

روش‌های وزن دهی احتمال معکوس و جانهی چندگانه برای تحلیل پاسخ در حالت گمشدگی

فرشته عثمانی، علی‌اکبر راسخی

گروه آمار زیستی، دانشکده علوم پزشکی، دانشگاه تربیت مدرس
تاریخ دریافت: ۱۳۹۵/۰۷/۰۲ تاریخ آخرین بازننگری: ۱۳۹۷/۰۴/۲۹

چکیده: ریزش داده‌ها و مقادیر گمشده از مشکلات معمول در تحلیل داده‌ها محسوب می‌شود. لذا اهمیت دارد که با برآورد مقادیر گمشده، داده‌ها کامل شده و در مسیری مناسب و صحیح برای تحلیل قرار داده شوند. دو روش معمول برای مقابله با داده‌های گمشده «جانهی چندگانه» و «وزن‌دهی احتمال معکوس» هستند. در این مقاله، رویکرد سومی معرفی خواهد شد که ترکیبی از دو روش جانهی چندگانه و وزن‌دهی احتمال معکوس است. با توجه به نتایج حاصل از مطالعه شبیه‌سازی روش ترکیبی مزایای بیشتری نسبت به سایر گزینه‌ها دارد. با توجه به وجود مقادیر گمشده در اکثر مطالعات به‌خصوص در حوزه پزشکی، نادیده گرفتن آن‌ها سبب تحلیل اشتباه شده و استفاده از روش‌های نیرومند، برای تحلیل صحیح، با اهمیت تلقی می‌شود.

واژه‌های کلیدی: جانهی چندگانه، وزن‌دهی احتمال معکوس، گمشدگی.

۱ مقدمه

داده‌های گمشده در بسیاری از مطالعات آماری از جمله مدل‌های رگرسیونی وجود دارند و باعث کاهش دقت برآورد می‌شوند. تاکنون روش‌های گوناگونی برای مقابله با مشکل داده‌های گمشده ابداع شده که عموماً بر داده‌های گمشده متغیر پاسخ متمرکز بوده است. حال آنکه متغیرهای پیش‌گو نیز می‌توانند دستخوش تغییر و از دست رفتن اطلاعات شوند. انتساب مقادیر مناسب به داده‌های گمشده یکی از چالش‌های موجود در پیش‌پردازش داده‌ها در بسیاری از حوزه‌ها است (یانگ و همکاران، ۲۰۰۸). طی سال‌های اخیر روش‌های

بسیاری برای غلبه بر این مشکلات ارائه شده است. حال آنکه استفاده از این روش‌ها در بسیاری از موارد به دلیل ارببی که وارد مسأله می‌کنند، بیش از آنکه کیفیت داده را افزایش دهد، باعث کاهش کیفیت داده می‌شود.

به‌طور کلی می‌توان چهار راهکار برای برخورد با داده‌های گمشده در نظر گرفت. اولین و ساده‌ترین راهکار، حذف واحدهای دارای مشاهدات ناقص از مطالعه و انجام تحلیل آماری بر اساس اطلاعات واحدهایی است که در تمام زمان‌های اندازه‌گیری در دسترس بوده‌اند. هرچند استفاده از روش تحلیل مبتنی بر موارد کامل (CC) به دلیل سادگی در بین پژوهشگران بسیار رایج است، اما در اکثر مواقع منجر به نتایج اربب و نامعتبر می‌شود. راهکارهای دیگر، روش‌های جانهی، روش‌های وزن‌دهی و روش‌های مبنی بر تابع درستنمایی هستند.

در حالت کلی هر سه روش یاد شده به جانهی مقادیر گمشده با مقادیری خاص می‌پردازند؛ با این تفاوت که جانهی مقادیر گمشده در روش اول به طور واضح و مستقیم، ولی در دو روش دیگر به‌طور غیرمستقیم انجام می‌شود. منظور از جانهی، جایگذاری مقادیری معقول به جای مقادیر گمشده است. بسته به برقراری فرضیات متفاوت در مورد الگوی داده‌های گمشده، روش مورد استفاده نیز متفاوت خواهد بود (شافر و همکاران، ۲۰۰۳). فعالیت‌های آماری انجام گرفته در معرفی و بررسی مزایا و معایب هر یک از راهکارهای فوق بسیار گسترده است. هم‌چنین کاملاً مشخص است که تحلیل‌های ساده از داده‌های موجود، سبب ایجاد ارببی‌هایی در برآورد پارامترها می‌شود. اما تاثیر دقیق این داده‌های ناکامل بستگی به تعداد مقادیر گمشده و هم‌چنین شدت ارتباط بین متغیرهای پاسخ و کمکی و نشانگرهای داده‌های گمشده دارد (هدکر و گیبنز، ۲۰۰۶).

عوامل مختلفی می‌توانند به دسترسی پاسخ و متغیرهای کمکی در زمان‌های بررسی تاثیرگذار باشند. معادلات برآورد تعمیم یافته وزنی احتمال معکوس (IPW-GEE) یک روش مناسب برای انجام این کار است. مطالعات تجربی نشان دادند که در نمونه‌های متوسط، ارببی برآوردهای تجربی به‌دست آمده خیلی کوچک هستند (شافر و همکاران، ۲۰۰۳). روش معادلات برآورد تعمیم یافته (GEE) برای اولین بار توسط لیانگ و زیگر (۱۹۸۶) به عنوان روشی برای تحلیل داده‌های طولی ارائه شد. تحت مکانیسم گمشدگی کاملاً تصادفی (MCAR) تحلیل‌ها براساس GEE برآوردهای سازگاری از پارامترهای رگرسیونی به‌دست می‌دهند. ولی اگر حالت‌های گمشدگی داده‌ها به‌صورت مکانیسم گمشدگی تصادفی (MAR) یا گمشدگی غیرتصادفی (MNAR) باشد، تحلیل‌ها براساس GEE برآوردهای ناسازگاری از پارامترهای رگرسیونی به‌دست می‌دهند (فیتسموریس و همکاران، ۲۰۰۹).

