

مقایسه‌ی دو رهیافت برای برخورد با متغیرهای کمکی گم‌شده در رگرسیون لوژیستیک

چاو- یینگ جوان پنگ[†] و جین ژو[‡]

[†] دانشگاه ایندیانا، بلومینگتون
[‡] شرکت فناوری ژنتیک

مترجم: شیرین گلچی و ساناز مهندسی

پژوهشکده‌ی آمار

چکیده: در طول ۲۵ سال گذشته پیشرفت‌های روش‌شناسانه‌ای در زمینه‌ی تیمار داده‌های گم‌شده صورت گرفته است. بیش‌تر مطالعه‌های انتشاریافته روی داده‌های گم‌شده در متغیرهای وابسته تحت شرایط گوناگون تمرکز داشته‌اند. مطالعه‌ی حاضر در پی آن است که این خلأ را با مقایسه‌ی دو رهیافت برای برخورد با داده‌های گم‌شده در متغیرهای کمکی رسته‌ای در رگرسیون لوژیستیک پر کند: روش امید ریاضی- ماکسیم‌سازی (EM) وزن‌ها و جانهای چندگانه (MI). داده‌های نمونه به‌صورت تصادفی از جامعه‌ای با مشخصه‌های معلوم انتخاب شده‌اند. داده‌های گم‌شده‌ی مربوط به متغیرهای کمکی تحت دو حالت شبیه‌سازی شده‌اند: گم‌شده‌ی کاملاً تصادفی و گم‌شده‌ی تصادفی با نرخ‌های گم‌شدگی متفاوت. یک مدل رگرسیونی لوژیستیک با استفاده از یکی از دو رهیافت EM یا MI بر هر نمونه برازش داده شده است. عملکرد این دو رهیافت با چهار معیار مورد مقایسه قرار گرفته است. اریبی، کارایی، پوشش و نرخ عدم پذیرش. نتیجه‌ها به‌طور کلی MI را بر

Peng, C.Y.J. and Zhu, J. (2008). Comparison of Two Approaches for Handling Missing Covariates in Logistic Regression, *Educational and Psychological Measurement*, 68, 58-77.

دریافت: ۱۳۸۸/۲/۲۵، پذیرش: ۱۳۸۸/۹/۲

EM ترجیح می‌دادند. مسائل عملی از جمله اجرا، گنجاندن متغیرهای کمکی پیوسته و اثر متقابل بین متغیرهای کمکی مورد بحث قرار گرفته‌اند.

واژگان کلیدی: رگرسیون لوژستیک؛ داده‌های گم‌شده؛ متغیرهای کمکی رسته‌ای گم‌شده؛ روش وزن‌ها؛ جانهای چندگانه؛ الگوریتم EM.

۱- مقدمه

وجود داده‌های گم‌شده در مطالعه‌های روان‌شناختی و آموزشی بیش از آن که استثناء باشد قاعده است. برای مثال، نوشتگان گوناگون نشان می‌دهند که یک نرخ ریزش ۲۰٪ به‌عنوان معیاری برای مطالعه‌های مربوط به جوانی، برنامه‌های مدرسه-مبنا، و پژوهش‌های بالینی در نظر گرفته می‌شود ([۸]، [۱۱]، [۱۵]). دلیل‌های بسیاری در گم شدن داده‌ها در پروژه‌های تحقیقاتی سهم دارند، از جمله ریزش واحدها در پیش‌آزمون-پس‌آزمون یا طرح‌های طولی، دستورالعمل‌های مبهم یا سؤال‌های مداخله‌گرانه در یک آمارگیری و داده‌دهی ضعیف یا ثبت ناقص آمار ثبتي در داده‌های بایگانی‌شده، که مورد اخیر نشان داده است که اربیبی پنهان در داده‌ها وارد می‌کند [۵]. دلیل (های) داده‌های گم‌شده هرچه باشد، اثر آن روی تحقیقات کمی یک نگرانی بزرگ برای روش‌شناسان بوده است.

در طول ۲۵ سال گذشته پیشرفت‌هایی در زمینه‌ی نظریه‌های مربوط به روش‌های داده‌های گم‌شده و کاربردهای آن در مجموعه‌های داده‌های واقعی صورت گرفته است ([۱۳]، [۲۴]). بیش‌تر مطالعه‌های انتشاریافته بر داده‌های گم‌شده در متغیرهای وابسته تحت حالت‌های مختلف تمرکز داشته‌اند ([۱۲]، [۲۱])، هرچند روش‌های تیمار مقادیر گم‌شده در متغیرهای مستقل (یا متغیرهای کمکی) تفاوت اساسی با روش‌های تیمار مقادیر گم‌شده در متغیرهای وابسته ندارند. تحت فرض‌های آماری حتمی، می‌توان توسط روش ساده‌ای چون حذف موردی به شکلی کارآمد با داده‌های گم‌شده در یک متغیر مستقل برخورد کرد [۷].

ابراهیم [۱۰] یک روش امید ریاضی-ماکسیم‌سازی (EM) وزن‌ها را برای به دست آوردن برآوردهای ماکسیم درستمایی (ML) ضریب‌های رگرسیونی برای مدل‌های خطی تعمیم‌یافته با متغیرهای کمکی رسته‌ای پیشنهاد کرد. هورتون و لرد [۹] روش ابراهیم را داخل یک الگوریتم اجرا کردند و آن را در چندین مثال سلامت عمومی نشان دادند.

مقاله‌ی آن‌ها بسط‌ها و محدودیت‌های روش ابراهیم را بدون مقابله با یک روش دیگر مانند جانهی چندگانه (MI) مورد بحث قرار می‌داد. لیتل [۱۲] روش‌های متفاوت برخورد با مقادیر گم‌شده در متغیرهای کمکی را در تحلیل رگرسیونی مورد بررسی قرار داد و نتیجه گرفت که روش‌های ارجح روش‌های برآورد مدل-مبنا هستند (یعنی، روش‌های ML). شیفر و گراهام ([۲۳]، ص ۱۷۳) هر دوی روش‌های مبتنی بر درست‌نمایی و روش‌های پارامتری MI را به‌عنوان نمایانگر آخرین تحولات عملی تحت فرض «گم‌شده‌ی تصادفی» در نظر می‌گیرند. در همان مقاله شیفر و گراهام پیشنهاد می‌کنند که پژوهش‌های پیش‌تری برای مقایسه‌ی رسمی عملکرد MI با روش‌های وزن‌دهی مانند روش EM ابراهیم انجام شود. دیدلز [۴] برآوردهای ML ضریب‌های رگرسیونی لوژستیک با مقادیر گم‌شده در متغیرهای کمکی را هنگامی که توزیع این متغیرها نامشخص است مورد بررسی قرار داد. کار او چنین نتیجه‌گیری می‌کرد که این رهیافت پارامتری در صورتی که توزیع مفروض شدیداً از توزیع واقعی انحراف داشته باشد می‌تواند به اریبی جدی منجر شود.

راگوناتان [۱۷] بزرگی اریبی حاصل از استفاده از (الف) یک روش حذف فهرست‌وار، (ب) یک فن وزن‌دهی و (پ) یک MI برای برخورد با داده‌های ناقص در تحلیل رگرسیونی لوژستیک را مورد مقایسه قرار داد. نتیجه‌های شبیه‌سازی مبتنی بر ۲۵۰۰ نمونه به اندازه‌ی ۱۰۰۰ نشان می‌داد که همان‌گونه که انتظار می‌رفت اریبی روش حذف فهرست‌وار از همه بیشتر بود. روش وزنی و MI هر دو برآوردهایی تولید کردند که توزیع‌های نمونه‌گیری آن‌ها روی مقدار واقعی جامعه، ۵/۰، متمرکز بود. راگوناتان با توجه به این‌که (الف) روش وزنی رهیافت ساده‌ای برای تصحیح اریبی بود، اما به اندازه‌ی MI کارا نبود؛ (ب) MI کاربردی‌تر و کاراتر از روش وزنی بود؛ و (پ) هر دو روش تنها تحت فرض گم‌شده‌ی تصادفی معتبر بودند، روی این نتیجه‌ها اظهار نظر کرد. او روش ML را به‌عنوان رهیافتی دیگر برای مجموعه‌ی محدودی از مدل‌ها، مانند مدل‌های رگرسیونی لوژستیک، مورد بحث قرار داد و به غیر عملی بودن به کارگیری روش ML به خاطر دشواری‌های فنی آن اشاره کرد. متأسفانه او در عمل روش ML را برای داده‌های شبیه‌سازی شده به کار نگرفت؛ لذا در مطالعه‌ی وی هیچ مقایسه‌ای بین ML و MI صورت نگرفته است.

مطالعه‌ی حاضر در پی آن است که خلأ موجود در نوشتگان را با مقایسه‌ی مستقیم و رسمی رهیافت MI با روش وزنی EM (یک رهیافت ML) پر کند. به‌طور خاص این دو

رهيافت براي برخورد با داده‌هاي گم‌شده در متغيرهاي كمكي رسته‌اي در رگرسيون لوژستيک مورد مقايسه قرار گرفته‌اند. در ابتدا داده‌هاي تصادفي در جامعه‌اي با مشخصه‌هاي معلوم انتخاب شدند. سپس داده‌هاي گم‌شده در متغيرهاي كمكي تحت دو حالت: گم‌شده‌ي کاملاً تصادفي (MCAR) و گم‌شدگي تصادفي (MAR) با نرخ‌هاي گم‌شدگي متفاوت شبیه‌سازي شدند. یک مدل رگرسيوني لوژستيک با استفاده از یکی از دو رهيافت EM يا MI براي برخورد با داده‌هاي گم‌شده در متغيرهاي كمكي بر هر نمونه برازش داده شده است. برآورد پارامترها و فاصله‌هاي اطمینان آن‌ها تحت هر رهيافت محاسبه شده است. عملکرد این دو رهيافت روی چهار معيار: اريبي، کارايي، پوشش و نرخ عدم پذيرش ارزشيابي شده است. چون تحليل موردی کامل (CC) تحت حالت MCAR معتبر است ([۱۲] و [۱۳])، نتیجه‌هاي آن ارايه شده و به‌عنوان مرجعی مورد استفاده قرار گرفته است که این دو رهيافت نسبت به آن مقايسه شده‌اند. یافته‌هاي حاصل از این مطالعه‌ي شبیه‌سازي می‌تواند کار هورتون و لرد [۹] را در جهت استفاده از رهيافت MI علاوه بر روش وزني EM براي تیمار متغيرهاي كمكي رسته‌اي گم‌شده در رگرسيون لوژستيک بسط دهد. یافته‌ها نتیجه‌هاي گراهام و شيفر [۷] را نیز با بررسی عملکرد MI براي هر دو مجموعه‌ي داده‌هاي بزرگ و کوچک با مقادير گم‌شده در متغيرهاي كمكي رسته‌اي، مورد تمجيد قرار می‌دهند.

