

بررسی روشهای برخورد با داده‌های گمشده

افشین آشفته^۱

چکیده

در این مقاله، روشهای جایگذاری داده‌های گمشده مانند جانهای چندگانه، الگوریتم EM^۲ و داده افزایی^۳ مرور و توضیح داده می‌شوند. شاید تا چندی پیش روش حذف معمول‌ترین روشی بود که برای برخورد با داده‌ها گمشده کاربرد داشت، ولی در حال حاضر، با وجود این روشهای جدید و قدرتمند آماری می‌توان برآوردهای خوبی را برای آنها به دست آورد.

واژه‌های کلیدی: داده گمشده تصادفی، داده گمشده کاملاً تصادفی، زنجیر مارکوف، الگوریتم EM، داده افزایی، استنباط جانهای چندگانه.

۱. مقدمه

داده گمشده در بیشتر تحقیقات علوم رفتاری و اجتماعی یک شکل آشنا است. در بیشتر مواردی که تحقیق توسط پرسشنامه انجام می‌گیرد، پاسخ دهندگان ممکن است تمایلی به پاسخ دادن نداشته باشند و یا به علل مختلف مانند کمی وقت، بیماری و یا مشکلات دیگر سوالات را بی‌پاسخ بگذارند. در این صورت، در مجموعه داده‌ها، داده گم شده به وجود می‌آید که برای جلوگیری از مشکلاتی مانند کمی حجم نمونه و یا عدم حضور واحدهایی از نمونه در مجموعه داده‌ها باید با استفاده از روشهای مناسب جایگذاری شوند.

۲. چند روش در برخورد با داده گمشده

۱- حذف داده گمشده از تحلیل.

۲- جایگذاری داده گمشده با مقادیر مناسب مانند میانگین.

۳- استفاده از برآوردهای رگرسیونی.

۲-۱- روش حذف داده‌های گمشده

در تجزیه و تحلیلهای چند متغیره با تعداد زیادی متغیر، روش حذف بسیار ناکارآمد است، چراکه تعداد کمی از واحدهای نمونه دارای داده کامل برای تمامی متغیرها می‌باشند. پس، تنها به علت بی‌پاسخ بودن یک متغیر یک واحد نمونه از تحلیل خارج می‌شود. از طرف دیگر، موجب اریب شدن برآوردها می‌گردد، چراکه هیچ تضمینی وجود ندارد که واحدهای باقیمانده نماینده خوبی برای کل داده‌ها باشند.

۲-۲- جایگذاری داده گمشده

جایگذاری داده گمشده با میانگین می‌تواند رابطه بین متغیرها را به طور جدی تحت تاثیر قرار دهد، ولی استفاده از برآوردهای رگرسیونی این

^۱ گروه آمار، دانشگاه اصفهان

^۲ Expectation- Maximization Alogorithm

^۳ Multiple Imputation

^۴ The EM Algorithm

^۵ Data augmentation

^۶ Multiple imputations

رابطه را حفظ می‌کند. حتی اگر جایگذاری به طوری انجام گیرد که رابطه بین متغیرها و توزیع احتمال آنها به طور کامل حفظ گردد، باز هم این روش مقادیر دقیقی از داده‌های گمشده را فراهم نمی‌کند، چرا که در این روش مقادیر برآورد شده جایگزین مقادیری ناشناخته می‌شوند که باید حالت احتمالی در برآورد آنها در نظر گرفته شود و نمی‌توان بدون در نظر گرفتن غیر حتمی برآوردهای داده‌های گمشده تقریب خوبی از آنها را جایگذاری کرد.

۲-۳- روشهای جایگذاری چندگانه

- الگوریتم EM
- داده افزایی
- استنباط جهانی چندگانه

۲-۳-۱- الگوریتم EM

در اواخر سال‌های ۱۹۷۰ روین، دمپستر و لارد^۷ الگوریتم EM را معرفی و گسترش دادند که روشی محاسباتی برای برآورد داده‌های گمشده بود. ثابت شد که این روش یک روش محاسباتی بسیار کارا است و نکته قابل ذکر دیگر تاثیر قابل توجه این روش بر دیدگاه آماردانان بود. تا آن زمان داده‌های گمشده به عنوان عاملی مزاحم بودند که می‌بایست با روشهایی مانند حذف یا جایگذاری از دست آنها خلاص شد، ولی از آن به بعد آماردانان به داده‌های گمشده به عنوان یک منبع تغییرات که نیاز به محاسبه مکرر و روشهای محاسبه احتمالی دارد، نگاه کردند.

