

جانه‌ی چندگانه‌ی غیر بیزی

جان جورن استد[†]

اداره‌ی آمار نروژ

مترجم: زهرا رضایی قهرودی

پژوهشکده‌ی آمار

چکیده: جانه‌ی چندگانه روشی است که به طور مشخص برای برآورد واریانس در حضور داده‌های گم‌شده طراحی شده است. فرمول ترکیبی روبین نیازمند آن است که روش جانه‌ی «مناسب» باشد، که اساساً به این معنی است که جانه‌ی‌ها، استخراج تصادفی از یک توزیع پسین در یک چارچوب بیزی باشند. در مؤسسات ملی آمار (NSI's National Statistical Institutes) جانه‌ی بی‌پاسخی استفاده می‌شوند، مانند بعضی از انواع جانه‌ی‌های بی‌درنگ طبقه‌بندی شده، نوعاً غیر بیزی هستند. بنا بر این روش جانه‌ی چندگانه‌ی روبین در مؤسسات ملی آمار معتبر نیست و نمی‌تواند استفاده شود. این مقاله به مسئله‌ی تعیین یک فرمول ترکیبی دیگر می‌پردازد که بتواند برای روش‌های جانه‌ی که نوعاً در مؤسسات ملی آمار استفاده می‌شود، به کار رود و رهیافتی برای مطالعه‌ی این مسئله پیشنهاد می‌کند. فرمول‌های ترکیبی دیگری برای مکانیسم‌های خاص پاسخ و روش‌های جانه‌ی از نوع بی‌درنگ تعیین شده‌اند.

[†] Bjørnstad, J. F. (2007). Non-Bayesian multiple imputation. *Journal of Official Statistics* 23, 433-452.

واژگان کلیدی: برآورد واریانس؛ آمارگیری نمونه‌ای؛ نمونه‌گیری طبقه‌بندی شده؛ رگرسیون لوئیستیک؛ بی‌پاسخی؛ جانه‌ی بی‌درنگ.

دریافت: ۱۳۸۸/۲/۲، پذیرش: ۱۳۸۸/۵/۱۷

۱ - مقدمه

جانهی چندگانه روشی است که به طور مشخص برای براورد واریانس در حضور داده‌های گم شده طراحی شده و توسط روین [۲] بسط یافته است. دو مرجع جدید دیگر با بحث‌ها و بررسی‌های بیشتر، روین [۳] و شیفر [۴] هستند. ایده‌ی اصلی این روش ایجاد m مقدار جانهی شده برای هر مقدار گم شده و ترکیب m مجموعه‌ی داده‌های کامل شده به‌وسیله‌ی فرمول ترکیبی روین برای براورد واریانس است. برای این‌که براوردگر معتبر باشد، جانهی باید سطح مناسبی از تغییرپذیری را نمایش دهد. در اصطلاح روین، یکی از الزامات روش جانهی «مناسب بودن» است. در مؤسسات ملی آمار روش‌هایی که برای جانهی بی‌پاسخی استفاده می‌شود، حتی اگر الزام «مناسب بودن» را داشته باشند، خیلی کم این نیاز را تأمین می‌کنند. اما ایده‌ی ایجاد جانهی چندگانه برای اندازه‌گیری عدم حتمیت جانهی و استفاده از آن برای براورد واریانس و محاسبه‌ی بازه‌ی اطمینان هنوز مورد علاقه است. مشکل این است که فرمول ترکیبی روین حالا دیگر برای جانهی‌های نامناسب معمول که توسط مؤسسات ملی آمار استفاده می‌شود، معتبر نیست. دلیل آن این است که تغییرپذیری در جانهی‌های نامناسب خیلی کم است و باید برای مؤلفه‌ی بین جانهی، وزن بیشتری در براورد واریانس در نظر گرفته شود. مشکل دیگر این است که مشخص شود این وزن چقدر باید باشد تا استنباط آماری معتبری ارائه دهد و همچنین برای چه نوع مکانیسم‌های بی‌پاسخی و مسایل براورد، امکان تعیین فرمول ترکیب ساده که به پارامترهای نامعلوم بستگی نداشته باشد، وجود دارد. این مقاله رهیافتی برای مطالعه‌ی این مسئله پیشنهاد می‌کند.

در بخش دوم رهیافتی برای تعیین ترکیب مجموعه‌ی داده‌های کامل شده از طریق جانهی پیشنهاد می‌شود. بخش سوم سه کاربرد با بی‌پاسخی تصادفی را بیان می‌کند:

(الف) براورد میانگین جامعه از روی نمونه‌های تصادفی ساده با استفاده از روش جانهی

بی‌درنگ؛

- (ب) براورد ضریب رگرسیونی در مدل نسبت، با استفاده از جانهی رگرسیونی مانده‌ها؛
- (پ) براورد ضریب رگرسیون خطی ساده با جانهی رگرسیونی مانده‌ها.