رایبیز و همکاران (۱۹۹۴) گروهی از معادلات برآوردی تعمیم یافته وزنی با احتمال معکوس را معرفی کردند که در آن وقتی حالت گمشدگی داده‌ها از نوع مکانیسم تصادفی باشد برآوردهای سازگاری به دست می‌دهد. در این روش وزن‌ها از مدل‌های گمشدگی به دست آمده‌اند که این مدل‌ها برای دستیابی به برآوردهای سازگار باید به درستی تعیین شوند. از طرفی دو روش رایج در برخورد با داده‌های گمشده، روش‌های MI و IPW هستند. روش IPW برای تعدیل کسرهای نامساوی از نمونه‌گیری قابل استفاده بوده است. به طور کلی روش MI با وجود پیچیدگی کاراتر از IPW است. به طوری که روش IPW فقط نیاز به یک مدل احتمال دارد، در حالی که روش MI نیاز به یک مدل توزیع توام از داده‌های گمشده به شرط داده‌های مشاهده شده نیز دارد. مناسب نبودن هر یک از این مدل‌ها در صورتی که حجم گمشدگی زیاد باشد، ممکن است منجر به اریبی قابل توجهی شود. در روش IPW مجدد تنها موارد کامل وارد تحلیل می‌شوند و وزن‌ها برای متعادل ساختن مجموعه موارد کامل به منظور نمایندگی از کل نمونه، استفاده می‌شوند.

با توجه به مطالب بیان شده، تمرکز مقاله حاضر بر مرور و مقایسه روش جانهای چندگانه و روش وزن‌دهی احتمال معکوس با استفاده از یک مطالعه شبیه‌سازی‌ریال خواهد بود.

۲ روش IPW-GEE

رایبیز و همکاران (۱۹۹۴) تعدیل اریبی ناشی از بی‌پاسخی واحدها را با روش وزن‌دهی معادله‌های برآورد تعمیم یافته مولنبرگ انجام دادند که در آن وزن برای هر واحد، به صورت یک ماتریس قطری در نظر گرفته شد. همچنین فیتسموریس و همکاران (۲۰۰۹) نیز نوع دیگری از وزن‌دهی را برای معادله‌های برآورد تعمیم یافته لیانگ و زیگر (۱۹۸۶) ارائه کردند که در آن وزن برای هر واحد نمونه‌ای تنها به صورت یک عدد در نظر گرفته می‌شد. وقتی گمشدگی MAR باشد، استنباط راجع به میانگین به هر نوع اشتباهی در تعیین توزیع توام بردار پاسخ، بسیار حساس است (دانیلز و هوگان، ۲۰۰۸). در روش GEE استاندارد لازم است، مدلی برای میانگین مشاهدات به شرط داشتن متغیرهای پیش‌گو وجود داشته باشد. به طور کلی، این مدل برای داده‌های مشاهده شده از نوع گمشدگی MAR برقرار نیست، بنابراین اعتبار تحلیل‌ها به خطر می‌افتد. با استفاده از روش GEE وزنی ساده، روش‌هایی برای تعدیل تحلیل‌ها ارائه می‌شود، که در آن وزن‌ها با استفاده از یک مدل برای احتمال گمشدگی برآورد می‌شوند و معمولاً روش‌های وزنی تمایل نامیده می‌شوند. در این روش وزن W_i را فقط برای افرادی که مطالعه را بدون انصراف به پایان رسانده‌اند، با استفاده از مقادیر پاسخ‌های مشاهده شده و متغیرهای کمکی محاسبه می‌کنند. این وزن برای هر فرد

به صورت

$$\begin{aligned} W_i &= P(D_i = n + 1)^{-1} \\ &= [P(D_i > 1 | D_i \geq 1) P(D_i > 2 | D_i \geq 2) \dots P(D_i > n | D_i \geq n)]^{-1} \\ &= [\pi_{i1} \pi_{i2} \dots \pi_{in}]^{-1} \end{aligned}$$

محاسبه می‌شود، که در آن

$$\pi_{ij} = P(R_{ij} = 1 | R_{ij} = \dots = R_{i,j-1} = 1, X_i, Y_{i1}, \dots, Y_{i,j-1}).$$

سپس هر فردی که احتمال باقی ماندن کمتری داشته باشد (احتمال انصراف بیشتری داشته باشد)، وزن بیشتری را به خود اختصاص می‌دهد. در نهایت روش GEE استاندارد را با این اوزان تطبیق می‌دهند تا براساس یک روش GEE موزون داده‌ها تحلیل شود. پس در این مدل، به افرادی که انصراف نداشته‌اند نیز وزنی برای انصراف داده می‌شود تا براساس این وزن، توزیع شرطی متغیر پاسخ در افرادی که انصراف نداشته‌اند همانند افرادی شود که انصراف داشته‌اند (فیتسموریس و همکاران، ۲۰۰۴). در روش GEE استاندارد، از معادلات برآوردگر

$$\sum_{i=1}^N D_i' V_i^{-1} (Y_i - \mu_i) = 0$$

استفاده می‌شود، اما در GEE موزون، از معادلات برآوردگر موزون

$$\sum_{i=1}^N D_i' V_i^{-1} W_i (Y_i - \mu_i) = 0$$

استفاده می‌شود.