در هر مجموعه داده گمشده، مقادیر مشاهده شده اطلاعات غیر مستقیمی را در مورد مقادیر احتمال داده‌های گمشده در اختیار می‌گذارند که اگر با فرضهای اساسی که توضیح داده خواهد شد، ترکیب گردند توزیع احتمال را برای مقادیر گمشده نتیجه می‌دهند که باید در تحلیلهای آماری میانگین‌گیری شود. روش EM این میانگین‌گیری را انجام می‌دهد.

انتخاب نام EM به این علت است که در هر تکرار الگوریتم یک مرحله امید ریاضی‌گیری^۸ و بعد از آن یک ماکسیم‌سازی^۹ انجام

می‌گیرد [۲]. الگوریتم EM از رابطه بین داده‌های گمشده و پارامترهای نامعلوم مدل داده‌ها استفاده می‌کند. بدین صورت که اگر داده‌های گمشده را بدانیم، در این صورت می‌توانیم با استفاده از آنها پارامترهای مدل را مستقیماً برآورد کنیم. همچنین، اگر پارامترهای مدل را داشته باشیم، می‌توان داده‌های گمشده را برآورد کرد. این رابطه بین پارامترهای مدل و داده‌های گمشده ما را بر آن می‌دارد که ابتدا مقادیری را برای پارامترها فرض کرده و مقادیر گمشده را برآورد کنیم. سپس با جایگذاری آنها پارامترها را برآورد کنیم و این کار را تا همگرایی در برآوردها ادامه دهیم. در این صورت برآوردهای به دست آمده با برآوردهای ماکسیم درستی با در نظر داشتن میانگینی از توزیع‌های مقادیر گمشده یکسان خواهند بود. اگر درصد داده‌های گمشده زیاد باشد، در این صورت همگرایی با تکرار بیشتری حاصل می‌شود.

یک راه برای بررسی همگرایی در EM آزمایش صعودی بودن تابع لگاریتم درستنمایی در هر مرحله است و همچنین بهتر است با مقادیر شروع مختلف و به دست آوردن نتایج یکسان اطمینان حاصل کنیم که لگاریتم درستنمایی یک ماکسیم منحصر به فرد دارد، زیرا ممکن است دارای مدهای چندگانه باشد و این در صورتی اتفاق می‌افتد که تعداد نمونه کم و یا تعداد داده‌های گمشده زیاد و یا تعداد پارامترهای مدل زیاد باشند.

۲-۳-۲- داده افزایی^{۱۰}

مانند روش EM روش داده افزایی یک محاسبه تکراری است که به طور متناوب داده‌های گمشده را جایگذاری کرده و پارامترهای ناشناخته را به صورت یک روند تصادفی پیش‌بینی می‌کند. داده افزایی ابتدا بر اساس مقادیر فرضی پارامترها یک جایگذاری ابتدایی را برای داده‌های گمشده در نظر می‌گیرد و سپس پارامترهای جدید را توسط توزیع پسین به دست آمده از داده‌های کامل برآورد می‌کند. این رویه شبیه‌سازی متناوب پارامترها و داده‌های گمشده یک زنجیر مارکوف تولید می‌کند که سرانجام تثبیت شده یا در توزیع همگرا می‌شود.

^۹ Maximization

^{۱۰} Data Augmentation

^۷ Rubin, Dempster, Laird

^۸ Expectation

در روش جانهای چندگانه هر مقدار گمشده با مجموعه‌ای از مقادیری که از توزیع پیش‌بینی شده، به دست آمده‌اند، جایگذاری می‌شوند.

تغییرات موجود بین m جایگذاری، نشان دهنده این است که فرض عدم وجود حتمیت برای مقادیر پیش‌بینی در نظر گرفته شده است. بعد از به کار بردن جانهای چندگانه به تعداد m مجموعه داده‌های کامل خواهیم رسید که هر کدام را می‌توان توسط روشهای آماری موجود تحلیل کرد و نتایج به دست آمده از هر مجموعه داده را توسط قوانین ساده ارائه شده روین در سال ۱۹۸۷ ترکیب کرد.

روش جانهای چندگانه به چند دلیل جالب است؛

□ این روش برای به کار بردن تحلیل‌های آماری و نرم‌افزارهایی که نیاز به مجموعه داده کامل دارند، بسیار مناسب است.

□ هر کدام از m جایگذاری می‌تواند در تحلیل‌های مختلف استفاده شود و نیازی به تولید دوباره مجموعه‌های داده‌ها توسط روش جانهای چندگانه برای انجام تحلیل‌های جدید وجود ندارد.