بخش چهارم به مسئله‌ی کلی جانهی چندگانه برای نمونه‌های طبقه‌بندی شده می‌پردازد. در بخش پنجم از نظریه‌ی بخش چهارم برای نمونه‌های طبقه‌بندی شده با بی‌پاسخی تصادفی درون طبقات استفاده می‌شود که موارد زیر را پوشش می‌دهد:

الف) براورد میانگین جامعه با استفاده از جانهی بی‌درنگ طبقه‌بندی شده؛

ب) براورد لگاریتم (نسبت بخت‌ها) در رگرسیون لوژستیک با وجود گم‌شدگی هم در متغیرهای وابسته و هم در متغیرهای تبیینی.

بخش ششم مسئله‌ی استفاده از قاعده‌ی ترکیبی یکسان برای همه‌ی مسائل براورد با یک روش جانهی مشخص و داده‌ها و مدل پاسخ معلوم را بیان می‌کند. یک نتیجه‌ی کلی برای جانهی بی‌درنگ و براورد خطی نیز ارائه شده است.

۲- رهیافتی برای تعیین یک فرمول ترکیبی دیگر برای براورد واریانس در جانهی چندگانه

فرض کنید $s = \{y_1, \dots, y_n\}$ معرف کل نمونه با داده‌های کل نمونه‌ی $(y_i : i \in S_r)$ از مقادیر متغیرهای تصادفی Y_1, \dots, Y_n باشند. در نمونه‌گیری از جامعه‌ی متناهی تحت یک مدل طرح، شماره‌گذاری مجدد واحدهای انتخاب شده، انجام می‌شود و ماهیت تصادفی بودن y به وسیله‌ی طرح نمونه‌گیری تعیین می‌شود. هدف براورد برخی از پارامترهای θ است. داده‌های مشاهده شده که به وسیله‌ی $y_{obs} = \{y_i : i \in S_r\}$ نمایش داده می‌شوند، بخش مشاهده شده‌ی y و نمونه‌ی پاسخ s از اندازه‌ی n_r است.

فرض کنید $\hat{\theta}$ براورده‌گری براساس داده‌های کل نمونه‌ی y باشد که با $\text{var}(\hat{\theta})$ به وسیله‌ی $\hat{V}(y)$ براورد شده است. برای $i \in S_r - S_r^*$, y_i^* را به وسیله‌ی بعضی روش‌ها جانهی می‌کنیم و y^* را که معرف داده‌های کامل شده است، در نظر می‌گیریم. $\hat{V}^* = \hat{V}(y^*)$. براساس y^* , داریم $\hat{\theta}^* = \hat{\theta}(y^*)$ و $\hat{V}^* = \hat{V}(y^*)$. جانهی چندگانه‌ی m جانهی تکرار شده، منجر به m مجموعه‌ی داده‌های کامل شده با m براورد $i = 1, \dots, m$ و $\hat{\theta}_i^*$ ، و براوردهای واریانس مربوط به آن‌ها، \hat{V}_i^* ، برای $i = 1, \dots, m$ می‌شود. براورد ترکیب شده به وسیله‌ی $\bar{\theta}^* = \sum_{i=1}^m \hat{\theta}_i^* / m$ تعیین می‌شود.

واریانس داخل جانه‌های $\bar{V}^* = \sum_{i=1}^m \hat{V}_i^* / m$ و مؤلفه‌ی بین جانه‌های $B^* = \sum_{i=1}^m (\hat{\theta}_i^* - \bar{\theta}^*)^2 / (m-1)$ تعریف می‌شوند. واریانس براورده‌ی کل $\bar{\theta}^*$ به‌وسیله‌ی

$$(1) \quad W = \bar{V}^* + \left(k + \frac{1}{m}\right) B^*.$$

پیشنهاد می‌شود. همچنین لازم است k را به‌گونه‌ای تعیین کنیم که

$$(2) \quad E(W) = \text{var}(\bar{\theta}^*).$$

روین [۲] نشان داد که $k = 1$ می‌تواند با جانه‌های مناسب استفاده شود، که اساساً به این معنی است که مقادیر جانه‌شده از یک توزیع پسین در یک چارچوب بیزی قرعه‌کشی شوند.