برای برآورد میانگین پاسخ (در حالت گمشدگی در پاسخ) روش IPW نیز قابل استفاده است. این تحلیل وقتی احتمال مشاهده شدن پاسخ فرد به پاسخ مربوطه بستگی نداشته باشد، برآورد سازگاری به دست می‌دهد. اما در غیر این صورت اریب خواهد بود (دانیلز و هوگان، ۲۰۰۸). برآوردگر IPW هورویترز-

تامسون روشی برای تصحیح این آریبی به دست می‌دهد، در این روش، مجدد تنها افرادی با مشاهدات پاسخ کامل وارد می‌شوند، اما وزن‌هایی برای متعادل کردن مجموعه موارد کامل استفاده می‌شوند. تا اینکه به عنوان نماینده‌ای از کل نمونه به حساب آیند. وزن هر فرد برابر با معکوس احتمال کامل کردن مطالعه است که این احتمال ناشناخته و نیاز به برآورد دارد (سیمن و همکاران، ۲۰۱۳).

۳ روش جانهای چندگانه

استفاده از روش‌های جانهای، برای پر کردن جاهای خالی در داده‌هاست؛ به این صورت که جاهای خالی حاصل از گمشدگی با مقادیر مناسب جانهای شده پر می‌شوند، یکی از مزایای این رویکرد، این است که با جایگذاری مقادیر گمشده، مجموعه داده‌ها به صورت کامل درآمده و برای تحلیل می‌توان از روش‌های مناسب مربوط به داده‌های کامل استفاده کرد (عثمانی و همکاران، ۲۰۱۴). در روش جانهای ساده، به جای مقادیر گمشده، یک مقدار جانهای شده تک قرار داده می‌شود و در تحلیل با مقدار جانهای شده به همان طریقی رفتار می‌کند که با مقادیری واقعی اندازه‌گیری شده رفتار می‌شود و عدم حتمیت این مقادیر گمشده در نظر گرفته نمی‌شود. به همین دلیل، رابینز و وانگ (۲۰۰۰) روشی به نام جانهای چندگانه را ابداع کردند، که در آن جای هر مقدار گمشده با دو یا چند مقدار جانهای شده پر شده و هر مجموعه از داده‌های به دست آمده با استفاده از روش‌های مربوط به داده‌های کامل تحلیل می‌شود. با ترکیب نتایج این تحلیل‌ها، استنباط‌هایی به دست می‌آیند که می‌توانند عدم حتمیت مقادیر گمشده را در نظر بگیرند (روبین، ۱۹۸۷). در این روش به تعداد جانهای صورت گرفته، مجموعه داده کامل تولید و بر مبنای هر یک از این مجموعه داده‌ها، برآوردی از مولفه‌های مدل و واریانس آن‌ها حاصل می‌شود. سپس در انتها این مقادیر با یکدیگر ترکیب می‌شوند تا برآوردی کلی حاصل شود. از روش‌های رایج جانهای چندگانه می‌توان به روش نمره تمایل، روش مبتنی بر مدل پیش بینی کننده و روش تطابق میانگین پیش‌بینی کننده اشاره کرد. روش جانهای چندگانه دارای مزایا و معایبی است. مزیت اصلی روش جانهای چندگانه همان‌طور که توسط رابینز و وانگ (۲۰۰۰) اشاره شده، این است که جانهای چندگانه به گردآورندگان داده امکان انعکاس عدم قطعیت مقادیر جانهای شده را می‌دهد. از جمله معایب این روش زمان بر بودن جانهای چندین مجموعه داده کامل، آزمون جداگانه هر کدام از این مجموعه داده‌ها و همچنین ترکیب نتایج این مدل‌ها است. در روش MI، داده‌های گمشده با داده‌های به دست آمده از مدل جانهای جایگزین می‌شوند. این عمل M بار تکرار می‌شود، در نتیجه M مجموعه داده کامل بوجود می‌آید. هر کدام از این مجموعه‌ها، جداگانه تحلیل می‌شوند و برآوردی از پارامترهای مدل θ به دست می‌دهند. فرض کنیم $\hat{\theta}$ برآوردگر θ از مجموعه داده کامل، \hat{V} واریانس برآورد شده و $\hat{V}_{(m)}$ و $\hat{\theta}_{(m)}$

مقادیر برآورد شده از m امین مجموعه داده جانهای شده باشند ($m = 1, \dots, M$). برآورد و واریانس آن به ترتیب عبارتند از

$$\hat{\theta}_M = \frac{1}{M} \sum_{i=1}^M \hat{\theta}_M,$$

$$\hat{V}_M = \frac{1}{M} \sum_{i=1}^M \hat{V}_{(M)} + (1 + M^{-1})(M - 1)^{-1} \sum_{i=1}^M (\hat{\theta}_{(m)} - \hat{\theta}_M)(\hat{\theta}_{(m)} - \hat{\theta}_M)^T.$$