۲- اصطلاح‌شناسي و روش

در این بخش دو عبارت را که به فراواني در نوشتگان داده‌هاي گم‌شده مورد استفاده قرار می‌گیرند تعريف می‌کنيم: MCAR و MAR. این‌ها به مکانيزمی که به داده‌هاي گم‌شده منجر شده است اشاره دارند [۱۳]. دو رهيافت داده‌هاي گم‌شده (روش وزني EM و MI) نیز به‌همین ترتیب تعريف شده و مورد بررسی قرار گرفته‌اند. این چهار عبارت برای درک مطالعه‌ي حاضر از اهمیت اساسی برخوردارند.

MCAR

همان‌طور که توسط لیتل و روبین [۱۳] تعريف شده است، «اگر احتمال یک پاسخ نه به مقدار مشاهده‌شده بستگی داشته باشد و نه به مقدار گم‌شده‌اي که می‌توانست گردآوری یا ثبت شود، داده‌هاي گم‌شده کاملاً تصادفي هستند» (ص ۱۴). به‌نحوی

رسمی‌تر فرض کنید Y نمایانگر یک ماتریس داده‌ی $n \times p$ تشکیل شده از دو بخش باشد: داده‌هایی که کاملاً مشاهده شده‌اند و داده‌هایی که به‌صورت بالقوه گم شده‌اند. به عبارت دیگر $Y = (Y_{\text{observed}}, Y_{\text{missing}})$. فرض کنید R یک ماتریس $n \times p$ از مقادیر صفر و ۱ متناظر با Y مشاهده‌شده یا گم‌شده باشد. روبین [۱۸] MCAR را به‌صورت زیر تعریف کرده است:

$$(۱) \quad \text{Probability}(R | Y_{\text{missing}}, Y_{\text{observed}}, \zeta) = \text{Probability}(R | \zeta)$$

به عبارت دیگر اگر احتمال این‌که Y گم‌شده باشد به خود مقدار گم‌شده، یا مقادیر مشاهده‌شده‌ی Y بستگی نداشته، بلکه در عوض به برخی پارامترهای نامعلوم ζ بستگی داشته باشد، آنگاه گفته می‌شود که داده‌ها MCAR هستند. تحت حالت MCAR، داده‌های گم‌شده را می‌توان به‌شکل یک زیرنمونه‌ی تصادفی از داده‌های بالقوه‌ی کامل، تیمار کرد؛ لذا رهیافت CC یا حذف فهرست‌وار، در این حالت معتبر است. به بیان رسمی، مکانیزم گم‌شدن داده‌ها که دلیل‌های گم شدن داده‌ها را در بر می‌گیرد می‌تواند در استنباط مبتنی بر نمونه‌گیری و درست‌نمایی تحت MCAR نادیده گرفته شود ([۱۳]، ص ۱۵).

MAR

فرض MAR اظهار می‌کند که:

$$(۲) \quad \text{Probability}(R | Y_{\text{missing}}, Y_{\text{observed}}, \zeta) = \text{Probability}(R | Y_{\text{observed}}, \zeta)$$

معادله‌ی (۲) حاکی از آن است که احتمال شرطی گم‌شده بودن Y به شرط هر دوی Y_{observed} و Y_{missing} برابر است با احتمال شرطی مقادیر گم‌شده‌ی Y به شرط مقادیر مشاهده‌شده‌ی Y . به بیان دیگر، MAR به این معنی است که احتمال گم‌شدن یک مشاهده ممکن است به مقادیر مشاهده‌شده بستگی داشته اما به خود مقدار گم‌شده بستگی نداشته باشد.

محدودیت MAR از MCAR کم‌تر است؛ لذا از MCAR به‌عنوان حالت خاص MAR یاد می‌شود. تحت حالت MAR، مکانیزم گم‌شدن برای استنباط مبتنی بر درست‌نمایی قابل چشم‌پوشی است ([۱۳]، ص ۱۵). اساساً قابل چشم‌پوشی به این معنی

است که لازم نیست مکانیزم گم شدن داده‌ها به صورت بخشی از فرایند برآورد مدل‌سازی شود. دو رهیافت مورد بررسی در این مقاله تحت هر دو حالت MCAR یا MAR کاربرد دارند و در نوشتگان نشان داده شده است که هر دوی این رهیافت‌ها تحت حالت گم شدن غیر تصادفی بهتر از رهیافت CC یعنی حذف فهرست‌وار، عمل می‌کنند.

روش وزنی EM

ابراهیم [۱۰] یک روش برآورد ML برای مدل‌های رگرسیونی خطی تعمیم‌یافته با متغیرهای کمکی رسته‌ای گم شده پیشنهاد کرده است. این کاربردی از الگوریتم EM همراه با روش وزنی است. الگوریتم EM یک رهیافت اصولی است که دو مرحله را به صورت متوالی تکرار می‌کند: مرحله E و مرحله M ([۳] و [۱۳]). در مرحله E، امید ریاضی لگاریتم درست‌نمایی داده‌های کامل نسبت به توزیع شرطی داده‌های گم شده به شرط داده‌های مشاهده شده و برآورد پارامترهای به دست آمده از تکرار قبل گرفته می‌شود. در مرحله M، برآوردهای جدید پارامتر توسط ماکسیم کردن لگاریتم درست‌نمایی مرحله E به دست می‌آیند، و برآوردهای جدید پارامتر در تکرار بعدی مرحله E مورد استفاده قرار می‌گیرند. الگوریتم EM تا زمانی که درست‌نمایی‌های مشاهده شده تولید شده در دو تکرار متوالی تقریباً معادل باشند ادامه پیدا می‌کند. در تکرار آخر که الگوریتم همگرا می‌شود، پارامترهای برآورد شده نهایی می‌شوند.

الگوریتم EM می‌تواند با وزن‌دهی مکرر مشاهده‌ها، مانند روش ابراهیم [۱۰]، اجرا شود. به خصوص وزن‌ها در مرحله E به منظور محاسبه لگاریتم درست‌نمایی با مقادیر گم شده در متغیرهای کمکی داخل یک مجموعه داده‌های ناقص قرار داده می‌شوند. فرض کنید یک مجموعه داده‌های ناقص دارای یک متغیر وابسته (Y) و سه متغیر کمکی دو حالتی (X_1, X_2, X_3) است، که در سمت چپ جدول ۱ آرایه شده‌اند. این مجموعه داده‌های ناقص دارای دو الگوی داده‌های گم شده است: یکی فاقد مقدار X_1 و دیگری فاقد مقادیر X_1 تا X_3 است. برای این دو الگوی داده‌های گم شده، وزن‌هایی برای افزوده شدن به مجموعه داده‌ها درج می‌شوند (سمت راست جدول ۱ را ببینید). اگر مشاهده‌ای هیچ مقدار گم شده‌ای نداشته باشد، وزن منتسب شده برابر ۱ است، مانند حالت‌های ۱، ۲، ۴، ۵ و ۶. در صورتی که مشاهده‌ای دارای یک مقدار گم شده باشد (مانند حالت‌های ۳ و ۷) وزن به صورت یک احتمال شرطی برآورد می‌شود. بنا بر این،

جدول ۱- تصویری از روش امید ریاضی- ماکسیم‌سازی وزن‌ها توسط هورتون و لرد

مجموعه‌ی داده‌های اصلی الف					مجموعه‌ی داده‌های افزوده						
حالت	Y	X _۱	X _۲	X _۳	حالت	Y	X _۱	X _۲	X _۳	w _i	
۱	۰	۰	۰	۰	۱	۰	۰	۰	۰	۱	
۲	۰	۰	۰	۱	۲	۰	۰	۰	۱	۱	
۳	۰	—	۰	۰	→ {	۳	۰	۰	۰	۰	w _{۳۱}
						۳	۰	۱	۰	۰	w _{۳۲}
۴	۰	۰	۱	۰	۴	۰	۰	۱	۰	1	
۵	۱	۱	۰	۱	۵	۱	۱	۰	۱	۱	
۶	۱	۰	۰	۰	۶	۱	۰	۰	۰	۱	
۷	۱	—	—	—	→ {	۷	۱	۰	۰	۰	w _{۷۱}
						۷	۱	۰	۰	۱	w _{۷۲}
						۷	۱	۰	۱	۰	w _{۷۳}
						۷	۱	۰	۱	۱	w _{۷۴}
						۷	۱	۱	۰	۰	w _{۷۵}
						۷	۱	۱	۰	۱	w _{۷۶}
						۷	۱	۱	۱	۰	w _{۷۷}
						۷	۱	۱	۱	۱	w _{۷۸}

الف. داده‌های گم‌شده با خط تیره نشان داده شده‌اند.

$$w_{۳۱} = P(x_1 = 0 | x_2 = x_3 = y = 0, \hat{\theta})$$

و

$$w_{۳۲} = P(x_1 = 1 | x_2 = x_3 = y = 0, \hat{\theta})$$

که در آن $\hat{\theta}$ براورد پارامتر به دست آمده از مرحله‌ی M قبلی است. وزن‌های منتسب به یک مشاهده باید تابع این قید باشند که جمع آن‌ها برابر ۱ شود. لذا $w_{۳۱} + w_{۳۲} = 1$ و

$$\sum_{i=1}^n w_i y_i = 1$$

لگاریتم درست‌نمایی مجموعه‌ی داده‌های افزوده سپس در مرحله‌ی بعدی E

تعیین می‌شود. الگوریتم EM به کار گرفته شده روی مجموعه‌ی داده‌های افزوده، بین مرحله‌های E و M تکرار می‌شود تا زمانی که براوردهای پارامترهای مورد نظر نهایی شوند.