□ انحراف معیارها و مقادیر احتمال^{۱۲} و غیره که از روش جانهای چندگانه به دست می‌آیند، معتبر خواهند بود. چرا که عدم وجود حتمیت برای داده‌های جایگذاری شده را در نظر گرفته‌ایم.

□ این روش در تعداد کمی از تکرار کارآیی بسیار بالایی برای به دست آوردن نتایج دقیق ارائه می‌دهد (در اغلب موارد فقط ۳ تا ۵ تکرار کافی است!).

شاید این تعداد تکرار کم در تولید مجموعه‌های داده‌ها باعث تعجب شود، ولی روین [۸، صفحه ۱۱۴] نشان داده است که کفایت یک برآورد بر اساس m جایگذاری تقریباً برابر است با

$$\left(1 + \frac{\gamma}{m}\right)^{-1}$$

که در آن، γ کسری از اطلاعات گمشده^{۱۳} با توجه به مقادیر برآورد شده است. یعنی، در صورت عدم وجود داده گمشده چه مقدار ممکن

توزیع پارامترها به یک توزیع پسین که در اصل میانگینی از داده‌های گمشده در آن نقش دارند، می‌رسد. همچنین، توزیع داده‌های گمشده نیز به سمت توزیع پیش‌بینی کننده‌ای که در حقیقت همان چیزی است که برای به دست آوردن جانهای چندگانه دقیق لازم است، پیش می‌رود. دقیقاً مانند EM در داده افزایشی نیز تعداد تکرار لازم برای همگرایی به میزان اطلاعات داده‌های گمشده ربط دارد.

توجه شود که همگرایی در داده افزایشی با همگرایی در روش EM متفاوت است. چرا که در EM همگرایی زمانی حاصل می‌شود که پارامترها از یک تکرار تا تکرار بعدی تغییر نکنند، ولی همگرایی در داده افزایشی زمانی است که توزیع پارامترها از مرحله‌ای به بعد در صورت تغییر مقادیر تصادفی پارامتر ثابت باقی بماند. به همین علت، تشخیص همگرایی در داده افزایشی پیچیده‌تر از روش EM است.

- تشخیص همگرایی در داده افزایشی

همگرایی در داده افزایشی به صورت تاخیری در همبستگی دنباله‌ای تعریف می‌شود. به این صورت که داده افزایشی در تکرار k ام همگرا خواهد شد. اگر مقدار هر پارامتر در تکرار t ام به صورت آماری مستقل از مقدار آن در تکرار $(t+k)$ ام برای $t = 1, 2, 3, \dots$ باشد.

ابتدا مقادیر پارامترها را در هر تکرار ثبت کرده، سپس با نمودار تابع خود همبستگی^{۱۱} سری زمانی می‌توان به راحتی ضریب همبستگی پیرسن تاخیر k ام را برای مقادیر مختلف k به دست آورد، یعنی همبستگی بین مقادیر شبیه‌سازی یک پارامتر در هر تکرار و مقادیر آن در k تکرار قبل. اگر مقادیر تابع خود همبستگی بعد از تاخیر k برای تمام پارامترها از k ام به صفر برسند، در این صورت می‌گوییم داده افزایشی در تکرار k ام همگرا است. نکته جالب اینجاست که روش داده افزایشی سریعتر از EM به همگرایی می‌رسد.

۲-۳-۳- استنباط در مورد جانهای چندگانه

روین در سال ۱۹۸۷ نمونه‌ای از جانهای چندگانه را گسترش داد که بر میانگین‌گیری بر اساس شبیه‌سازی استوار بود.

^{۱۲} P-Value

^{۱۳} Fraction of Missing Information

^{۱۱} ACF

دو نکته را باید مورد توجه قرار داد

۱. مدل در نظر گرفته شده باید با تحلیلی که می‌خواهیم بر مجموعه‌های داده‌ها انجام دهیم، سازگار باشد (حداقل به طور تقریبی).
۲. مدل در نظر گرفته شده باید روابط متغیرها را حفظ کند.

برای توضیح بیشتر مثال ساده زیر را در نظر می‌گیریم.

مثال: فرض کنید Y با یک مدل نرمال و با در نظر گرفتن متغیر X_1 جایگذاری شده است. بعد از جایگذاری می‌خواهیم از رگرسیون خطی برای پیش‌بینی Y توسط X_1 و متغیر دیگری مانند X_2 که در مدل جایگذاری حضور نداشته است، استفاده کنیم. در این صورت، برآورد ضریب رگرسیونی X_2 به سمت صفر میل می‌کند. چرا که Y تنها بر اساس متغیر اول جایگذاری و پیش‌بینی شده است. پس لازم است که در مدل جایگذاری تمامی متغیرهای مورد نیاز در تحلیل حضور داشته باشند.