برای به دست آوردن برآورد پارامترها در این روش از قوانین ادغام رویین (۱۹۸۷) استفاده شده و برآورد واحدی برای پارامتر رگرسیونی با معدل گرفتن از M برآورد حاصله به دست می‌آید. واریانس برآورد نیز با ترکیب واریانس درون جانهای و واریانس بین جانهای به دست می‌آید. یعنی

$$\hat{\beta} = \frac{1}{M} \sum_{k=1}^m \hat{\beta}^{(k)} = \bar{\beta}$$

و

$$\frac{1}{M} \sum_{k=1}^m \text{Cov}(\hat{\beta}^k) = (1 + \frac{1}{m}) \frac{1}{m-1} \sum_{k=1}^m (\hat{\beta}^k - \bar{\beta})(\hat{\beta}^k - \bar{\beta})^T.$$

روش‌های IPW و MI در حالت گمشدگی تصادفی و در صورتی که مدل وزندهی و جانهای به درستی مشخص شده باشد، برآوردگرهای سازگاری از θ به دست می‌دهند. در روش MI وقتی که $\hat{\theta}$ برآوردگر MLE، ماتریس \hat{V} معکوس ماتریس اطلاع فیشر و همچنین داده‌های گمشده از توزیع پسین بیزی نمونه‌گیری شده باشند، استنباط می‌شود که $\hat{\theta}_M$ به طور مجانبی دارای توزیع نرمال با واریانس V و \hat{V}_M برآوردگر نارایب V_M است و در صورتی که $M = \infty$ ، این برآوردگر سازگار نیز خواهد بود (فیتسموریس و همکاران، ۲۰۰۴).

اگر مدل جانهای به درستی مشخص شده باشد، روش MI اغلب به روش IPW به علت کارایی بیشتر ترجیح داده می‌شود. ولی در صورتی که مقادیر گمشده جانهای شده زیاد باشند، ممکن است منجر به آریبی قابل توجهی شود؛ اما اگر متغیرهای کمی دارای مقادیر گمشده باشند، جانهای این مقادیر نسبت به خارج کردن آن‌ها رضایت‌بخش‌تر خواهد بود. از طرف دیگر، اگر متغیرهای زیادی از افراد دارای مقادیر گمشده باشند؛ پیچیدگی مدل و احتمال انتخاب نادرست مدل جانهای ممکن است نگران‌کننده باشد. این موقعیت ممکن است به عنوان مثال در مطالعات طولی، زمانی که کلیه بلوک‌های داده‌ها در بعضی از افراد به دلیل از

دست دادن ملاقات‌ها گمشده باشند یا در مطالعه‌ای که بعضی از افراد همه سوال‌های مربوطه را پاسخ نمی‌دهند، پیش‌آید که در چنین موقعیت‌هایی تحلیل‌گر اعتماد بیشتری به استفاده از روش IPW دارد (شافر، ۱۹۹۷).

۴ مطالعه شبیه‌سازی

در این بخش، عملکرد تجربی روش‌های ارائه شده از طریق یک مطالعه شبیه‌سازی بررسی و با سایر روش‌ها در یک مدل رگرسیونی استفاده خواهد شد. همچنین از روش IPW/MI برای جانهی پاسخ در یک مدل رگرسیونی استفاده خواهد شد.

۱.۴ مدل شبیه‌سازی

مدل تحلیل، تنها برای افرادی با متغیرهای کمکی کامل، برازش داده شده و متغیرپاسخ در این افراد جانهی خواهد شد. برای تحلیل نمونه تولید شده دو مرحله گمشدگی در X و Y وجود دارد. در مرحله اول یا می‌توان افراد با X ناکامل را حذف کرد یا اینکه مقادیر گمشده X را جانهی نمود. به طور مشابه، هر فردی با پاسخ گمشده که در مرحله اول حذف نشده، می‌تواند در مرحله دوم حذف یا جانهی شود. قابل ذکر است که در هر مرحله اگر حذف انجام شود می‌توان از روش IPW برای تعدیل داده‌ها استفاده کرد. بنابراین در هر مرحله سه حالت و در کل $3 \times 3 = 9$ حالت وجود دارد. اگر راهکارهای بکار گرفته شده در این دو مرحله با نماد ST1/ST2 نشان داده شود، این راهکارها می‌توانند یکی از روش‌های مبتنی بر موارد کامل (حذف بدون وزن دهی)، IPW (حذف با وزن دهی بقیه) یا MI باشند.

در روش IPW/MI که هدف اصلی این مقاله است افراد با متغیرهای گمشده X حذف می‌شوند و تعدیل داده‌ها با وزن دهی سایر داده‌ها صورت می‌گیرد. همچنین در افراد با X کامل اما با گمشدگی در متغیر پاسخ، Y جانهی می‌شود و در روش CC/CC هم تنها افرادی که X و Y کامل دارند بدون وزن دهی تحلیل می‌شوند و در روش IPW/IPW هم همین افراد استفاده می‌شوند، با این تفاوت که در این حالت، افراد به اندازه معکوس احتمال باقی ماندن در مطالعه وزن دهی می‌شوند، همچنین در روش MI/MI همه مقادیر گمشده جانهی می‌شوند. به علاوه روش‌های CC/MI، CC/IPW و IPW/CC هم در نظر گرفته شدند.

