه گیریهاي متعلق به آزمودنی يکسان را درنظر می گيرند. در اين مقاله نشان می دهيم که درنظر نگرفتن اين وابستگي، هر چند روی برآورد پارامترها تأثیری ندارد ولی موجب بیش یا کم برآورد کردن خطای استاندارد برآورد کننده های پارامترهاي مدل می شود.
دریخش بعد، داده های دودویی داشتن یا نداشت آسم، [۶] و دریخش ۳ مدل پرویت دو متغیر را معرفی کرده ايم. دریخش ۴ تابع درستنمایی بيان و در نهایت، دریخش ۵، تحلیل آماری داده های بخش ۶، با و بدون درنظر گرفتن همبستگي آورده شده است.

^۱ گروه آماری دانشگاه شهید بهشتی

^۲ دانشجوی کارشناسی ارشد، دانشگاه شهید بهشتی

^۳ Longitudinal Data

^۴ Paired Data^۰

^۵ Agresti^۷

^۶ Diggle^۸

در این مدل a ، b و c پارامترهای مدل می‌باشند که باید برآورده شوند. در حالت کلی خطاهای ε_{i1} و ε_{i2} وابسته‌اند و

$$\text{var}(\varepsilon_{it}) = \sigma^2 t = 1,2$$

(در داده‌های دودویی، واریانس دلخواه برای خطاهای، در معادله (۱) قابل برآورده نیست [۵]) و

$$\text{cov}(\varepsilon_{i1}, \varepsilon_{i2}) = \rho$$

در حالت کلی پارامتر ρ باید برآورده شود. اگر فرض کنیم $\rho = 0$ ، می‌توان نشان داد چه اتفاقی در برآورده پارامترها و خطاهای استاندارد برآورده کننده‌های پارامتر مدل، روی می‌دهد.

۴.تابع درستنمایی

برای نیل به هدف نهایی، دو حالت کلی را درنظر می‌گیریم. ابتدا دو متغیر Y_{i1} و Y_{i2} را مستقل فرض کرده و در مرحله بعد همبستگی بین این دو متغیر را در نظر می‌گیریم. تابع درستنمایی برای مدل کلی که همبستگی را در نظر می‌گیرد به صورت زیر است:

$$L(a, b, c, \rho | y_{i1}, y_{i2}, G) = \quad (۲)$$

$$\prod_{i=1}^n P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}) =$$

$$\prod_{i=1}^n p_{00}^{(1-y_{i1})(1-y_{i2})} p_{01}^{y_{i1}(1-y_{i2})} p_{10}^{y_{i2}(1-y_{i1})} p_{11}^{y_{i1}y_{i2}}$$

که در آن:

$$p_{00} = P(Y_1 = 0, Y_2 = 0) = \Phi_r(-a - bG, -a - bG - c, \rho),$$

$$p_{01} = P(Y_1 = 0, Y_2 = 1) = \Phi_r(-a - bG, a + bG + c, -\rho),$$

$$p_{10} = P(Y_1 = 1, Y_2 = 0) = \Phi_r(a + bG, -a - bG - c, -\rho),$$

$$p_{11} = P(Y_1 = 1, Y_2 = 1) = \Phi_r(a + bG, a + bG + c, \rho),$$

واضح است که

$$\sum_{i=0}^1 \sum_{j=0}^1 p_{ij} = 1$$

شده است)، برای استیباط آماری کافی است.

جدول ۱ داده‌های کامل و بدون مقادیر گمشده برای متغیرهای پاسخ در این داده‌ها را نشان می‌دهد. همان طور که جدول ۱ نشان می‌دهد ۵۷۰ نفر پسر و ۵۹۰ نفر دختر به هر دو متغیر پاسخ (داشتن یا نداشتن آسم در ۹ و ۱۳ سالگی) جواب داده‌اند. با استفاده از این داده‌ها، نشان می‌دهیم که درنظر گرفتن همبستگی پاسخها در سن ۹ و ۱۳ سالگی، چه تأثیری بر استیباط پارامترها و خطای استاندارد برآورده کننده‌های پارامترها دارد.

۳. مدل مورد استفاده

برای بررسی تأثیر سن و زمان بر متغیر پاسخ (داشتن آسم)، مدل زیر را درنظر می‌گیریم:

$$y_{it}^* = a + bG_i + cI_{\{age\geq 13\}} + \varepsilon_{it} \quad (1)$$

وقتی $t = 1, 2$ ، که در آن برای آزمودنی آم، G جنس و $I_{\{age\geq 13\}}$ به صورت زیر است:

$$I_{\{age\geq 13\}} = \begin{cases} 1 & \text{اگر فرد ۱۳ ساله باشد,} \\ 0 & \text{اگر فرد ۱۳ ساله نباشد,} \end{cases}$$

و y_{it}^* برای $t = 1, 2$ خطای اندازه گیری متغیر y_{it} است. به متغیرهای y_{it} و y_{it}^* ، متغیرهای پنهان می‌گوییم. این متغیرها خود قابل مشاهده نیستند ولی باعث می‌شوند که پاسخهای گستته قابل مشاهده باشند. برای مثال داشتن یا نداشتن آسم یک آزمودنی خاص (y) به علت وجود متغیر پیوسته دیگری (y^*) است که خود مشاهده نشده است ولی وقتی از آستانه خاصی گذر کند، فرد دچار بیماری آسم می‌شود. در اینجا مقدار آستانه، صفر فرض شده است. در حالتی که مقدار آستانه مثلاً مقدار d باشد، می‌توان از متغیر $d - y^*$ به عنوان متغیر پنهان استفاده کرد. بنابراین ثابت d ، قابل برآورده کردن نیست [۵].