روش EM وزن‌های ابراهیم به‌طور گسترده توسط هورتون و لرد بررسی و نشان داده شده است [۹]. همان‌طور که توسط هورتون و لرد (۱۹۹۸، ص ۴۰) نشان داده شده است. خطاهای استاندارد نمونه‌ی بزرگ برآوردهای پارامتر می‌توانند با استفاده از روش لوییس [۱۴] محاسبه شوند. اجرای روش وزنی EM ابراهیم توسط هورتون و لرد برای مطالعه‌ی حاضر تصحیح و تطبیق داده شده است.

MI

MI یک فن مبتنی بر شبیه‌سازی است که مقادیر گم‌شده را چندین بار جانمی می‌کند [۱۹]. برخلاف روش جانمی تکي، MI یک معیار عدم قطعیت فراهم می‌کند که به‌وسیله‌ی آن داده‌های گم‌شده از روی داده‌های مشاهده‌شده پیش‌بینی می‌شوند. به خاطر سادگی و کلیت آن، MI به‌صورت یک روش محبوب در میان محققان علوم اجتماعی در آمده است. برای مقدمه‌ای بر MI مقاله‌ی [۲۲] را توصیه می‌کنیم.

ایده‌ی اولیه‌ی MI پر کردن مقادیر گم‌شده به دفعات متعدد به‌منظور ساختن مجموعه‌های داده‌های کامل چندگانه است. از لحاظ تحلیلی، MI شامل سه مرحله است: جانمی، تحلیل و ادغام. در مرحله‌ی اول، جانمی، هر مقدار گم‌شده توسط $m > 1$ مقدار شبیه‌سازی شده جایگزین می‌شود. در فرایند پر کردن، یک توزیع توأم از داده‌های کامل (مشاهده‌شده و گم‌شده) و یک توزیع پیشین از پارامترها برای الگوریتم داده‌افزایی به‌منظور شبیه‌سازی عددهای تصادفی برای داده‌های گم‌شده، فرض می‌شوند. برای مجموعه‌های داده‌های دارای متغیرهای رسته‌ای، یک توزیع چندجمله‌ای برای داده‌های کامل و یک توزیع پیشین دیریکله برای پارامترها فرض می‌شود. تحت فرض قابل چشم‌پوشی MAR، می‌توان m عدد تصادفی مستقل از توزیع شرطی ایستای داده‌های گم‌شده، به شرط داده‌های مشاهده‌شده، مطابق رهیافت بیزی، شبیه‌سازی کرد. بعد از مرحله‌ی جانمی، m مجموعه‌ی داده‌های کامل ساخته می‌شود. در مرحله‌ی دوم، تحلیل، هر یک از m مجموعه‌ی داده‌های کامل توسط یک روش داده‌های کامل استاندارد، مانند تحلیل رگرسیونی معمولی کم‌ترین توان‌های دوم، تحلیل می‌شود. سرانجام در مرحله‌ی سوم، ادغام، نتایج m تحلیل با استفاده از قاعده‌ی رابین (۱۹۸۷، نشان داده‌شده در [۱۶]) ادغام می‌شوند تا یک نتیجه‌ی نهایی، مثلاً یک برآورد فاصله‌ای از ضرایب رگرسیون، یک $-p$ مقدار از آزمون فرض صفر، یا یک آماره‌ی آزمون نسبت درست‌نمایی به دست آید.

تعداد مجموعه‌های داده‌های جان‌هی شده لازم نیست خیلی زیاد باشد. معمولاً ۳ تا ۱۰ مجموعه‌ی داده‌های جان‌هی شده نتایج مطلوبی ارائه می‌دهند ([۱۶]، [۲۱]، [۲۲] و [۲۴]).

طرح

هدف اولیه‌ی این مطالعه مقایسه‌ی روش وزنی ابراهیم (۱۹۹۰) (یک رهیافت EM) با رهیافت MI برای برخورد با داده‌های گم‌شده در متغیرهای کمکی رسته‌ای مدل‌های رگرسیونی لوژستیک بوده است. برای مطالعه‌ی عملکرد EM و MI پنج مرحله طراحی و اجرا شده است. در مرحله‌ی اول، یک جامعه‌ی تجربی تعریف شده است. یک مدل رگرسیون لوژستیک به داده‌های جامعه برازش داده شده و لذا ضرایب رگرسیون لوژستیک حقیقی تولید شده‌اند. در مرحله‌ی دوم، ۱۰۰۰ نمونه‌ی تصادفی به دو اندازه از جامعه انتخاب شده‌اند. در مرحله‌ی سوم، داده‌های گم‌شده در نمونه تحت شرایط MCAR و MAR با نرخ‌های گم‌شدگی مختلف به وجود آمده‌اند. در مرحله‌ی چهارم، همان مدل رگرسیون لوژستیک برازش داده شده به جامعه برای هر نمونه به کار رفته است؛ ضریب‌های رگرسیونی و خطاهای استاندارد آن‌ها با استفاده از یکی از دو رهیافت EM یا MI برای تیمار داده‌های گم‌شده در نمونه، برآورد شده‌اند و در مرحله‌ی پنجم، نتیجه‌ها برای ۱۰۰۰ نمونه‌ی شبیه‌سازی شده تحت هر حالت گردآوری شده است. عملکرد EM و MI روی چهار معیار ارزیابی، کارایی، پوشش و نرخ‌های عدم پذیرش ارزیابی شده است. به علاوه، نتیجه‌های تحلیل CC ارائه شده و به عنوان مرجع‌هایی به کار رفته‌اند که این دو رهیافت نسبت به آن‌ها مقایسه شده‌اند.

جامعه

مجموعه داده‌هایی که شامل اطلاعاتی پیرامون رضایت مشتریان از خدمات یک کارت اعتباری به خصوص می‌شد به عنوان جامعه به کار رفته است. داده‌ها شامل ۲۱۵۰۴ مشتری (یا صاحب کارت اعتباری) با اطلاعات کامل در مورد ۲۰ متغیر است. دو نوع متغیر جمع‌آوری شده است: (a) پاسخ‌های مشتریان به سؤال‌های آمارگیری (b) متغیرهای متعلق به وجهه‌ی اعتباری شخصی. سؤال‌های آمارگیری اطلاعاتی در مورد استفاده‌ی مشتری از کارت اعتباری می‌پرسیدند، مانند «در سال گذشته چند بار برای

کارت اعتباری درخواست کردید؟». نوع دوم متغیرها در جستجوی اطلاعاتی درباره‌ی وجهی اعتباری شخصی بودند، از قبیل محدودیت اولیه‌ی اعتبار. ما چهار متغیر کمکی $(X_4 - X_1)$ را که ممکن است روی رضایت مشتریان از خدمات کارت اعتباری، یعنی متغیر وابسته‌ی Y ، تأثیر بگذارد انتخاب کردیم. برای متغیرهای Y ، X_1 تا X_4 ، پاسخ مورد نظر با مقدار ۱ و پاسخ دیگر با صفر کدبندی شده است. به‌علاوه یک متغیر کمکی پنجم، X_5 ، کاملاً نامرتب با Y ، X_1 تا X_4 به‌صورت تصادفی از توزیع برنولی با احتمال (P) ، $0/2$ یا $0/4$ تولید شده است. این متغیر مستقل پنجم برای ارزیابی عملکرد الگوریتم‌های EM و MI و نسبت به نرخ خطای نوع I مورد استفاده قرار گرفته است. انتخاب این دو پارامتر برنولی توجیه نظری ندارد، اگرچه $0/2$ افراطی تر از $0/4$ است.

برای به دست آوردن پارامترهای واقعی جامعه، یک مدل رگرسیون لوژستیک با ۵ متغیر کمکی به داده‌های جامعه برازش داده شد. همه‌ی ضریب‌های رگرسیونی در سطح $\alpha = 0/5$ یا کوچک‌تر از نظر آماری معنی‌دار بودند، به جز β_5 که امری قابل درک است. ضریب‌های همبستگی پی‌یرسون بین متغیر وابسته و ۵ متغیر کمکی نشان داد که رضایت مشتریان از یک کارت اعتباری به خصوص (Y) دارای همبستگی مثبت با حد اولیه‌ی اعتبار (X_4) ، و همبستگی منفی با تعداد دفعاتی که در سال گذشته برای کارت اعتباری درخواست داده‌اند (X_1) و تعداد سال‌های داشتن کارت اعتباری (X_2) و مجموع حد اعتبار همه‌ی کارت‌های اعتباری آن‌ها (X_3) است. رضایت مشتریان به متغیر کمکی تصادفی برنولی ارتباط نداشت، که نباید داشته باشد. نتیجه‌های رگرسیون لوژستیک و ضریب‌های همبستگی پی‌یرسون را می‌توان از مؤلف اول به‌دست آورد.

۳- روش نمونه‌گیری

به‌محض این که جامعه مشخص شد، ۱۰۰۰ نمونه‌ی تصادفی به اندازه‌ی ۲۰۰ یا $n = 400$ انتخاب شدند. نرخ‌های گم‌شدگی و مکانیزم‌ها برای هر نمونه‌ی تصادفی دست‌کاری شدند. از آن‌جا که بیش‌تر ارزیابی‌های MI بر اساس نمونه‌های بزرگ بوده (به‌عنوان مثال، [۶]) و نظریه‌ی آماری حول MI و روش وزنی EM تا اندازه‌ای بر اساس تقریب‌های نمونه‌ی بزرگ است، دو اندازه‌ی نمونه تعیین شدند تا به پیروی از طرح مطالعه‌ی [۷] نمایانگر نمونه‌های بزرگ و کوچک باشند. اندازه‌های نمونه‌ی کوچک‌تر از

۲۰۰ با توجه به نرخ‌های گم‌شدگی و مکانیزم‌های گم‌شدن انجام‌شده در این مطالعه میسر نبودند.