۲.۴ تولید داده‌ها

انجام این شبیه‌سازی با سه هدف دنبال شد: الف- بررسی ناریب بودن \hat{V} برای روش IPW/MI. ب- بررسی کارایی روش IPW/MI نسبت به روش IPW/IPW. ج- نشان دادن اینکه روش MI/MI زمانی که مدل جانهای به‌درستی تعیین نشود، برآوردگرهای اربیی را به‌دست می‌دهد و IPW/MI به‌طور مجانبی ناریب و یا حداقل از روش MI/MI دارای اربیی کمتری است. با توجه به این اهداف فرایند تولید داده‌ها به صورت زیر در نظر گرفته شد.

بردار $X = (x_1, \dots, x_5)$ و متغیر Y برای $M = 1000$ فرد تولید شد. برای هر فرد، x_1 با احتمال ۰/۵ مقدار یک و در غیر این صورت صفر را اختیار می‌کند. مشاهدات (x_2, x_3, x_4) به‌صورت مستقل از توزیع $N(0, 1)$ تولید شدند. همچنین x_5 از توزیع $N(x_2x_3, 1)$ تولید شد. متغیر پاسخ نیز از رابطه

$$Y = -3 + x_1x_2 + x_1x_3 + 0/5x_2x_3 + x_4 + 0/5x_5 + \epsilon, \quad \epsilon \sim N(0, 1) \quad (1)$$

محاسبه شد. از بین متغیرها x_1 برای همه N نفر مشاهده شده بود ولی با $X = (x_2, x_3, x_4, x_5)$ احتمال $0/8 - 0/6x_1$ مشاهده شدند. از طرفی اگر بردار $X = (x_2, x_3, x_4, x_5)$ مشاهده می‌شد متغیر پاسخ Y نیز با احتمال $1 + \exp(-1/5 + 0/6x_2x_3)$ مشاهده و در غیر این صورت گمشده در نظر گرفته شد. مدل تحلیل نیز

$$Y = \theta_0 + \theta_2x_2 + \theta_3x_3 + \theta_{23}x_2x_3 + e$$

بود، که در آن $E(e|x_2, x_3) = 0$. بدین ترتیب بردار Z به صورت $Z = (1, x_2, x_3, x_2x_3)$ در نظر گرفته شد که با انتگرال‌گیری از رابطه (۹۹) نسبت به x_1, x_4, x_5 مدل تحلیل به‌درستی مشخص خواهد شد. مقادیر θ هم به‌صورت $(\theta_0, \theta_2, \theta_3, \theta_{23}) = (-3, 0/5, 0/5, 1)$ در نظر گرفته شد. این فرایند تولید داده به سه دلیل انتخاب شده است:

۱- اثرات متقابل x_1x_2, x_1x_3 در رابطه (۹۹) نشان دهنده ارتباط بین Y و x_2x_3 در دو طبقه x_1 است. هم‌چنین احتمال مشاهده در یک طبقه x_1 برابر ۰/۲ و در طبقه دیگر برابر ۰/۸ است. بنابراین ارتباط بین Y و x_2x_3 در افراد با X های کامل و افراد با X های ناکامل با یکدیگر فرق می‌کند. از طرفی تعدیل نکردن گمشدگی در مرحله اول باعث اربیی θ_2 و θ_3 می‌شود. بنابراین روش‌های CC/MI، CC/IPW و CC/CC اربیب خواهند بود.

۲- برای افرادی با بردار مشاهده شده (x_2, x_3, x_4, x_5) ، احتمال مشاهده Y به x_4 بستگی دارد که در مدل تحلیل نیست ولی با Y در ارتباط است؛ این مسئله باعث می‌شود که ارتباط بین X و Y در مدل تحلیل، در مجموعه داده کامل با مجموعه داده با X کامل ولی با گمشدگی در Y متفاوت باشد. بخصوص به دلیل اینکه احتمال گمشدگی در Y به x_2, x_4 بستگی دارد. ارتباط بین Y و x_2 در دو مجموعه متفاوت خواهد بود. از طرفی تعدیل نشدن گمشدگی در مرحله ۲ منجر به اریبی بخصوص در θ_2 خواهد شد.

۳- مقدار x_5 که در فرایند تولید داده برای Y بکار رفته بود، نشان می‌دهد که استفاده از روش MI در مرحله یک در صورتی که مدل جانهی X به درستی تعیین نشده باشد می‌تواند باعث اریبی شود (نتایج MI/MI در جدول ۱). یک مجموعه ۱۰۰۰ تایی داده تولید شده و این هفت روش برای هر یک از مجموعه داده‌ها بکار گرفته شد. برای هر یک از $(\theta_0, \theta_2, \theta_3, \theta_{23})$ و هر یک از روش‌ها، میانگین و واریانس ۱۰۰۰ پارامترهای برآورد شده محاسبه شدند. مقدار SE نیز به صورت انحراف معیار برآورد پارامترها محاسبه شد. هم‌چنین در روش‌هایی که شامل روش جانهی بودند ده بار جانهی انجام شد.