۳-۲- توزیع پسین

تئوری جانهای چندگانه قوانین اساسی از احتمال که به نظریه بیز مشهور است را شامل می‌شود. برای این منظور نیاز به توزیع پسین برای پارامترهای وارد شده در مدل داریم که در صورت بزرگ بودن حجم نمونه تقریباً تمامی توزیع‌های توجیه‌پذیر پسین منجر به نتایج یکسان می‌شوند. در غیر این صورت همانطور که می‌دانیم هر توزیع پسینی نتایج متفاوتی خواهد داد که باید به طور دقیق تشخیص داده شود (برای مطالعه روش انتخاب توزیع پسین می‌توان به مرجع [۱۰] مراجعه کرد).

۳-۳- مکانیزم مقادیر بی پاسخ

در روشهای جایگذاری فرض می‌شود که سیستم از دست رفتن داده‌ها قابل چشم‌پوشی^{۱۸} [۸] و یا دارای فرضیات گمشدگی تصادفی^{۱۹} است [۶].

است (با توجه به برآورد موجود) دقت برآورد بیشتر شود. روش محاسبه γ در انتها نشان داده خواهد شد. کارآیی به دست آمده از مقادیر مختلف γ و m در جدول ۱ نشان داده شده است و می‌توان دید که کارآیی بالایی بعد از تعداد کمی تکرار به دست می‌آید.

روش جانهای چندگانه در حدود ۲۰ سال قبل [۷] برای اولین بار ارائه شد، ولی در طی این سالها این روش به علت عدم وجود نرم‌افزار مناسب برای تولید جانهای چندگانه چندان مورد استفاده قرار نگرفت. چرا که به دست آوردن توزیع احتمال، که داده‌ها باید از آن تولید شوند، پیچیده و بفرنج به نظر می‌رسید. اخیراً با معرفی روشهای شبیه‌سازی که به نام زنجیر مارکوف مونت کارلو^{۱۴} معروف هستند، تحولی در مدل‌های پارامتری کاربردی ایجاد شده است [۴].

حال به بررسی دقیق‌تر روشهای جایگذاری و فرضیات لازم می‌پردازیم.

۳. فرضیات

برای بررسی فرضیات لازم است تا اطلاعات کاملی نسبت به موارد زیر داشته باشیم.

- مدل داده‌ها
- توزیع پسین برای پارامترهای مدل
- چگونگی و خصوصیات مقادیر بدون پاسخ.

۳-۱- مدل داده‌ها

تحلیل گران می‌دانند که احتمال برآزنده شدن مدل‌های ساده‌ای مانند نرمال چند متغیره به داده‌ها بسیار کم است. ولی در روشهای جانهای چندگانه نشان داده شده است که حتی با این فرض تقریبهای بسیار خوبی به دست می‌آید.

برای مثال، فرض کنید بر روی متغیرهایی ترتیبی یا دو حالتی کار می‌کنیم. اغلب قابل قبول است که مقادیر را با فرض نرمال برآورد و سپس به نزدیکترین گروه یا عدد گرد کنیم. متغیرهای با توزیع شدیداً چوله را می‌توان توسط تبدیلاتی مانند لگاریتم به نرمال نزدیک کرد.

مطالعات شبیه‌سازی برای روشهای استوار^{۱۵} جانهای چندگانه توسط ازتی و رایس^{۱۶} [۳] و شیفر^{۱۷} [۱۰] انجام گرفته است.

^{۱۸} Ezzati and Rice

^{۱۷} Schafer

^{۱۸} Ignorable

^{۱۹} Missing at Random (MAR)

^{۱۴} Makov chain Monte Carlo (MCMC)

^{۱۵} Robust

ممکن است در نگاه اول از گمشدگی تصادفی این گونه استنباط شود که مقادیر گمشده باید یک ریز نمونه تصادفی از مجموعه کل داده‌ها باشند. در صورتی که این تعریف مربوط به داده‌های گمشده کاملاً تصادفی است و اغلب در دنیای واقعی اینگونه نیست. داده‌های گمشده را کاملاً تصادفی گوئیم، اگر احتمالات متغیر پاسخ، به داده‌های گمشده یا مشاهده شده ربطی نداشته باشد.