نرخ‌های گم‌شدگی و مکانیزم‌های گم‌شدن

مقادیر گم‌شده در X_3 طبق یکی از دو مکانیزم MCAR یا MAR تولید شدند. ۴ متغیر کمکی دیگر دست نخورده نگاه داشته شدند. برای مکانیزم MCAR متغیر کمکی X_3 دارای نرخ گم‌شدگی ۴۰٪ در هر نمونه بود. این حالت از این پس (MCAR(۰/۴) خوانده می‌شود. برای مکانیزم MAR دو نرخ گم‌شدگی ایجاد شد. اولی ۱۰٪ داده‌ی گم‌شده در X_3 تولید می‌کرد اگر متغیر وابسته‌ی Y کد ۱ گرفته بود و ۲۰٪ اگر Y کد صفر گرفته بود. این الگوی داده‌های گم‌شده از این پس (MAR(۰/۱,۰/۲) خوانده می‌شود. نرخ گم‌شدگی دوم در X_3 بیشتر بود، یعنی ۳۰٪ و ۴۰٪ به ترتیب به ازای Y با کد ۱ یا صفر. این حالت (MAR(۰/۳,۰/۴) خوانده می‌شود. ترکیب دو اندازه‌ی نمونه، سه الگوی گم‌شدگی، و دو احتمال برای متغیر تصادفی برنولی (X_5)، ۱۲ حالت به دست می‌دهد. هر حالت ۱۰۰۰ بار شبیه‌سازی شده است. هر نمونه‌ی شبیه‌سازی‌شده توسط هر دو روش EM و MI تحلیل شده‌اند تا برآورد پارامترها و خطاهای استاندارد آن‌ها به دست آید.

۴- تحلیل و جانمایی

رهیافت EM با استفاده از برنامه‌ی S-Plus نوشته شده توسط هورتون و لرد [۹] برای هر نمونه‌ی شبیه‌سازی‌شده به کار رفته است. این برنامه روش وزن‌های EM را برای به دست آوردن برآوردهای ML ضریب‌های رگرسیون لوژستیک با ۵ متغیر کمکی X_1 تا X_5 اجرا می‌کند. ضریب‌های رگرسیونی برآورد شده و خطاهای استاندارد آن‌ها به دست آمده و بین ۱۰۰۰ نمونه‌ی شبیه‌سازی‌شده همگردانی شده است.

رهیافت MI با استفاده از برنامه‌ی CAT شیفر، که در S-Plus نیز نوشته شده (قابل دسترسی از نشانی <http://www.stat.psu.edu/~jls/misoftwa.html>) به کار گرفته شده است. پیش از مدل‌بندی رگرسیون لوژستیک، مقادیر گم‌شده در هر نمونه $m = 10$ بار جانمایی شده‌اند. جانمایی بر اساس یک مدل چندجمله‌ای برای X_1 تا X_5 و Y بوده است. ده مجموعه از ضریب‌ها برآورد شده و خطاهای استاندارد آن‌ها با توجه

به قاعده‌ی رابین [۱۹] ترکیب شده‌اند. انتخاب $m = 10$ مناسب به نظر می‌رسد زیرا کارایی برآورد پارامترها بر اساس m جانهی عبارت است از $\left(1 + \frac{\gamma}{m}\right)^{-1}$ ، که در آن γ نرخ اطلاعات گم‌شده است ([۱۹]، ص ۱۱۴) برای نرخ‌های پایین داده‌های گم‌شده (۲۰% یا کم‌تر)، بیش‌تر از $m = 3$ جانهی لازم نیست تا حد اقل ۹۴% کارا باشد. نتیجه‌های مدل‌بندی رگرسیون لوژیستیک با استفاده از هر دو رهیافت EM و MI برای ۱۰۰۰ نمونه‌ی شبیه‌سازی شده تحت هر یک از ۱۲ حالت ثبت شده است. نتیجه‌ها شامل برآوردهای ضریب‌های رگرسیونی، فاصله‌های اطمینان ۹۵% برآوردها و آزمون t برآوردها در سطح $\alpha = 0.05$ هستند.

۵- معیار عملکرد

عملکرد رهیافت‌های EM و MI از طریق ۴ معیار ارزیابی شد: اریبی، کارایی، پوشش، و نرخ عدم پذیرش. این چهار معیار [۷] برای ارزیابی عملکرد MI در مجموعه‌های کوچک داده‌های چندمتغیره و [۲۳] برای بحث عملکرد نسبی روش‌های گوناگون داده‌های گم‌شده شامل حذف فهرست‌وار، جانهی‌های تکی، MI و ML به کار رفته‌اند. این چهار معیار به‌صورت عملیاتی برای این مطالعه به‌صورت زیر تعریف شده‌اند:

اریبی: اریبی به‌صورت تفاضل بین میانگین ۱۰۰۰ برآورد ضریب و ضریب جامعه‌ی متناظر تعریف شده است. یک رهیافت بهتر مقدار جامعه را به‌طور متوسط با اریبی کم‌تری تولید می‌کند.

کارایی: کارایی به‌صورت تغییرپذیری برآوردها حول ضریب واقعی جامعه تعریف شده است. کارایی در این مطالعه توسط پهنای متوسط فاصله‌ی اطمینان ۹۵% اندازه‌گیری می‌شود. فاصله‌ی ۹۵% تقریباً ۴ برابر مقدار خطای استاندارد است. لذا، فاصله‌ی پهن‌تر حاکی از رهیافتی با کارایی کم‌تر است.

پوشش: پوشش به‌صورت درصد فاصله‌های اطمینان ۹۵% است که شامل ضریب واقعی جامعه در بین ۱۰۰۰ فاصله‌ی اطمینان این چنین است. یک رهیافت بهتر

باید فاصله‌های اطمینان ۹۵ درصدی تولید کند که در حدود ۹۵٪ اوقات شامل پارامتر مزبور باشند.

نرخ عدم پذیرش: نرخ عدم پذیرش به صورت درصد دفعاتی که فرض صفر درباره‌ی ضریب رگرسیونی جامعه‌ی X_h در میان ۱۰۰۰ نمونه رد شده است، تعریف می‌شود. چون X_h عمداً به گونه‌ای تولید شده است که با چهار متغیر کمکی دیگر و Y ناهمبسته باشد، ضریب رگرسیونی واقعی جامعه برای X_h در واقع صفر است. لذا، فرض صفر ضریب رگرسیونی جامعه‌ی X_h برابر صفر باید ۵ درصد اوقات در سطح α برابر ۰/۰۵ رد شود.

زمانی که اندازه‌ی نمونه از ۲۰۰ به ۴۰۰ افزایش می‌یافت انتظار می‌رفت که کارایی EM و MI طبق تعریف آماری استاندارد «کارایی»، افزایش یابد و آزمون t مرتبط با این برآوردها پرتوان‌تر گردد.

۶- نتیجه‌ها

عملکرد EM و MI از لحاظ اریبی، کارایی، پوشش و نرخ عدم پذیرش در جدول‌های ۲ و ۳ برای پارامتر دو جمله‌ای $p = 0/2$ ارایه شده است. به خاطر تشابه یافته‌ها، نتیجه‌ها برای پارامتر دو جمله‌ای $p = 0/4$ ارایه نشده‌اند اما از مؤلف اول در دسترس‌اند. تفاوت نتیجه‌ها بین $p = 0/2$ و $p = 0/4$ در بخش ۷ (بحث) ذکر شده است. جدول ۲ نتیجه‌ها را برای $n = 400$ و سه الگوی داده‌های گم‌شده: $MAR(0/1, 0/2)$ ، $MAR(0/3, 0/4)$ و $MCAR(0/4)$ به ترتیب نشان می‌دهد. جدول ۳ نتیجه‌ها را برای $n = 200$ و همان سه مکانیزم داده‌های گم‌شده نشان می‌دهد. دو متغیر X_p و X_h در هر جدول مشخص شده‌اند زیرا X_p شامل داده‌های گم‌شده و X_h متغیر تصادفی دو جمله‌ای است.

مقایسه‌ی EM و MI در مقابل CC

اریبی: برای هر متغیر کمکی و عرض از مبدأ در جدول‌های ۲ و ۳، بزرگ‌ترین اریبی که بدترین اریبی هم هست مشخص شده است. برای هر سه رهیافت، زمانی که نرخ گم‌شدگی بالا است، در $MAR(0/3, 0/4)$ ، نسبت به زمانی که نرخ گم‌شدگی پایین است، در $MAR(0/1, 0/2)$ ، اریبی بدتر است. رهیافت CC بدترین اریبی را با بیش‌ترین

جدول ۲- ايني، کاري، پوشش و نرخ عدم پذيرش روشهاي MI، EM، CC و براي $n = 40$ و P برابري 0.05

نرخ عدم پذيرش مقدار اسمي $\alpha = 0.05$	پوشش						کاري						ايني						
	MI		EM		CC		MI		EM		CC		MI		EM		CC		
	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	واقعي	%95	
MAR(0.1,0.4)	عرض از مبدا	0.974	0.970	0.976	0.970	0.976	0.974	0.970	0.976	0.974	0.970	0.976	0.974	0.970	0.976	0.974	0.970	0.976	
	X_1	0.953	0.965	0.951	0.965	0.951	0.953	0.965	0.951	0.953	0.965	0.951	0.953	0.965	0.951	0.953	0.965	0.951	
	X_2	0.962	0.958	0.962	0.958	0.962	0.962	0.958	0.962	0.962	0.958	0.962	0.962	0.958	0.962	0.962	0.958	0.962	0.958
	X_3	0.952	0.971	0.978	0.971	0.978	0.952	0.971	0.978	0.971	0.978	0.952	0.971	0.978	0.971	0.978	0.952	0.971	0.978
	X_4	0.966	0.960	0.966	0.960	0.966	0.966	0.960	0.966	0.966	0.960	0.966	0.966	0.960	0.966	0.966	0.960	0.966	0.960
	X_5	0.953	0.971	0.956	0.971	0.956	0.953	0.971	0.956	0.953	0.971	0.956	0.953	0.971	0.956	0.953	0.971	0.956	0.953
MAR(0.3,0.4)	عرض از مبدا	0.927	0.968	0.929	0.968	0.927	0.927	0.968	0.929	0.968	0.927	0.927	0.968	0.929	0.968	0.927	0.927	0.968	
	X_1	0.927	0.969	0.926	0.969	0.927	0.927	0.969	0.926	0.969	0.927	0.927	0.969	0.926	0.969	0.927	0.927	0.969	
	X_2	0.957	0.925	0.925	0.925	0.925	0.957	0.925	0.925	0.925	0.925	0.957	0.925	0.925	0.925	0.925	0.925	0.925	
	X_3	0.944	0.968	0.984	0.968	0.984	0.944	0.968	0.984	0.968	0.984	0.944	0.968	0.984	0.968	0.984	0.944	0.968	
	X_4	0.952	0.965	0.925	0.965	0.925	0.952	0.965	0.925	0.965	0.925	0.952	0.965	0.925	0.965	0.925	0.965	0.925	
	X_5	0.927	0.974	0.955	0.974	0.955	0.927	0.974	0.955	0.974	0.955	0.927	0.974	0.955	0.974	0.955	0.927	0.974	
MCAR(0.4)	عرض از مبدا	0.957	0.974	0.929	0.974	0.929	0.957	0.974	0.929	0.974	0.929	0.957	0.974	0.929	0.974	0.929	0.974	0.929	
	X_1	0.951	0.966	0.955	0.966	0.951	0.951	0.966	0.955	0.966	0.951	0.951	0.966	0.955	0.966	0.951	0.951	0.966	
	X_2	0.964	0.975	0.969	0.975	0.964	0.964	0.975	0.969	0.975	0.964	0.964	0.975	0.969	0.975	0.964	0.964	0.975	
	X_3	0.957	0.970	0.989	0.970	0.989	0.957	0.970	0.989	0.970	0.989	0.957	0.970	0.989	0.970	0.989	0.957	0.970	
	X_4	0.949	0.965	0.925	0.965	0.925	0.949	0.965	0.925	0.965	0.925	0.949	0.965	0.925	0.965	0.925	0.965	0.925	
	X_5	0.948	0.971	0.950	0.971	0.950	0.948	0.971	0.950	0.971	0.950	0.948	0.971	0.950	0.971	0.950	0.948	0.971	