چون داده‌های مشاهده شده، در دو حالت وجود دارند، لذا برای $t = 1, 2$ ، متغیرهای مشاهده شده y_{it} به صورت زیر تعریف می‌شوند:

$$y_{it} = \begin{cases} 1 & ; y_{it}^* > 0 \\ 0 & ; y_{it}^* \leq 0 \end{cases}$$

y_{it}^* مربوط به سن ۹ سالگی و y_{it} مربوط به سن ۱۳ سالگی است. در ۱۳ سالگی متغیر $I_{\{age\geq 13\}}$ ، وارد مدل شده است که پارامتر منتبه به آن (c) نشان دهنده تأثیر زمان بر ابتلا یا عدم ابتلا به آسم است.

همچنین در این مدل داده‌ها نشان می‌دهند که جنس در ابتلا به آسم تأثیر دارد ولی این تأثیر چندان قوی نیست ($p = 0.43$). مقدار.

مقایسه بین مدل استقلال و عدم استقلال حاکی از اهمیت برآورد ρ است. هرچند برآوردهای پارامترها در این دو مدل به هم نزدیک‌اند ولی به طور کلی خطای استاندارد برآوردها در دو مدل متفاوت است. خطای استاندارد پارامترهایی که متغیر توصیفی متناظر با آنها با زمان، تغییر نمی‌کنند (متغیر توصیفی مانا)، مانند b در این مثال، تمایل به کم برآورد شدن دارند (کم برآورد شدن p مقدار را نتیجه می‌دهد) و خطای استاندارد پارامترهایی که متغیر توصیفی متناظر با آنها با زمان تغییر می‌کند (متغیر توصیفی ناما)، مانند c در این مثال، تمایل به بیش برآورد شدن دارند (زیاد برآورد شدن p مقدار را نتیجه می‌دهد).

مانده‌هایی که با $= 0$ به دست می‌آیند، فرض استقلال بین دو پاسخ را مدنظر قرار می‌دهند. اگر $0 \neq \rho$ مانده‌ها نیز باید به گونه‌ای تصحیح شوند که همبستگی بین پاسخها را در نظر گیرید. بررسی مانده‌ها به عنوان کار بعدی پیشنهاد می‌شود.

پارامتر مربوط به ρ همبستگی بین پاسخ در ۹ سالگی و ۱۳ سالگی است که باید برآورده شود. این پارامتر در زمان برآورد در فاصله ۱- تا ۱ محدود می‌شود. Φ تابع توزیع نرمال دو متغیره را نشان می‌دهد، که به صورت زیر تعریف می‌شود:

$$\Phi_2(q_1, q_2, \rho) = \int_{-\infty}^{q_1} \int_{-\infty}^{q_2} f(x_1, x_2, \rho) dx_1 dx_2$$

که در آن $f(x_1, x_2, \rho)$ تابع چگالی نرمال دو متغیره استاندارد شده است. از آنجا که تابع $f(x_1, x_2, \rho)$ در *S-Plus ms* در *S-Plus* مجموعه‌ای از توابع غیرخطی را برای برآوردهای پارامترها با استفاده از روش بهینه‌سازی شبیه-نیوتون [۲] کمینه می‌کند، برای برآوردهای پارامترها، منهای لگاریتم تابع درستنمایی را با استفاده از تابع *ms* کمینه کرده‌ایم.

۵. نتایج استفاده از دو مدل

برای آزمون معنی‌داری پارامترها، دو حالت $\rho = 0$ و $\rho \neq 0$ را در نظر گرفته، نتایج را در جدول (۲) خلاصه می‌کیم. همان‌طور که برآوردهای پارامترها در مدل کامل (عدم استقلال) نشان می‌دهد، داده‌ها شواهد کافی بر تأثیر زمان در ابتلا به آسم دارند ($p = 0.001$ مقدار).

جدول ۱: داده‌های مربوط به آسم

وضعیت آسم			۱۳ سالگی		
پسرها	۹ سالگی	دارد ندارد	دارد	ندارد	جمع کل
			۱۵	۵۱۴	۵۲۹
دخترها	۹ سالگی	جمع کل	۳۷	۵۲۰	۵۵۷
		دارد ندارد	۱۳	۳	۱۶
		جمع کل	۱۳	۵۶۱	۵۷۴
		جمع کل	۲۶	۵۶۴	۵۹۰

جدول ۲: برآورد پارامترها و p مقدار

پارامترها	p مقدار	مدل استقلال	p مقدار	مدل عدم استقلال
a	-۰/۹۰۰	-	-۰/۸۹۰	-
b	۰/۲۴۰	۰/۰۰۹	۰/۲۳۰	۰/۰۴۳
c	۰/۱۷۰	۰/۰۰۵۷۰	۰/۱۷۰	۰/۰۰۱
ρ	-	-	۰/۹۳۰	۰/۰۰۰

مراجع:

- [1] Agresti,A., 1990, *Categorical Data Analysis*, New York, John Wiley.
- [2] Chambers, J.M. and Hastie, T.J., 1992, *Statistical Models in S*, Chapter 10: Nonlinear models, Pacific Groves, CA: Wadsworth Brooks /Cole.
- [3] Diggle, P.J., Liang, K. and Zeger, S., 1996, *Analysis of Longitudinal Data*, Oxford Science publication.
- [4] Little, R.J. and Rubin, D., 1987, *Statistical Analysis with Missing Data*, New York, Wiley.
- [5] Long , S.J., 1997, *Regression Models for Categorical and Limited Dependent Variables*, London, SAGE.
- [6] Ronitzky, A. and Wypij, D., 1994, *A Note on the Bias of Estimation with Missing Data*, *Biometrics* 147, 87-99.

علم آمار تحقق اندیشه انسان در گشودن درهای بررسی غیر ملموس حقایق است.