گمشدگی تصادفی به این معنی است که احتمالات از دست رفتن داده‌ها با مقادیر داده‌هایی که مشاهده شده‌اند، ربط داشته ولی با مقادیر گمشدگی رابطه‌ای نداشته باشند. برای روشن تر شدن موضوع مجموعه داده‌های دو متغیره را در نظر بگیرید که X همیشه مشاهده می‌شود و Y ممکن است گاهی دارای داده گمشده باشد. تحت فرض گمشدگی تصادفی می‌توان Y را به وسیله X و روش رگرسیون خطی پیش‌بینی کرد و داده‌های گمشده را با استفاده از این مدل رگرسیون جایگذاری نمود. در واقع، گمشدگی تصادفی کمک می‌کند تا احتمالات از دست رفتن مشاهدات Y را تنها با مقادیر داده‌هایی از X که مشاهده شده‌اند، مرتبط دانسته و مطمئن باشیم که مقادیر گمشده تاثیری در احتمالات نمی‌گذارند.

در این صورت داده‌های گمشده دارای خاصیت گمشدگی تصادفی هستند. اگر برای هر مقدار ϕ ، $g_\phi(\tilde{m}|u)$ مقدار یکسانی برای تمام $u_{(0)}$ اختیار کند. یعنی برای هر مقدار ممکن پارامتر ϕ ، احتمال شرطی الگوی داده گمشده به شرط داده‌های گمشده، به علاوه مقادیر مشاهده شده برای تمام مقادیر داده‌های گمشده یکسان باشد.

داده‌ها دارای خاصیت مشاهده شدگی تصادفی^{۲۱} هستند، اگر برای هر مقدار ϕ و $u_{(0)}$ ، $g_\phi(\tilde{m}|u)$ مقدار یکسانی را برای تمامی $U_{(1)}$ اختیار کند.

حال فرض کنید که

۱. داده‌های گمشده دارای خاصیت گمشدگی تصادفی هستند

۲. داده‌ها دارای خاصیت مشاهده شدگی تصادفی هستند.

در این صورت با فرض $k_{\theta,\phi}(\tilde{m}) > 0$ داریم،

$$\int f_\theta(u) du_{(0)} = \int \left\{ \frac{f_\theta(u) g_\phi(\tilde{m}|u)}{k_{\theta,\phi}(\tilde{m})} \right\} du_{(0)}$$

که در آن، $k_{\theta,\phi}(\tilde{m}) = \int f_\theta(u) g_\phi(\tilde{m}|u) du$ می‌باشد. یعنی احتمال حاشیه‌ای $M = \tilde{m}$

سمت راست تساوی نشان دهنده توزیع نمونه مورد بررسی ما با در نظر گرفتن جریان به وجود آورنده داده گمشده است و سمت چپ تساوی نشان دهنده توزیع نمونه‌ای با توجه به داده‌های مشاهده شده می‌باشد. دیده می‌شود که با دو شرط (۱) و (۲) می‌توان این دو توزیع را برابر در نظر گرفت. چرا که با توجه به این دو شرط $g_\phi(\tilde{m}|u)$ برای هر ϕ مقدار یکسانی را برای تمامی u ها اختیار می‌کند (توجه کنید که این مطلب استقلال توزیع‌های U و M را فقط در صورتی که موضوع برای تمامی \tilde{m} برقرار باشد، نتیجه می‌دهد). بنابراین،

در حقیقت، گمشدگی تصادفی فرضیه‌ای است که به ما اجازه می‌دهد تا ابتدا برآوردی از ارتباط بین متغیرها برای داده‌های مشاهده شده به دست آورده و سپس این ارتباط را برای به دست آوردن برآوردهایی نااریب از مقادیر گمشده به کار ببریم (این مطلب طی چند قضیه با فرض شرایط گمشدگی تصادفی در مرجع [۶] آمده است).

نمادهای زیر را که در بیشتر مقالات مورد استفاده قرار گرفته‌اند، در نظر بگیرید؛

$U = (U_1, U_2, \dots, U_n)$ یک بردار متغیرهای تصادفی با تابع چگالی احتمال f_θ باشد. هدف به دست آوردن استنباطی در مورد θ (بردار پارامترهای تابع چگالی احتمال f_θ است).
 $M = (M_1, M_2, \dots, M_n)$ بردار متغیرهای تصادفی نشانگر داده گمشده است که در آن M_i دو مقدار صفر و یک را اختیار می‌کند. احتمال اینکه M برابر با $m = (m_1, m_2, \dots, m_n)$ باشد، به شرط اینکه U مقدار $u = (u_1, u_2, \dots, u_n)$ را اختیار کند، برابر با $g_\phi(m|u)$ است که در آن ϕ بردار پارامترهای مزاحم توزیع است. توزیع شرطی

^{۲۰} Process that causes missing data

^{۲۱} Observed at Random (OAR)

خطی یا لوزستیک را توسط نرم‌افزارهای آماری انجام می‌دهیم. هر مدل باید m بار یعنی به تعداد جایگذاریهای داده‌های گمشده برآزش شود. در این صورت، نتایج به دست آمده از مجموعه‌های کامل داده‌ها، انعکاسی از عدم وجود حتمیت برای داده‌های گمشده را دربر خواهد داشت. حال، برای به دست آوردن ضرایب و انحراف معیارهای برآورد شده از قوانین روبین برای ترکیب نتایج و به دست آوردن برآوردهای نهایی استفاده می‌شود [۸].