توضيح: بيست و نيم ايني و کاري براي هر برآورد از روش عدم پذيرش بزرگ نشان داده شده است. $MCAR =$ گمگشدي کاملاً تصادفي، $MAR =$ گمگشدي تصادفي، $CC =$ مورد کامل، $EM =$ ماکسيمم ساري اميد رياضي، $MI =$ جابهي چنگانه الف. نرخ عدم پذيرش به صورت درصد دفعاتي که فرض صفر دريادي ضريب گرسيني جامعي X_i در ميان 1000 نمونه رد شده است تعريف مي شود. اين نرخ بايد با مقدار اسمي $\alpha = 0.05$ مقايسه شود. بدين ترتيب نرخ عدم پذيرش به صورت بزرگ نشان داده شده است.

جدول ۲- آرتیبی، پوشش و نرخ عدم پذیرش روش‌های EM، MI و CC برای $m = 20$ و $n = 100$ و p برابری $n/2$

نرخ عدم پذیرش مقدار اسمی $\alpha = 0.05$	پوشش						کارایی								
	سطح واقعی = ۹۵٪			سطح واقعی = ۹۵٪			EM			MI			CC		
	EM	MI	CC	EM	MI	CC	EM	MI	CC	EM	MI	CC	EM	MI	CC
MAR(0/1,0/1)	عرض از مبدا	-0/0754	0/0666	0/0517	0/018	1/857	0/977	0/967	0/939	0/977	0/955	0/952	0/977	0/955	0/952
	X_1	0/0004	-0/0028	-0/0110	1/982	1/913	1/859	0/977	0/955	0/936	0/979	0/952	0/979	0/952	0/936
	X_2	-0/0405	-0/0300	-0/0654	0/308	1/926	0/500	0/980	0/983	0/952	0/980	0/983	0/980	0/983	0/952
	X_3	0/0456	-0/0797	0/0027	0/283	1/906	1/879	0/980	0/983	0/952	0/980	0/983	0/980	0/983	0/952
	X_4	-0/0256	-0/0199	-0/0377	1/955	1/996	1/829	0/976	0/957	0/951	0/976	0/957	0/976	0/957	0/951
	X_5	-0/0563	-0/0533	-0/0856	0/153	0/064	0/428	0/991	0/968	0/964	0/991	0/968	0/991	0/968	0/964
MAR(0/3,0/4)	عرض از مبدا	-0/0759	0/052	0/0267	0/222	1/788	0/976	0/967	0/930	0/976	0/967	0/930	0/976	0/967	0/930
	X_1	-0/0170	-0/0241	-0/0499	0/114	1/926	0/260	0/977	0/952	0/977	0/952	0/936	0/977	0/952	0/936
	X_2	-0/0777	-0/0551	-0/0925	0/369	1/960	0/341	0/978	0/969	0/969	0/978	0/969	0/978	0/969	0/969
	X_3	0/0281	-0/0775	0/0873	0/683	0/063	0/472	0/982	0/995	0/966	0/982	0/995	0/982	0/995	0/966
	X_4	-0/0299	0/034	-0/0579	1/984	1/900	0/120	0/969	0/961	0/963	0/969	0/961	0/969	0/961	0/963
	X_5	-0/0689	-0/0660	-0/1647	0/685	0/081	0/340	0/989	0/967	0/960	0/989	0/967	0/989	0/967	0/960
MCAR(0/4)	عرض از مبدا	-0/054	0/0107	-0/0920	0/255	1/813	0/982	0/978	0/956	0/982	0/978	0/956	0/982	0/978	0/956
	X_1	-0/0158	-0/0225	-0/0201	0/328	1/923	0/341	0/979	0/956	0/979	0/956	0/956	0/979	0/956	0/956
	X_2	-0/0571	-0/0335	-0/0843	0/533	1/955	0/504	0/982	0/967	0/967	0/982	0/967	0/982	0/967	0/967
	X_3	0/0677	-0/0997	0/0048	0/718	0/079	0/522	0/984	0/998	0/963	0/984	0/998	0/984	0/998	0/963
	X_4	-0/0103	0/0196	-0/0638	1/992	1/992	0/293	0/971	0/968	0/963	0/971	0/968	0/971	0/968	0/963
	X_5	-0/0550	-0/0537	-0/0900	0/227	0/076	0/560	0/988	0/970	0/961	0/988	0/970	0/988	0/970	0/961

توضیح: بیش‌ترین آرتیبی و کارایی برای هر پروتکل از آن‌جایی به‌صورت بزرگ نشان داده شده است. MCAR = گم‌شدگی کاملاً تصادفی؛ MAR = گم‌شدگی کاملاً تصادفی؛ MI = ماکسیمم‌سازی امید ریاضی؛ CC = موزع کامل؛ EM = مقایسه‌ی نرخ عدم پذیرش بدترین نرخ عدم پذیرش به‌صورت بزرگ. الف. نرخ عدم پذیرش به‌صورت درصد دهانه‌ای که فرض صفر دربارگی ضریب رگرسیونی جامعه‌ی X_i در میان ۱۰۰ نمونه رد شده است تعریف می‌شود. این نرخ باید با مقدار اسمی $\alpha = 0.05$ مقایسه شود. بدترین نرخ عدم پذیرش به‌صورت بزرگ نشان داده شده است.

فراوانی به دست می‌دهد. بدترین عملکرد CC روی اریبی در سه مکانیزم داده‌های گم‌شده و دو اندازه‌ی نمونه نفوذ کرده است. EM در مورد اریبی اندکی بهتر از MI عمل کرده است. برآوردهای MI ضریب‌های رگرسیونی X_p و X_h شامل اریبی کم‌تری از EM‌ها در الگوهای MAR بوده و برآوردهای ضریب‌های رگرسیونی X_p و X_h دارای اریبی کم‌تری از EM‌ها در الگوی MCAR بوده‌اند. برای ۶ حالت نشان داده شده در جدول‌های ۲ و ۳، MI اریب‌ترین برآورد را برای X_p و کم‌ترین برآورد اریبی را برای X_h تولید کرده است. عرض از مبدأ Y توسط MI درست‌تر از EM برآورد شد.

کارایی: از آن‌جا که فاصله‌ی پهن‌تر حاکی از رهیافتی با کارایی کم‌تر بوده، پهن‌ترین و در واقع بدترین فاصله‌ی ۹۵٪ مشخص شده است. MI به‌صورت یکنواخت بهترین رهیافت از لحاظ کارایی بود. CC بدترین رهیافت تحت $MAR (0/3, 0/4)$ و MCAR بوده و EM بدترین رهیافت تحت $MAR (0/1, 0/2)$ ، بدون در نظر گرفتن اندازه‌ی نمونه بوده‌اند. به‌علاوه، EM در برآورد ضریب X_p به‌طور یکنواخت رهیافتی با کم‌ترین کارایی بوده است.

پوشش: هر سه رهیافت نرخ‌های پوشش مشابهی ارائه دادند. بنا بر [۲۳]، پوشش کم‌تر از ۹۰٪ نشان‌دهنده‌ی یک پوشش جداً سطح پایین است، زیرا ۹۰٪ برابر است با دو برابر نرخ خطای اسمی (۵٪). با این توضیح، هر سه رهیافت به‌طور مساوی فواصل اطمینان خوبی برای هم‌ی متغیرهای کمکی تحت همه‌ی حالت‌ها ارائه دادند.

نرخ عدم پذیرش: این معیار تنها برای X_h دارای کاربرد است زیرا این متغیر کمکی نامرتب با X_p تا X_h و Y تولید شده و لذا ضریب واقعی آن در جامعه برابر صفر است. جدول‌های ۲ و ۳ نشان می‌دهند که EM به‌طور یکنواخت بدترین رهیافت است زیرا فرض صفر را محتاطانه‌تر از هر دو رهیافت MI و CC رد می‌کند. نرخ عدم پذیرش محتاطانه توسط EM برای اندازه‌ی نمونه‌ی ۲۰۰ بدتر از ۴۰۰ بود. زمانی که $n = 200$ ، نرخ عدم پذیرش EM به کمی ۰/۰۱ است، یعنی تنها ۲۰٪ سطح اسمی $\alpha (0/05)$. قابل درک است که همه‌ی رهیافت‌ها بدترین نرخ عدم پذیرش را برای $n = 200$ و الگوی داده‌های گم‌شده‌ی $MAR (0/1, 0/2)$ ارائه داده‌اند. هنگامی که اندازه‌ی نمونه به $n = 400$ افزایش پیدا کرد همه‌ی رهیافت‌ها نرخ‌های عدم پذیرش بهتری ارائه دادند. در $n = 400$ بهترین نرخ عدم پذیرش تحت $MAR (0/3, 0/4)$ ، سپس MCAR و پس

از آن $(0/1, 0/2)$ MAR به دست آمد، هر چند سه نرخ عدم پذیرش برای هر رهیافت بیش از ۰/۱۲ اختلاف نداشتند.