به‌طور کلی، باید عبارت (۲) را تعیین کرد. یک راه آزمایش و انجام آن، استفاده از امید ریاضی دوگانه به شرط y_{obs} ، است به‌طوری که $E(W) = E\{E(W|Y_{\text{obs}})\}$ و $\text{var}(W) = E\{\text{var}(W|Y_{\text{obs}})\} + \text{var}\{E(W|Y_{\text{obs}})\}$

$$(3) \quad E(\bar{V}^*) \approx \text{var}(\hat{\theta})$$

و $E(B^*|y_{\text{obs}}) = \text{var}(\hat{\theta}^*|y_{\text{obs}})$. بنا بر این به طور تقریبی

$$(4) \quad E(W) = \text{var}(\hat{\theta}) + \left(E(k) + \frac{1}{m}\right) E \text{ var}(\hat{\theta}^*|Y_{\text{obs}})$$

است. علاوه بر این، $E(\bar{\theta}^*|y_{\text{obs}}) = E(\hat{\theta}^*|y_{\text{obs}})$ و $\text{var}(\bar{\theta}^*|y_{\text{obs}}) = \frac{\text{var}(\hat{\theta}^*|y_{\text{obs}})}{m}$ این حاکی از آن است که $\text{var}(\bar{\theta}^*) = m^{-1} E\{\text{var}(\hat{\theta}^*|Y_{\text{obs}})\} + \text{var}\{E(\hat{\theta}^*|Y_{\text{obs}})\}$. با استفاده از روابط (۳) و (۴)، معادله‌ی (۲) به‌صورت

$$\text{var}(\hat{\theta}) + E(k)E[\text{var}(\hat{\theta}^* | Y_{\text{obs}})] = \text{var}\{E(\hat{\theta}^* | Y_{\text{obs}})\}$$

در می‌آید و عبارت کلی زیر را ارائه می‌دهد.

$$(5) \quad E(k) = \frac{\text{var} E(\hat{\theta}^* | Y_{\text{obs}}) - \text{var}(\hat{\theta})}{E \text{var}(\hat{\theta}^* | Y_{\text{obs}})}$$

اگر مقدار k مورد علاقه‌ی ما باشد باید حد اقل به طور تقریبی، مستقل از پارامترهای نامعلوم تعیین شود. علاوه بر این، لازم است که رابطه‌ی (۳) برقرار باشد. برای تشریح این که چگونه رابطه‌ی (۵) می‌تواند استفاده شود، در بخش بعدی سه حالت خاص با بی‌پاسخی تصادفی را مورد بررسی قرار می‌دهیم.

۳- سه کاربرد برای بی‌پاسخی تصادفی

۱- براورد میانگین جامعه با جانهی بی‌درنگ

یک نمونه‌ی تصادفی ساده از جامعه‌ی متناهی به اندازه‌ی N را در نظر می‌گیریم، که هدف آن براورد میانگین جامعه، مل از متغیر y است. همچنین فرض بی‌پاسخی کاملاً تصادفی نیز درنظر گرفته شده است. در اصطلاح‌شناسی روبین [۲] و لیتل و روبین [۱]، مکانیسم گمشدگی، گم شدن کاملاً تصادفی (Missing Completely at Random) (MCAR) نامیده می‌شود. گم شدن کاملاً تصادفی به این معنا است که نماگرهای پاسخ R_1, \dots, R_N با احتمال‌های پاسخ مشابه $(1 = \Pr(R_i = 1))$ مستقل هستند. روش جانهی، روش بی‌درنگ است که y^* به تصادف از y_{obs} با جایگذاری قرعه‌کشی شده است و براورد، میانگین نمونه است. فرض کنید \bar{y}_r میانگین نمونه‌ی مشاهده شده و

$$\bar{Y}^* = \frac{1}{n_r - 1} \sum_{i \in s_r} (y_i - \bar{y}_r)$$

میانگین نمونه بر اساس جانهی برای نمونه‌ی کامل شده است، و براوردگر ترکیب شده به وسیله‌ی $\bar{Y}^* = \sum_{i=1}^m \bar{Y}_i^* / m$ تعیین می‌شود. فرض کنید \bar{Y}_s معرف میانگین نمونه بر اساس کل نمونه باشد. بنا بر این $\text{var}(\bar{Y}_s) = \sigma^2 \left(\frac{1}{n} - \frac{1}{N} \right)$ است که در آن