۳.۴ تحلیل داده‌ها با روش ترکیبی جانهی چندگانه

در روش MI/MI مدل جانهی در مرحله ۱ به صورت

$$(X_2, X_3, X_4) \sim N((\gamma_2, \gamma_3, \gamma_4), \Sigma_1)$$

$$X_5 | X_2, X_3 \sim N(\gamma_5 + \gamma_6 X_2 + \gamma_7 X_3 + \gamma_8 X_2 X_3, \Sigma_2)$$

در نظر گرفته شد. از پیشین‌های ویشارت معکوس و نرمال استفاده شد که باعث می‌شود پسین‌ها هم نرمال و ویشارت معکوس به دست بیاید. برای روش‌های CC/MI و IPW/MI مدل جانهی استفاده شده در مرحله ۲ بصورت زیر است.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5$$

$$+ \beta_{12} X_1 X_2 + \beta_{13} X_1 X_3 + \beta_{23} X_2 X_3 + \beta_{123} X_1 X_2 X_3 + \epsilon$$

۴.۴ تحلیل داده‌ها با روش ترکیبی وزن‌دهی احتمال معکوس

برای روش‌های IPW/MI، IPW/CC و IPW/IPW وزن‌هایی با برازش مدل گمشدگی برای مرحله ۱ برآورد شدند، توجه شود که چون X_1 دو حالتی است تابع

$$W = (\delta_0 + \delta_1 X_1)^{-1} = \delta_0^{-1} - \delta_1 X_1 [\delta_0 (\delta_0 + \delta_1)]^{-1}$$

یک تابع خطی از X_1 است. از این رو مدل جانهی مرحله ۲ شامل $X_1 Z$ و Z است که شامل WZ است. برای روش‌های CC/IPW و IPW/IPW وزن‌هایی با استفاده از مدل صحیح مشخص شده در مرحله ۲ به صورت

$$\text{logit } P(Y \text{ observed} | X \text{ observed}) = \delta_2 + \delta_3 X_2 + \delta_4 X_4 + \delta_5 X_2 X_5$$

برآورد می‌شود. در روش IPW/IPW احتمال کامل بودن یک فرد، برابر ضرب این دو احتمال IPW است. برآورد میانگین پارامترها و مقادیر انحراف معیار و هم‌چنین جذر میانگین واریانس‌های برآورد شده در جدول ۱ نشان داده شده است. با توجه به نتایج به دست آمده می‌توان گفت بطور مجانبی برآوردگرهای ناراییی از پارامترها و انحراف معیارها به دست آمده است.

۵.۴ مقایسه روش‌های ترکیبی

همان‌طور که در بخش پیش بیان شد، روش‌های CC/IPW، CC/MI و IPW/CC برای حداقل یکی از پارامترها اریب هستند. ولی می‌توان گفت روش‌های IPW/IPW و MI/MI تقریباً ناریب هستند. مدل MI/MI کارایی کمتری نسبت به روش IPW/MI دارد زیرا مدل جانهی در راهکار دوم از اطلاعات کمکی نظیر متغیرهای کمکی به‌ویژه (X_4, X_5) که در مدل تحلیل وارد نشدند، استفاده می‌کند. کاراترین روش ناریب با توجه به نتایج به دست آمده روش MI/MI است. جانهی، زمانی که مدل جانهی به درستی مشخص شده باشد بهترین روش است ولی زمانی که مدل جانهی در مرحله ۱ یا ۲ به درستی مشخص نشود، روش MI/MI ممکن است اریب باشد. همان‌طور که نیمی از افراد دارای X های ناکامل بودند با برازش مدل تحلیل به کل نمونه برآورد θ_{23} در حدود ۰/۷۵ به دست آمد. این موضوع در سطر MI*/MI در جدول ۱ نشان داده شده است (نماد MI* مدل جانهی نادرست مشخص شده را نشان می‌دهد). از طرفی با فرض

اینکه مدل جانہی در مرحلہ ۱ بہ درستی مشخص شدہ باشد ولی در مرحلہ ۲ بہ دلیل حذف $\beta_{23}X_2X_3$ ، $\beta_{123}X_1X_2X_3$ و β_5X_5 بہ درستی مشخص نشدہ باشد، مقادیر گمشدہ Y با این فرض کہ هیچ اثر متقابلی بین X_2 و X_3 نیست، جانہی خواهد شد. حدود ۶۰ درصد از مقادیر Y گمشدہ هستند و بہ همین دلیل مقادیر θ_{23} کم برآورد خواهند شد. این نتایج در ردیف MI/MI^* جدول ۱ نشان دادہ شدہ است. ردیف IPW/MI^* نتایج مربوط بہ روش IPW/MI با مدل جانہی بہ اشتباہ مشخص شدہ در مرحلہ ۲ رانشان می دہد. این روش بہ طور قابل توجہی نسبت بہ روش MI/MI^* اریبی کمتری خواهد داشت، زیرا مقادیر Y کمتری جانہی شدہ اند. بنابراین قسمت IPW در روش IPW/MI یک روش محافظہ کار در برابر انتخاب نادرست مدل جانہی می باشد.

بسیاری از نرم افزارہای آماری، دارای بستہ ہا و روش ہایی برای تحلیل مقادیر گمشدہ هستند کہ در این مقالہ از نرم افزار R نسخہ ۳/۱/۰ و نرم افزار SAS استفادہ شد.