اثر مکانیزم داده‌های گم‌شده، نرخ‌های گم‌شدگی، و اندازه‌های نمونه

MAR، نرخ‌های گم‌شدگی، و اندازه‌های نمونه. با در نظر گرفتن مکانیزم‌های داده‌های گم‌شده‌ی MAR به تنهایی و اندازه‌ی نمونه‌ی بزرگ‌تر (یعنی، $n = 400$)، نتیجه‌های جدول ۲ نشان می‌دهند که اریبی برآوردهای حاصل از EM، MI و CC به‌طور کلی زمانی که نرخ گم‌شدگی زیاد بود $(0/1, 0/2)$ کم‌تر از زمانی بود که نرخ گم‌شدگی کم $(0/3, 0/4)$ بود. مورد استثناء X_4 بود که زمانی که نرخ گم‌شدگی زیادتر بود اریبی آن کم‌تر بود. همچنین، برآوردهای EM برای X_1 و X_2 زمانی که نرخ‌های گم‌شدگی بیش‌تر بود اریبی کم‌تری داشتند. EM و CC زمانی که نرخ گم‌شدگی کم بود $(0/1, 0/2)$ هر دو کارآمدتر از زمانی بودند که این نرخ زیاد $(0/3, 0/4)$ بود. اختلاف کارایی بین دو نرخ گم‌شدگی برای CC فاحش‌تر از EM است. MI تقریباً برای هر دو نرخ گم‌شدگی به یک اندازه کارآمد است.

همان‌طور که در بالا گفته شد هر سه رهیافت پوشش قابل قبولی از پارامترهای جامعه به دست دادند. به‌علاوه، به‌نظر می‌رسید که پوشش مستقل از نرخ گم‌شدگی است. نرخ‌های عدم پذیرش توسط MI و CC زمانی که نرخ گم‌شدگی زیاد بود، به سطح اسمی ۰/۰۵ نزدیک‌تر از زمانی بودند که نرخ گم‌شدگی کم بود. این روند برای EM، که بدترین و محافظه‌کارانه‌ترین نرخ عدم پذیرش را به دست می‌داد، معکوس می‌شد. برای الگوی داده‌های گم‌شده‌ی MAR و اندازه‌ی نمونه‌ی کوچک (یعنی، $n = 200$)، نتیجه‌های نشان داده‌شده در جدول ۳ مشابه یافته‌های نمونه‌ی بزرگ‌تر روی هر ۴ معیار است. چند استثناء در زیر ذکر شده است:

- اریبی در X_2 و X_4 برای EM زمانی که نرخ گم‌شدگی بالاتر است کم‌تر از زمانی است که نرخ گم‌شدگی کم‌تر است.
- اریبی در عرض از مبدأ Y برای CC در نرخ گم‌شدگی بالا کم‌تر از نرخ گم‌شدگی پایین است.

- کارایی EM در برآورد ضریب‌های X با نرخ گم‌شدگی پایین‌تر کم‌تر از زمانی است که نرخ گم‌شدگی زیادتر است.
- نرخ‌های عدم پذیرش توسط سه رهیافت برای هر دو نرخ گم‌شدگی تقریباً یکسان است.

اگر نرخ گم‌شدگی ثابت نگاه داشته شود، نتیجه‌های حاصل از نمونه‌ی بزرگ از لحاظ کارایی و نرخ عدم پذیرش بهتر از نمونه‌ی کوچک است؛ $MAR (0/1, 0/2)$ و $MAR (0/3, 0/4)$ را از جدول‌های ۲ و ۳ مقایسه کنید. پوشش توسط سه رهیافت مستقل از اندازه‌ی نمونه به نظر می‌رسد. اریبی تولید شده توسط MI و CC برای عرض از مبدأ Y با نمونه‌های کوچک کم‌تر از نمونه‌های بزرگ است. میزان اریبی در متغیرهای کمکی دیگر مستقل از اندازه‌ی نمونه برای MI و CC به نظر می‌رسد.

MCAR و اندازه‌های نمونه. تحت مکانیزم MCAR و برای هر دو اندازه‌ی نمونه‌های بزرگ و کوچک، همان‌طور که در جدول‌های ۲ و ۳ نشان داده شده است، EM برای بیش‌تر متغیرهای کمکی اریبی کم‌تری نسبت به MI و CC داشته است. MI و CC تنها زمانی که اندازه‌ی نمونه بزرگ بود نرخ‌های عدم پذیرش نزدیک‌تر به سطح اسمی $\alpha (0/05)$ ارائه دادند. زمانی که اندازه‌ی نمونه به ۲۰۰ کاهش یافت نرخ‌های عدم پذیرش توسط هر سه رهیافت بدتر شد. نرخ عدم پذیرش EM فارغ از اندازه‌ی نمونه از همه بدتر و محافظه‌کارانه‌تر بود.

مشابه یافته‌هایی که تحت MAR به دست آمد، نتیجه‌های نمونه‌های بزرگ بر مبنای MCAR بر حسب کارایی و نرخ عدم پذیرش بهتر از نتیجه‌های نمونه‌های کوچک بود. پوشش توسط هر سه رهیافت مستقل از اندازه‌ی نمونه به نظر می‌رسید. همان‌طور که قبلاً اشاره شد، هر سه رهیافت پوشش قابل قبولی برای هر دو $n = 200$ و $n = 400$ به دست دادند. اریبی در مجموع با نمونه‌های بزرگ کم‌تر از نمونه‌های کوچک بود، به استثناء اریبی تولید شده توسط MI و CC برای عرض از مبدأ Y .

نتیجه‌گیری:

بر اساس نتیجه‌های به دست آمده می‌توان به نتیجه‌گیری‌های زیر رسید:

- ۱- به نظر می‌رسد پوشش پارامتر واقعی توسط سه رهیافت مستقل از اندازه‌ی نمونه، مکانیزم داده‌های گم‌شده، و نرخ‌های گم‌شدگی است. هر سه رهیافت به‌طور یکنواخت پوشش قابل قبولی تحت همه‌ی شرایط به دست دادند که هیچ یک کم‌تر از ۹۰٪ نبود.
- ۲- MI از لحاظ کارایی تحت هر دو مکانیزم MAR و MCAR برای هر دو اندازه‌ی نمونه به‌صورت سازگار بهتر از EM عمل کرده است. زمانی که نرخ گم‌شدگی پایین بود، کارایی EM بدتر از CC بوده است. تحت همه‌ی شرایط دیگر، برآوردهای CC کم‌ترین کارایی را داشته‌اند.
- ۳- MI از نظر نرخ عدم پذیرش تحت هر دو مکانیزم MAR و MCAR برای هر دو اندازه‌ی نمونه به‌صورت سازگار بهتر از EM عمل کرده است. نرخ عدم پذیرش EM زمانی که اندازه‌ی نمونه کوچک بود، فارغ از مکانیزم گم‌شدگی به شدت محافظه‌کارانه بود. نرخ عدم پذیرش MI خیلی نزدیک به CC بود. هر دو نرخ عدم پذیرش زمانی که اندازه‌ی نمونه بزرگ بود به نرخ اسمی ۰/۰۵ نزدیک‌تر از زمانی بودند که اندازه‌ی نمونه کوچک بود.
- ۴- از نظر اریبی، رهیافت CC بدترین اریبی را با بیش‌ترین فراوانی داشته است. EM زمانی که اندازه‌ی نمونه بزرگ ($n = 400$) و مکانیزم گم‌شدگی MCAR بود بهتر از MI عمل کرد. برای شش حالت تحت مطالعه، MI اریب‌ترین برآورد را برای X_3 و برآوردی با کم‌ترین اریبی را برای X_4 به دست داد. عرض از مبدأ Y توسط MI صحیح‌تر از EM برآورد شد. میزان اریبی تولید شده توسط MI و CC برای عرض از مبدأ Y به اندازه‌ی نمونه بستگی داشت به این شکل که در نمونه‌های کوچک فارغ از مکانیزم داده‌های گم‌شده، کم‌تر بود. اریبی در دیگر متغیرهای کمکی تولید شده توسط MI و CC به اندازه‌ی نمونه و مکانیزم‌های داده‌های گم‌شده بستگی داشت.
- ۵- نتیجه‌های حاصل از نمونه‌ی بزرگ از لحاظ کارایی و نرخ عدم پذیرش، فارغ از مکانیزم گم‌شدگی داده‌ها و نرخ‌های گم‌شدگی بهتر از نتایج نمونه‌ی کوچک بود. از لحاظ اریبی و پوشش، عملکرد EM و MI هر دو به‌طور کلی مستقل از اندازه‌ی نمونه و مکانیزم داده‌های گم‌شده بود.

۶- دو الگوی گم‌شدگی MCAR (۰/۴) و MAR (۰/۳, ۰/۴) هیچ یک از دو عملکرد EM و MI را به صورت متفاوت تحت تأثیر قرار ندادند. این پدیده می‌تواند توسط این واقعیت که نرخ‌های گم‌شدگی در هر دو الگو مشابه بوده، MCAR حالت خاصی از MAR است و در هر دو الگوریتم قاعده‌ی بیز به کار رفته، توضیح داده شود.

۷- زمانی که هر چهار معیار همزمان در نظر گرفته شوند، EM و MI در حالت اندازه‌ی نمونه‌ی بزرگ ($n = 400$) و نرخ گم‌شدگی پایین در (MAR (۰/۱, ۰/۲) بهترین عملکرد را داشتند. از آنجا که هر دو رهیافت بر مبنای درست‌نمایی‌اند، این امر مورد انتظار بود.