قوانین روبین به صورت زیر است؛

فرض کنید \hat{Q} برآورد کیفیت مورد علاقه جامعه و U برآورد واریانس آن باشد. برای مثال، می‌توان \hat{Q} را برآورد ضریب رگرسیونی و U را واریانس آن در نظر گرفت.

بعد از انجام تحلیل یکسان بر روی هر مجموعه داده کامل شده به تعداد m برآورد $\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_m$ و همچنین واریانسهای آنها U_1, U_2, \dots, U_m خواهیم داشت. برآورد جهانی چندگانه به صورت

$$\bar{Q} = \frac{1}{m} \sum_{i=1}^m \hat{Q}_i \text{ است.}$$

واریانس کل برآورد دارای دو مؤلفه است که یکی تغییرات درونی هر مجموعه داده‌ها (واریانس درون-جایگذاری) و دیگری تغییرات بین مجموعه‌های داده (واریانس بین-جایگذاری) را مشخص می‌کند.

واریانس درون-جایگذاری برابر است با، $\bar{U} = \frac{1}{m} \sum_{i=1}^m U_i$ که میانگین واریانسهای برآورد شده است.

واریانس بین-جایگذاری به صورت $B = \sum_{i=1}^m (\hat{Q}_i - \bar{Q})^2 / (m-1)$ است که واریانس نمونه‌ای خود برآوردها می‌باشد. واریانس کل (T) مجموع این دو مؤلفه به علاوه ضریب تصحیحی است که برای خطای شبیه‌سازی \bar{Q} در نظر گرفته می‌شود.

$$T = \bar{U} + \left(1 + \frac{1}{m}\right) B$$

توجه کنید که اگر داده گمشده نداشته باشیم $\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_m$ مشخص و برابرند و $B = 0$ و T برابر با \bar{U} خواهد بود که \bar{U} با هر کدام از U_i ها برابر است.

$$k_{\theta, \phi}(\bar{m}) = \int f_{\theta}(\bar{m}|u) du \\ = g_{\phi}(\bar{m}|u) \int f_{\theta}(u) du = g_{\phi}(\bar{m}|u) = c$$

$$\int \left\{ \frac{f_{\theta}(u) g_{\phi}(\bar{m}|u)}{k_{\theta, \phi}(\bar{m})} \right\} du_{(o)} = \int \left\{ \frac{f_{\theta}(u) g_{\phi}(\bar{m}|u)}{g_{\phi}(\bar{m}|u)} \right\} du_{(o)} \\ = \int f_{\theta}(u) du_{(o)}$$

بر اساس فرضیه گمشدگی تصادفی، چند نکته را باید در نظر داشت.

□ فرض گمشدگی تصادفی بر اساس چگونگی روابط متغیرهای موجود در مجموعه‌های داده‌ها قابل قبول است. بدین معنی که اگر متغیر X با داده‌های گمشده متغیرهای دیگر و همچنین مقادیر مشاهده شده آنها ارتباط داشته باشد و در مجموعه‌های داده‌ها موجود نباشد، فرض گمشدگی تصادفی برای استفاده از روش جایگذاری درست نیست. به همین دلیل در یک روند جایگذاری تمامی متغیرهای مرتبط با مقادیر گمشده را می‌بایست در نظر گرفت.

□ فرض گمشدگی تصادفی را نمی‌توان با داده‌های موجود آزمون کرد، مگر اینکه داده‌های گمشده را نیز داشته باشیم. بنابراین، نمی‌توان برقرار بودن فرض گمشدگی تصادفی را به راحتی بررسی کرد، ولی به علت اینکه بتوان از روشهایی آماری برای پیش‌بینی مقادیر گمشده بر اساس داده‌های موجود استفاده کرد، در بیشتر مواقع این فرض را برقرار در نظر می‌گیریم.