بحث و نتیجہ گیری

مشکل وجود دادہ ہای گمشدہ بخصوص در مطالعات پزشکی، در سالہای اخیر بہ شدت مورد توجہ قرار گرفتہ است (عثمانی و ہمکاران، ۲۰۱۵). این مقالہ با مقایسہ رویکردہای مختلف بہ دنبال یافتن بہترین رویکرد در مواجہہ با گمشدگی است. در تحقیق صورت گرفتہ، ہدف مقایسہ ترکیبی روش ہای معادلہ برآورد وزنی معکوس و روش جانہی چندگانہ بود. بہ بیان دیگر، در این مقالہ از ترکیب رویکردہای دادہ ہای کامل، جانہی چندگانہ و وزندہی احتمال معکوس با ہدف مقایسہ این روش ہادر تحلیل دادہ ہا استفادہ شد. جدول ۱ برازش مدل رگرسیونی مفروض را بر دادہ ہای تولید شدہ با استفادہ از رویکردہای یاد شدہ نشان می دہد. همان طور کہ مشاہدہ می شود ترکیب رویکردہای مختلف، نتیجی متفاوت را از لحاظ اریبی برآورد پارامترہا نشان می دہند.

بطور کلی، رویکرد استفادہ از دادہ ہای کامل، روشی مناسب برای حل مشکل گمشدگی در دادہ ہا نیست، زیرا باعث کاهش شدید حجم نمونہ و در نتیجہ، کارایی برآوردہا می شود (مولنبرگ و ہمکاران، ۲۰۰۸). رابینز و وانگ (۲۰۰۰) فرمولی کلی برای واریانس برآوردگر MI براساس برآوردگر دادہ ہای کامل با حل مجموعہ ای از معادلات برآوردگر ارائه کردند. این فرمول زمانی کہ جانہی نامناسب و مدل جانہی پارامتری باشد، استفادہ می شود. در روش IPW/MI نیز این رابطہ می تواند استفادہ شود.

جدول ۱: میانگین برآورد پارامتر، جذر میانگین واریانس برآورد شده (aSE) و برآورد تجربی (eSE) برای چهار پارامتر و ده روش تحلیل

eSE	$\theta_3 = 1$		$\theta_4 = -0.5$		$\theta_5 = -0.5$		$\theta_6 = -2$		روش
	aSE	میانگین	eSE	میانگین	eSE	میانگین	eSE	میانگین	
۰/۰۹۱	۰/۰۸۲	۱/۰۰۵	۰/۰۸۷	۰/۰۸۰	۰/۰۸۷	۰/۰۷۸	۰/۰۷۹	۰/۰۸	CC/CC
۰/۱۰۷	۰/۰۹۶	۱/۰۰۴	۰/۰۸۹	۰/۰۸۶	۰/۰۹۱	۰/۰۸۹	۰/۰۷۹	۰/۰۸۲	CC/IPW
۰/۰۸۶	۰/۰۸۴	۱/۰۰۴	۰/۰۸۳	۰/۰۷۹	۰/۰۸۳	۰/۰۸۱	۰/۰۷۵	۰/۰۷۵	CC/MI
۰/۱۱۹	۰/۱۱۳	۱/۰۰۷	۰/۱۱۴	۰/۱۰۹	۰/۱۱۱	۰/۱۱۰	۰/۱۰۴	۰/۱۰۲	IPW/CC
۰/۱۳۲	۰/۱۲۱	۱/۰۰۶	۰/۱۱۶	۰/۱۱۲	۰/۱۲۲	۰/۱۱۸	۰/۱۰۸	۰/۱۰۶	IPW/IPW
۰/۱۱۲	۰/۱۰۹	۱/۰۰۸	۰/۱۰۷	۰/۱۰۴	۰/۱۰۳	۰/۱۰۳	۰/۰۹۴	۰/۰۹۷	IPW/MI
۰/۰۸۲	۰/۰۹۲	۱/۰۰۶	۰/۰۸۸	۰/۰۹۰	۰/۰۸۷	۰/۰۹۲	۰/۰۷۹	۰/۰۸۹	MI/MI
۰/۰۸۴	۰/۱۰۲	۰/۰۴۹	۰/۰۹۴	۰/۰۹۴	۰/۰۹۶	۰/۰۹۵	۰/۰۸۳	۰/۰۹۲	MI*/MI
۰/۰۵۵	۰/۰۹۱	۰/۰۳۹	۰/۰۵۱	۰/۰۸۸	۰/۰۵۲	۰/۰۸۸	۰/۱۰۲	۰/۱۰۷	MI/MI*
۰/۱۲۷	۰/۱۳۲	۰/۰۷۶	۰/۱۱۵	۰/۱۱۷	۰/۰۹۵	۰/۱۱۹	۰/۱۰۱	۰/۱۰۶	IPW/MI*

در مورد رگرسیون خطی با اعمال روش MI برای پاسخ گمشده، برآوردگر واریانس روبین برای روش IPW/MI وقتی M به سمت بینهایت میل کند، سازگار است. شافر و همکاران (۲۰۰۳) نیز اشاره به این نکته کرده است که اگرچه اثبات عملکرد مناسب روش MI مشکل است، ولی برخی از پژوهشگران ممکن است استفاده از روش MI/MI را به شرطی که مدل‌های جانهای به‌درستی مشخص شده باشند، ترجیح بدهند. در نتیجه این روش کارآمدتر از روش IPW/MI خواهد بود (شافر و همکاران، ۲۰۰۳). با این حال، اگر نتایج حاصل از IPW/MI و MI/MI بسیار متفاوت باشند، تحقیق بیشتر در راستای تصحیح مدل جانهای لازم خواهد بود، زیرا تشخیص نادرست مدل جانهای در MI/MI معمولاً باعث اریبی می‌شود.