یافته‌ها به‌طور کلی MI را بر EM ترجیح می‌دادند. MI از لحاظ کارایی و نرخ عدم پذیرش به‌طور سازگار بهتر از EM عمل می‌کرد. پوشش MI در سطح ۹۵٪ به‌طور یکنواخت از EM به ۹۵٪ نزدیک‌تر بود، به استثنای متغیر کمکی X_3 با داده‌های گم‌شده. از نظر اریبی EM کمی بهتر از MI عمل کرد. برآورد MI برای X_3 به‌طور سازگار اریب‌ترین برآورد تحت همه‌ی حالت‌های مورد بررسی بود. EM زمانی که اندازه‌ی نمونه بزرگ ($n = 400$) و پارامتر برنولی (p) کوچک (۰/۲) بود برآوردهایی با اریبی کم‌تر از MI به دست داد. زمانی که اندازه‌ی نمونه به ۲۰۰ کاهش پیدا کرد، برآوردهای MI برای X_4 (متغیر دوجمله‌ای) و عرض از مبدأ Y فارغ از مکانیزم داده‌های گم‌شده یا پارامتر برنولی کم‌تر از برآوردهای EM اریبی داشتند. از آنجا که هر دو الگوریتم به صورت مجانبی سازگارند، اریبی موضوعی جدی نیست مگر زمانی که اندازه‌ی نمونه کوچک باشد. همان‌گونه که انتظار می‌رفت، هر دو رهیافت زمانی که اندازه‌ی نمونه بزرگ ($n = 400$) و نرخ گم‌شدگی پایین (MAR (۰/۱, ۰/۲) تحت MAR بود، نتیجه‌های بهتری به دست دادند. زمانی که نرخ گم‌شدگی بالا و اندازه‌ی نمونه کوچک بود، برای EM و MI میسر نبود که اطلاعات کافی از داده‌های مشاهده‌شده به دست آورند تا داده‌های گم‌شده را پیش‌بینی کنند. در نتیجه، هر دو رهیافت، با کارایی مجانبی، با $n = 200$ عملکرد ضعیفی داشتند.

۷- بحث

یافته‌های مبتنی بر پارامتر دو جمله‌ای $0/4$ بسیار شبیه نتیجه‌های حاصل از $p = 0/2$ بودند؛ این یافته‌ها هفت نتیجه‌گیری ارائه شده در بالا را تقویت می‌کنند. کارایی و پوشش تحت تأثیر پارامترهای دو جمله‌ای قرار نداشتند. پوشش به دست آمده توسط سه رهیافت از دو پارامتر دو جمله‌ای به طور مساوی قابل قبول بود. ولی هنوز تفاوت‌های کوچکی باقی می‌ماند:

- با $n = 400$ بیش‌تر آریبی‌ها با $p = 0/4$ کمی بزرگ‌تر از $p = 0/2$ برای همه‌ی مکانیزم‌های داده‌های گم‌شده بودند.
- نرخ عدم پذیرش EM که بدترین در میان سه رهیافت است، با $p = 0/4$ به سطح اسمی $0/5$ نزدیک‌تر از $p = 0/2$ برای دو اندازه‌ی نمونه بود.

دشواری‌های محاسباتی

چندین مشکل محاسباتی در مطالعه‌ی حاضر رخ داده که بیش‌ترین آن‌ها در الگوریتم EM بوده است. اول، زمانی که $n = 200$ بود، برای برنامه‌ی S-Plus محاسبه‌ی خطاهای استاندارد در اجرای EM دشوار بود به خصوص برای X_p که شامل داده‌های گم‌شده بود. دلیل این امر آن بود که اندازه‌ی نمونه در مقایسه با تعداد پارامترها $(= 31 = 1 - 2 \times 2 \times 2 \times 2 \times 2)$ در یک مدل لوژستیک با ۵ متغیر دو حالتی، خیلی کوچک بود. به علاوه، نسبت Y ‌های برابر ۱ (ناراضی از کارت اعتباری) در جامعه تنها $0/148112$ بود. در نتیجه مشاهداتی با Y برابر ۱ در نمونه‌هایی به اندازه‌ی 200 کم بود. باور بر این است که این عوامل در دشواری EM برای محاسبه‌ی خطاهای استاندارد سهم داشته‌اند.

دومین دشواری رهیافت EM پیچیدگی محاسباتی آن بود. محققان علاقه‌مند به استفاده از روش وزنی EM باید کدهای برنامه‌نویسی خاص خود را در MATLAB، S-Plus یا زبان R خود رأساً بنویسند زیرا این الگوریتم EM در یک بسته‌ی آماری برای مقاصد عام مانند SPSS یا SAS گنجانده نشده است. نرم افزارهای رایگانی مانند LEM یا نسخه‌ی جدیدتر آن Latent Gold 4.0 می‌توانند برآوردهای ML پارامترهای رگرسیونی لوژستیک را به سادگی ارائه دهند. هرچند از وب سایت آن در ۲۵ فوریه ۲۰۰۷

(http://www.statisticalinnovations.com/products/latentgold_v4.html) مشخص نیست که آیا این برنامه‌ها می‌توانند برای اجرای روش وزن‌های ابراهیم [۱۰] مطابقت داده شوند. مشکل سوم رهیافت EM خاص بودن محاسباتی آن است. اجرای روش ابراهیم برای مدل‌هایی با متغیرهای کمکی گسسته و داده‌های گم‌شده‌ای که تنها به برآمدها و داده‌های مشاهده‌شده، یعنی مکانیزم MAR ([۹]) وابسته‌اند، آسان است. اگر مدلی شامل متغیرهای کمکی پیوسته، آمیزه‌ای از متغیرهای کمکی پیوسته و گسسته، و/یا نرخ گم‌شدگی بالا باشد، الگوریتم EM ساخته‌شده توسط هورتون و لرد باید برای مطابقت با مدل تصحیح شود. این کار به هیچ وجه آسان یا ساده نیست. سرانجام، اگر درصد داده‌های گم‌شده بالا باشد نرخ همگرایی می‌تواند به شدت کند باشد ([۲] و [۱۳]). این امر به این دلیل است که نرخ همگرایی برای روش‌های تکرار مبتنی بر درستنمایی مانند EM و روش وزن‌های ابراهیم با درصد داده‌های کاملاً مشاهده‌شده در یک مجموعه‌ی داده‌ها متناسب است ([۱۳]، ص ۱۳۰).

ولی درباره‌ی امکان تعمیم روش وزن‌های EM به متغیرهای کمکی پیوسته، هورتون و لرد ([۹]، صص ۴۳-۴۴) اظهار داشته‌اند که چنین فنونی، با یا بدون یک مدل پارامتری برای متغیرهای کمکی پیوسته موجود است؛ ولی از نظر محاسباتی پیچیده‌اند. همچنین گسترش رهیافت EM به مدل‌هایی با اثر متقابل از نظر محاسباتی پیچیده است، اگرچه غیر ممکن نیست [۲۳]. لذا کاربرد روش وزن‌های EM محدود است.

این پدیده‌ها به یک مزیت واضح رهیافت MI اشاره می‌کنند؛ یعنی محققان تنها یک برنامه‌ی S-Plus برای MI نیاز دارند. زمانی که چند مجموعه‌ی داده‌های جانمایی شده ساخته می‌شود، می‌توان هر مدل آماری را برای این مجموعه‌ی داده‌های جانمایی شده در هر نرم‌افزار آماری به کار برد. در مقابل، رهیافت EM به یک برنامه‌ی مخصوص که خاص اجرای مدل رگرسیونی لوژستیک یا یک مدل انتخابی نوشته شده باشد نیاز دارد. برنامه‌های رایانه‌ای مدل-ویژه به‌طور کلی ترجیح داده نمی‌شوند مگر این که پژوهش‌گران زمان و منابع کافی در اختیار داشته باشند. به‌علاوه، مدل جانمایی و مدل تحلیل برای MI لازم نیست به مدل‌هایی با متغیرهای کمکی دو حالتی و/یا یک متغیر برآمد دو حالتی محدود شود. این مدل‌ها می‌توانند شامل اثر متقابل بین متغیرهای کمکی نیز باشند. در مجموع، MI برای هر مدل خطی تعمیم‌یافته‌ی مفروض برای پژوهش‌های آموزشی و روان‌شناختی کاربرد دارد. لذا، اجرای MI نسبتاً ساده است. سرانجام، MI در مطالعه‌ی

حاضر و مطالعه‌های دیگر عملکرد خوبی داشته است [۷]. حتی در وضعیت‌هایی که MI بهتر از الگوریتم‌های مدل-ویژه یا مشکل-ویژه عمل نمی‌کند، هنوز هم رهیافتی راحت و کارایی برای برخورد با داده‌های گم‌شده در مجموعه‌ی داده‌های تجربی محسوب می‌شود.

مدل جانهی در MI

کلید موفقیت رهیافت MI تعیین مدل جانهی است. MI فرض می‌کند که مدل جانهی با مدل تحلیل که متعاقباً برای تحلیل داده‌ها به کار خواهد رفت یکسان است. هر چند در عمل نیازی نیست که این دو مدل یکسان باشند. بنا بر [۲۱]، مدل جانهی باید شامل پیشگوهایی باشد که دارای اهمیت ذاتی باشند گم‌شدگی و متغیرهایی با داده‌های گم‌شده را به‌خوبی پیشگویی کنند، و ویژگی‌های به خصوص طرح نمونه (به‌عنوان مثال، آمارگیری‌های احتمالاتی) را بازتاب دهند. مدل جانهی نیازی نیست که از نظر مفهومی معنی داشته باشد. در عمل، این مدل نوعاً شامل متغیرهای ذاتاً مورد علاقه و همچنین پیشگوهایی است که منعکس‌کننده‌ی گم‌شدگی پنداشته می‌شوند. شواهد تجربی نشان می‌دهند که مدل جانهی حتی اگر درست تعیین نشده باشد [۲۱] می‌تواند در کاهش اریبی مفید باشد.

تعداد جانهی‌ها

در مورد سؤال «چند جانهی برای MI مورد نیاز است؟» رابین ([۱۹]، ص ۱۱۴) نشان می‌دهد که کارایی یک برآورد پارامتر بر اساس m جانهی برابر است با $\left(1 + \frac{\gamma}{m}\right)^{-1}$ که در آن γ نرخ گم‌شدگی اطلاعات است. نرخ اطلاعات گم‌شده به‌صورت افزایش واریانس برآورد پارامتر ناشی از داده‌های گم‌شده در یک نمونه‌ی بزرگ تعریف می‌شود؛ که می‌تواند از نرخ مشاهده‌های گم‌شده کم‌تر یا بیش‌تر باشد. برای نرخ‌های پایین اطلاعات گم‌شده (۲۰٪ یا کم‌تر)، برآوردی مبتنی بر $m = 3$ جانهی دارای کارایی برابر حد اقل ۹۴٪ است. برای نرخ‌های بالاتر اطلاعات گم‌شده، جانهی‌های اضافی مورد نیاز است. شیوه‌ی MI در SAS نسخه‌ی ۸ و ۹ پیش‌فرض تعداد جانهی‌ها را برابر ۵ قرار می‌دهد (SAS Institute, 2004, [۲۰]). به‌طور کلی، توصیه می‌شود که مقادیر گم‌شده پیش از ادغام نتایج ۵ بار جانهی شوند [۱۶].