در حال حاضر، روشی مبنی بر عدم نادیده گرفتن داده‌های گمشده وجود ندارد و تا زمانی که روشی پیشنهاد گردد، بهتر است از فرضیه‌هایی مانند گمشدگی تصادفی که داده‌های گمشده را نادیده گرفته و پیش‌بینی را بر اساس داده‌های موجود انجام می‌دهد، استفاده گردد. در ضمن اینکه گمشدگی تصادفی برآوردهای بسیار بهتری نسبت به روشهایی مانند حذف یا جایگذاری میانگین پیشنهاد می‌کند (برای توضیحات بیشتر در مورد روشهای عدم نادیده گرفتن داده‌های گمشده می‌توان به مرجع [۱۰] صفحات ۲۷ تا ۲۸ مراجعه کرد).

۴. قوانین جهانی چندگانه

با استفاده از قوانین جهانی چندگانه می‌توان تمامی تحلیل‌های آماری که نیاز به مجموعه داده کامل دارند را انجام داد. فرض کنیم رگرسیون

با روش حذف برآوردها شدیداً اریب خواهند شد. به علاوه، همانطور که گفته شد در تحلیلهای چند متغیره روش حذف شدیداً حجم نمونه را کاهش می‌دهد.

در حقیقت روش جانهای چند گانه شباهت نزدیکی با الگوریتم EM و روشهای دیگر عددی دارد که برآوردهای ماکسیمم درستنمایی را بر اساس تنها مقادیر مشاهده شده محاسبه می‌کنند. این روشها تابع درستنمایی که در اصل متوسطی است بر توزیعهای پیش‌بینی شده، برای داده‌های گمشده را محاسبه می‌کنند. روش جانهای چند گانه شبیه همین متوسط‌گیری را به جای روشهای عددی به وسیله مونت کارلو انجام می‌دهد.

شاید این سوال مطرح شود که آیا جانهای چند گانه ربطی به زنجیر مارکف مونت کارلو دارد یا خیر؟ در جواب این سوال باید گفت، زنجیر مارکف مونت کارلو یک روش شبیه‌سازی است که برای به دست آوردن مقادیر تصادفی از توزیعهای غیر استاندارد به وسیله زنجیر مارکف به کار می‌رود. این روش با شبیه‌سازی پارامترها از یک توزیع پسین بیزی برای یک مدل پارامتری پیچیده در اصل این امکان را به ما می‌دهد تا در روش جانهای چند گانه تعدادی مقادیر مستقل برای داده‌های گمشده توسط توزیعی که برآورد شده است، به دست آوریم و از این مقادیر برای استنباط جانهای چند گانه استفاده کنیم.

توجه کنید که فرض گمشدگی تصادفی را به صورت ساده می‌توان اینگونه تصور کرد که این فرض به ما این اطمینان را می‌دهد که داده‌های گمشده دارای هیچ اطلاعاتی از توزیع احتمال مقادیر گمشده نبوده‌اند. می‌توان گفت که اگر بی‌پاسخ ماندن متغیری به داده‌های فعلی یا قبلی آن ربط نداشته باشد، فرض گمشدگی تصادفی قابل دفاع است. در آخر می‌توان گفت که روش جانهای چند گانه یک روش ساختن از هیچ چیز است. در اصل، این روش نیز مانند روشهای دیگر شبیه‌سازی تنها فرض احتمالی داده‌های جایگزین را در نظر گرفته و نمی‌توان گفت که معجزه‌ای برای دیدن داده‌هایی که در اصل گمشده‌اند، رخ داده است. البته نباید فراموش کرد که در نظر گرفتن همین فرض نیز در محاسبات تأثیر به‌سزایی دارد و بسیار مهم است.

فاصله اطمینان ۹۵ درصد به صورت $\bar{Q} \pm 2\sqrt{T}$ به دست می‌آید که به صورت دقیق‌تر برابر است با

$$\bar{Q} \pm t_{df} \sqrt{T}, \quad df = (m-1) \left(1 + \frac{m\bar{U}}{(m+1)B} \right)^2.$$

توجه کنید که اگر تعداد تکرارها بی‌شمار باشد ($m = \infty$) واریانس کل به مجموع دو مولفه واریانس کاهش می‌یابد و فاصله اطمینان بر اساس توزیع نرمال خواهد بود. اگر مقدار درجه آزادی خیلی کوچک باشد (کمتر از ۱۰) با افزایش تعداد تکرارها می‌توان دقت برآوردها را افزایش و طول فاصله اطمینان را کاهش داد.

رویین نشان داده است که یک برآورد برای نسبت اطلاعات گمشده برای کمیت جامه Q برابر است با

$$r = \frac{r + \frac{2}{df + 3}}{r + 1}$$

و در آن $r = (1 + m^{-1})B/\bar{U}$ که مشخص می‌کند چه مقدار برآورد Q ممکن است توسط داده‌های گمشده تحت تأثیر قرار گیرد [A].