علاوه بر این، روش IPW/MI وقتی مدل وزن‌دهی در مقایسه با مدل جانهای نسبتاً ساده‌تر باشد، مطلوب‌تر خواهد بود. روش جایگزین دیگری برای IPW/MI روش IPW/IPW است که با وجود ساده‌تر بودن، دارای این نقطه ضعف است که اگر یک فرد حتی فقط یک متغیر از او گمشده باشد از مطالعه و تحلیل حذف می‌شود (رابینز و گیل، ۱۹۹۷). از سوی دیگر IPW/MI امکان استفاده از یک مجموعه از وزن‌ها را می‌دهد.

یافته‌های این مطالعه نشان داد که ترکیب IPW/MI برای تحلیل داده‌های گمشده در صورتی که مدل وزن‌دهی به نسبت مدل جانهای ساده‌تر باشد، بهتر خواهد بود. از طرفی با توجه به اینکه مقادیر گمشده همواره مشاهده نشده باقی خواهند ماند، تمام روش‌های موجود برای تحلیل داده‌های گمشده شامل فرضیات غیرقابل شناسایی و غیرقابل اثبات هستند، بنابراین بهتر است در برخورد با این داده‌ها صرفاً به نتایج حاصل از یک روش اتکا نشود و در این راستا تحلیل حساسیت مناسب نیز انجام شود.

تقدیر و تشکر

نویسندگان از پیشنهادات ارزنده داوران گرامی و ویراستار مجله که باعث ارائه بهتر و بهبود مقاله شده است کمال تشکر و قدردانی را اعلام می‌دارند.

مراجع

Daniels, M. J. and Hogan, J. W. (2008), *Missing Data in Longitudinal Studies*, London, Chapman & Hall.

- Diggle, P., Heagerty, P., Liang K. Y. and Zeger. S. (2002), *Analysis of Longitudinal Data*, Second Edition, Oxford University Press.
- Fitzmaurice, G. M., Laird, N. M. and Ware, J. H. (2004), *Applied Longitudinal Analysis*, Wiley-Interscience.
- Fitzmaurice, G. M., Davidian M., Verbeke, G. and Molenberghs, G. (2009), *Longitudinal Data Analysis*, Chapman Hall/CRC.
- Hedeker, D. and Gibbons R. D. (2006), *Longitudinal Data Analysis*, New York, Wiley.
- Liang, K. and Zeger, S. L. (1986), Longitudinal Data Analysis Using Generalized Linear Models, *Biometrika*, **73**, 13-22.
- Lipsitz S. R., Fitzmaurice J. G., Ibrahim M., Sinha D. M. Parzen, M. and Lipshultz, S. (2009), Joint Generalized Estimating Equations for Multivariate Longitudinal Binary Outcomes with Missing Data: An Application to Acquired Immune Deficiency Syndrome Data. *Journal of the Royal Statistical Society, Series A*, **1**, 3-20.
- Molenberghs, G., Kenward, M. G. and Goetghebeur, E. (2001), Sensitivity Analysis for Incomplete Contingency Tables: The Slovenian Plebiscite Case. *Journal of Applied Statistics*, **50**, 15-29.
- Molenberghs, G., Beunckens, C., Sotito, C. and Kenward, M. G. (2008), Every Missingness Not at Random Model Has a Missingness at Random Counterpart with Equal Fit, *Journal of the Royal Statistical Society, Series B*, **70**, 371-88.
- Robins, J., Rotnitzky, A. and Zhao, L. (1994), Estimation of Regression Coefficients When Some Regressors Are Not Always Observed. *Journal of the American Statistical Association*, **89**, 846-866.
- Robins, J. M. and Gill, R. D. (1997), Non-response Models for the Analysis of Non-monotone Ignorable Missing Data. *Statistics in Medicine*, **16**, 39-56.
- Osmani, F., Hajizadeh, E. and Mansoori, P. (2014), Analysis of Risk Factors for Psoriasis Recurrence Using Proportional Rates Model, *Journal of Skin and Stem Cell*, **3**, 129-137.
- Osmani, F., Hajizadeh, E. and Mansoori, P. (2015), Estimation of Seasonal Effect on the Psoriasis Recurrence Using Time Dependent Coefficient Rates Model for Recurrent Events, *Dermatology and Cosmetic*, **6**, 23-30.
- Robins, J. and Wang, N. (2000), Inference for Imputation Estimators, *Biometrika*, **87**, 113-24.

- Rubin, D. B. (1987), *Multiple Imputation for Nonresponse in Surveys*, Wiley, New York.
- Seaman, S. R., White, I. R., Copas, A. J. and Li, L. (2012), Combining multiple imputation and inverse-probability weighting, *Biometrics*, **68**, 129-137.
- Seaman, S. R. and White, I. R. (2013), Review of Inverse Probability Weighting for Dealing with Missing Data, *Statistical Methods in Medical Research*, 278-295.
- Schafer J. L. (1997), *Analysis of Incomplete Multivariate Data*, Chapman and Hall, London.
- Schafer, J. L. (2003), Multiple Imputation in Multivariate Problems When the Imputation and Analysis Models Differ. *Statistica Neerlandica*, **57**, 19-35.
- Vansteelandt, S., Carpenter, J. and Kenward, M. G, (2010), Analysis of Incomplete Data Using Inverse Probability Weighting and Doubly Robust Estimators, *Methodology* **6**, 37-48.
- Yang X., Li J. and Shoptaw S. (2008), Imputation-based Strategies for Clinical Trial Longitudinal Data with Nonignorable Missing Values. *Statistics in Medicine*, **27**, 2826-2849.