نگرانی‌های دیگر موجود در MI

علی‌رغم نقطه‌های قوت آن، نگرانی‌های بسیاری در مورد MI در نوشتگان مستندسازی شده است. اولین آن‌ها این است که هر کاربرد MI مقادیر جانمایی شده و آماره‌های مرتبط کمی متفاوت تولید می‌کند. در نتیجه، زمانی که MI برای برخورد با داده‌های گم‌شده مورد استفاده قرار می‌گیرد، همیشه نمی‌توان نتیجه‌ها را تکرار کرد. با مساوی نگه داشتن بقیه‌ی موارد، این مسئله با زیاد شدن مقدار داده‌های گم‌شده اهمیت بیش‌تری می‌یابد. دوم این‌که، آلیسون ([۱]) مشکلات MI را زمانی که برآورد اثر متقابل مد نظر است مورد بحث قرار می‌دهد. سومین نگرانی محدودیت دسترسی به نرم‌افزار MI در نرم‌افزارهای آماری معمول است. در کنار دو شیوه‌ی (MI، MIANALYZE) مورد اجرا در SAS نسخه‌ی ۸/۲ و بالاتر، نرم‌افزار رایگان MI، NORM که توسط شیفر ([۲۱]) نوشته شده است نیز می‌تواند برای اجرای MI تحت فرض نرمال چندمتغیره به کار گرفته شود. این نرم‌افزار به همراه سه بخش اضافی: CAT، MTX و PAN از طریق <http://www.Stat.Psu.edu/~jls/misotwa.html> در دسترس است. MI، CAT را برای داده‌های رسته‌ای تحت فرض توزیع لگ خطی (Log-Linear) اجرا می‌کند، MIX، MI را برای آمیزه‌ای از داده‌های پیوسته و رسته‌ای تحت فرض توزیع مکانی عام اجرا می‌کند، و PAN، MI را برای داده‌های پانلی یا داده‌های خوشه‌ای تحت یک توزیع چند متغیره‌ی خطی با اثرهای آمیخته اجرا می‌کند. با این‌که نرم‌افزار شیفر به کارگیری MI را در انواع گوناگونی از حالت‌های داده‌های گم‌شده میسر می‌سازد، پژوهشگران ممکن است هنوز این برنامه‌ها را کمتر از شیوه‌های MI و MIANALYZE در SAS کاربردوست ببینند.

پیامدهایی برای پژوهشگران

مطالعه‌ی حاضر در پی گسترش کار هورتون و لرد در [۹] با گنجاندن رهیافت MI برای برخورد با متغیرهای کمکی دو حالتی گم‌شده در مدل‌های رگرسیونی لوژستیک بوده است. این مطالعه به تمجید از کار گراهام و شیفر، [۷] با بررسی عملکرد MI برای هر دو مجموعه‌ی داده‌های بزرگ و کوچک پرداخته است. به ویژه، روش وزن‌های EM و MI از لحاظ آریبی، کارایی، پوشش و نرخ عدم پذیرش مقایسه شده‌اند. یافته‌های ما نشان می‌دهند که زمانی که داده‌های گم‌شده به‌طور تصادفی رخ می‌دهند (یعنی MAR)، MI

می‌تواند یک رهیافت ماندنی برای برخورد با متغیرهای کمکی دو حالتی گم‌شده در رگرسیون لوژستیک باشد. روش وزن‌های EM، که نیازمند یک برنامه‌ی رایانه‌ای است که ویژه‌ی هر مدل انتخابی پژوهشگر نوشته شده باشد، نه کاربردی است و نه برتر از MI تلقی می‌شود.

با کدبندی دو حالتی متغیر وابسته و هر ۴ متغیر کمکی، یافته‌های مطالعه‌ی حاضر با معیارهای روان‌شناختی و آموزشی مناسبت دارد. به ویژه، مدل‌های لوژستیک به وفور برای ربط دادن اندازه‌های برآمد دو جمله‌ای (یعنی، قبول یا رد، صحیح یا ناصحیح) با متغیرهای کمکی آن‌ها، اعم از رسته‌ای (نوع مدرسه)، پیوسته (سن)، یا آمیزه‌ای از هر دو مورد استفاده قرار می‌گیرند. MI حتی با وجود اثر متقابل می‌تواند برای برخورد مناسب با داده‌های گم‌شده در متغیرهای کمکی به کار رود. هورتون و لرد [۹] گسترشی از روش وزن‌های EM را برای پرداختن به خطاهای بدرده‌بندی به‌عنوان خطاهای اندازه‌گیری برای برآمد دو حالتی و متغیرهای کمکی در مدل‌های خطی تعمیم‌یافته، حتی با متغیرهای رده‌ای پنهان، نشان دادند. در واقع، پارامترهای مرتبط با متغیرهای پنهان می‌توانند با استفاده از MI برآورد شوند، زیرا مقادیر متغیرهای پنهان همیشه گم شده‌اند و لذا می‌توان آن‌ها را با یک رهیافت ماندگار مانند MI تیمار کرد.

پژوهش‌های آینده باید این مطالعه‌ی مقایسه‌ای را با گنجانیدن متغیرهای پیوسته و رسته‌ای هر دو در مدل‌های لوژستیک، با به کارگیری MI و EM برای داده‌های چندمتغیره‌ی گم‌شده تحت حالت‌های دیگری به جز MCAR یا MAR و با بررسی نتایج حاصل از بد مشخص کردن مدل‌های لوژستیک و/یا مدل‌های شامل اثر متقابل بین متغیرهای کمکی، گسترش دهند.

مرجع‌ها

- [1] Allison, P.D. (2001). *Missing Data* (Sage University Paper Series on Quantitative Applications in the Social Sciences, 07-136) Thousand Oaks, CA: Sage.
- [2] Buu, Y.P.A. (1999). *Analysis of Longitudinal Data with Missing Values: A Methodological Comparison*. Unpublished doctoral dissertation, Indiana University.

- [3] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977). Maximum likelihood estimation from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society*, **39**, 1-38.
- [4] Didelez, V. (2002). ML and semiparametric estimation in logistic models with incomplete covariate data. *Statistica Neerlandica*, **56**, 330-345.
- [5] Dworkin, R.J. (1987). Hidden bias in the use of archival data. *Evaluation & the Health Professions*, **10**, 173-185.
- [6] Ezatti-Rice, T.M., Johnson, W., Khare, M., Little, R.J.A., Rubin, D.B. and Schafer, J.L. (1995). A simulation study to evaluate the performance of model-based multiple imputations in NCHS health examination surveys. In *Proceedings of the Annual Research Conference* (PP. 257-266). Washington, DC: Bureau of the Census.
- [7] Graham, J.W. and Schafer, J.L. (1999). On the performance of multiple imputation for multivariate data with small sample size. In R. Hoyle (Ed), *Statistical Strategies for Small-Sample Research* (PP. 1-29). Thousand Oaks, CA: Sage.
- [8] Hall, J. (1993). *Skills Training for Pregnant and Parenting Adolescents* (NIDA Monograph No. 156, NIH Publication No. 95-3908). Rockville, MD: National Institute on Drug Abuse.
- [9] Horton, N.J. and Laird, N.M. (1998). Maximum likelihood analysis of generalized linear models with missing covariates. *Statistical Methods in Medical Research*, **8**, 37-50.
- [10] Ibrahim, J.G. (1990). Incomplete data in generalized linear models. *Journal of the American Statistical Association*, **85**, 765-769.
- [11] Kellam, S., Rebok, G., Ialongo, N. and Mayer, L. (1994). The course and malleability of aggressive behavior from early first grade into middle school: Results of a developmental epidemiologically-based preventive trial. *Journal of Child Psychology and Psychiatry*, **35**, 259-281.
- [12] Little, R.J.A. (1992). Regression with missing X's: A review. *Journal of the American Statistical Association*, **87**, 1227-1237.

- [13] Little, R.J.A., and Rubin, D.B. (1987). *Statistical Analysis with Missing Data*. Wiley, New York.
- [14] Louis T.A. (1982). Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society, Series B*, **44**, 226–233.
- [15] Mason, M.J. (1999). A Review of procedural and statistical methods for handling attrition and missing data in clinical research. *Measurement and Evaluation in Counseling and Development*, **32**, 111–118.
- [16] Peng, C.Y.J., Harwell, M., Liou, S.M. and Ehman, L.H. (2006). Advances in missing data methods and implications for educational research. In S. Sawilowsky (Ed.), *Real Data Analysis* (PP. 31–78). Greenwich, CT:Information Age Publishing.
- [17] Raghunathan, T.E. (2004). What do we do with missing data? Some options for analysis of incomplete data. *Annual Review of Public Health*, **25**, 99–117.
- [18] Rubin, D.B. (1976). Inference and missing data. *Biometrika*, **63**, 581–592.
- [19] Rubin, D.B. (1987). *Multiple Imputation for Nonresponses in Surveys*. Wiley, New York.
- [20] SAS Institute, Inc. (2004). SAS/STAT[®] 9.1 User's Guide [Computer manual]. Cary, NC: Author.
- [21] Schafer, J.L. (1997). *Analysis of Incomplete Multivariate Data*. Chapman and Hall, London.
- [22] Schafer, J.L. (1999). Multiple Imputation : A primer. *Statistical Methods in Medical Research*, **8**, 3–15.
- [23] Schafer, J.L. and Graham, J.W. (2002). Missing data: Our view of the state of the art. *Psychological Methods*, **7**, 147–177.
- [24] Schafer, J.L. and Olsen, M.K. (1998). Multiple imputation for multivariate missing-data problems: A data analyst's perspective. *Multivariate Behavioral Research*, **33**, 545–571.

شیرین گلچی

فوق لیسانس آمار

تهران، خیابان فاطمی، خیابان باباطاهر، خیابان سرتیپ فکوری، شماره ۱۴۵، پژوهشکده‌ی آمار.

رایانشانی: golchishirin@yahoo.com

ساناز مهندسی

لیسانس آمار

تهران، خیابان فاطمی، خیابان باباطاهر، خیابان سرتیپ فکوری، شماره ۱۴۵، پژوهشکده‌ی آمار.

رایانشانی: s.mohandesi@src.ac.ir