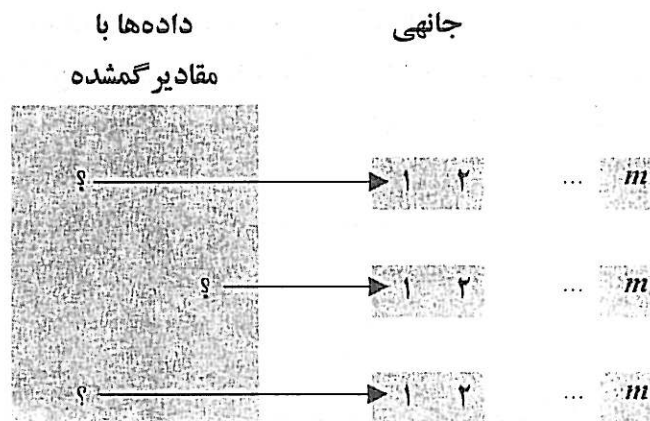
۵. نکاتی در مورد روش جانهای چند گانه

روش جانهای چند گانه یکی از روشهای برخورد با داده‌های گمشده است و لزوماً در تمامی شرایط بهترین روش نیست. برای مثال، زمانی که مدل‌های کاملاً پارامتری در اختیار داریم، اغلب برآوردهای ماکسیمم درستنمایی را می‌توان به راحتی با استفاده از روشهای عددی مانند EM محاسبه کرد.

همان طور که می‌دانیم روش حذف داده‌های گمشده بسیار ساده‌تر از روش جانهای چند گانه است. اگر داده‌های مشاهده شده نماینده خوبی از کل داده‌ها باشد و سهم داده‌های گمشده از کل داده‌ها زیاد نباشد، شاید روش حذف مناسب به نظر برسد. لیکن، در کل روش حذف داده‌های گمشده تنها زمانی کارآیی لازم را دارد که دارای داده‌های گمشده کاملاً تصادفی باشیم. به بیان دیگر، روش حذف فرض می‌کند که واحدهای حذف شده یکی زیر نمونه تصادفی هستند. وقتی که واحدهای حذف شده به طور سیستماتیک با بقیه داده‌ها متفاوت باشند،

جدول ۱

γ					
m	۰/۱	۰/۳	۰/۵	۰/۷	۰/۹
۳	۹۷	۹۱	۸۶	۸۱	۷۷
۵	۹۸	۹۴	۹۱	۸۸	۸۵
۱۰	۹۹	۹۷	۹۵	۹۳	۹۲
۲۰	۱۰۰	۹۹	۹۸	۹۷	۹۶



ماتریس داده‌های چند متغیره با مقادیر گمشده و

جانهی چند گانه

مراجع

[1] Ashofteh, A., Jahanshahi, M.A. and Bozorgnia, A., 2002. *Estimation of Systematic Missing Values by Simulation*, Submitted.

[2] Dempster, A.P., Laird, N.M. and Rubin, D.B., 1977. *Maximum Likelihood from Incomplete Data Via the EM Algorithm*, Journal of the Royal Statistical Society, Series B, 39, pp.1-38.

[3] Ezzati-Rice, TM., 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. Bureau of the Census, Washington, D.C.

[4] Gilks, W.R., Richardson, S. and Spiegelhalter, D.J., (Eds.), 1996. *Markov Chain Monte Carlo in Practice*, Chapman & Hall, London.

[5] Meng, X.L. and Rubin, D.B., 1993. *Maximum Likelihood Estimation Via the ECM Algorithm. A General Framework*, Biometrika, 80, pp.267-278.

[6] Rubin, D.B., 1976. *Inference and Missing Data*, Biometrika, 63, pp.581-592.

[7] Rubin, D.B., 1978. *Multiple Imputations in Sample Surveys- A Phenomenological Bayesian Approach to Nonresponse*, Proceedings of the Survey Research Methods Section, American Statistical Association, pp.20-34.

[8] Rubin, D.B., 1987. *Multiple Imputations for Nonresponse in Surveys*, John Wiley & Sons, New York.

[9] Rubin, D.B., 1996. *Multiple Imputations After 18+ Years (with discussion)*, J. A. Stat. Asso., 91, pp. 473-489.

[10] Schafer, J.I., 1997. *Analysis of Incomplete Multivariate Data*, Chapman & Hall, London.

[11] Tanner, M.A. and Wong, W.H., 1987. *The Calculation of Posterior Distributions by Data Augmentation (with discussion)*, J. Amer. Statist. Asso., 82, pp. 528-550.