$\sigma^* = (N - 1)^{-1} \sum_{i=1}^N (y_i - \mu)^*$ به عنوان واریانس جامعه است. همچنین با توجه به اینکه $\text{var}(Y_i^* | y_{\text{obs}}) = \hat{\sigma}_r^* (n_r - 1)/n_r$ و $E(Y_i^* | y_{\text{obs}}) = \bar{y}_r$ است، $\text{var}(\bar{Y}^* | y_{\text{obs}}) = \{(n - n_r)/n\} \{(n_r - 1)/n_r\} \hat{\sigma}_r^*$ و $E(\bar{Y}^* | y_{\text{obs}}) = \bar{y}_r$ خواهد بود. در این حالت، $\hat{V}^* = \hat{\sigma}_r^* (\frac{1}{n} - \frac{1}{N})$ است که در آن $\hat{\sigma}_r^* = \frac{1}{n-1} \left(\sum_{s_r} (y_s - \bar{y}^*)^2 + \sum_{s \neq s_r} (y_s^* - \bar{y}^*)^2 \right)$. میتوان نشان داد که $E(\hat{\sigma}_r^* | y_{\text{obs}}) = \hat{\sigma}_r^* \left(1 - \frac{1}{n_r} \right) \left(1 + \frac{n_r}{n(n-1)} \right) \approx \hat{\sigma}_r^*$ برقرار است. با استفاده از رابطه‌ی (۵) داریم

$$\begin{aligned} E(k) &= \frac{\text{var}(\bar{Y}_r) - \sigma^* \left(\frac{1}{n} - \frac{1}{N} \right)}{E \left(\frac{n - n_r}{n} \cdot \frac{n_r - 1}{n_r} \right) E(\hat{\sigma}_r^* | n_r)} \\ &= \frac{\sigma^* \left(E \left(\frac{1}{n_r} \right) - \frac{1}{N} \right) - \sigma^* \left(\frac{1}{n} - \frac{1}{N} \right)}{E \left(\frac{n - n_r}{n} \cdot \frac{n_r - 1}{n_r} \right) \sigma^*} \approx \frac{(1 - p_r)/p_r}{1 - p_r} = \frac{1}{p_r} \end{aligned}$$

که با در نظر گرفتن $f = (n - n_r)/n$ به عنوان نرخ بی‌پاسخی و قرار دادن $k = 1/(1-f)$ به طور تقریبی روابط مورد نظر برقرار است.

۳-۳- برآورد ضریب رگرسیون در مدل نسبت با جانبه‌ی مانده‌ها

فرض می‌کنیم بی‌پاسخی همانند بخش ۳-۱ کاملاً تصادفی است. یک مدل نسبت، یعنی رگرسیونی که از مبدأ می‌گذرد، $Y_i = \beta x_i + \varepsilon_i$ ، با $\text{var}(\varepsilon_i) = \sigma^* x_i$ ؛ $i = 1, \dots, n$ ، را در نظر می‌گیریم. فرض شده است که همه‌ی x_i ‌ها، همچنین در نمونه‌ی بی‌پاسخ

معلوم‌اند. برآوردهای داده‌های کامل β ، به وسیله‌ی $\hat{\beta} = \sum_{i=1}^n Y_i / \sum_{i=1}^n x_i$ مشخص می‌شود. برآوردهای نالریب $\sigma^* = \sum_{i=1}^n \frac{1}{x_i} (y_i - \hat{\beta} x_i)^2 / (n-1)$ به وسیله‌ی مشخص می‌شود.

روش جانهی رگرسیون مانده‌ها در نظر گرفته شده است. فرض کنید $\hat{\beta}_r$ ، برآورد $\hat{\beta}$ بر اساس نمونه‌ی مشاهده شده‌ی s_r باشد. مانده‌های استاندارد شده را به صورت $e_i^* = (y_i - \hat{\beta}_r x_i) / \sqrt{x_i}$ برای $i \in s_r$ تعریف می‌کنیم. برای $i \in s - s_r$ مقدار e_i^* را به تصادف و با جایگذاری از روی مجموعه‌ی مانده‌های مشاهده شده‌ی $i \in s_r$ ، $e_i^* = y_i - \hat{\beta}_r x_i + e_i^* \sqrt{x_i}$ قرعه‌کشی می‌کنیم. مقدار y_i جانهی شده به وسیله‌ی مشخص می‌شود.

$X_{nr} = \sum_{i \in s - s_r} x_i = X - X_r$ و $X_r = \sum_{i \in s_r} x_i$ ، $X = \sum_{i=1}^n x_i$ فرض کنید باشد. همه‌ی فرضیات از اکنون به بعد به شرط n_r و X_r است و هدف ما تعیین مقدار k به طور مستقیم از رابطه‌ی (۵) می‌باشد. نسبت x -کل در گروه بی‌پاسخی به وسیله‌ی رابطه‌ی $f_X = X_{nr} / X$ مشخص شده است. بنابراین $\hat{\beta}^* = (\sum_{s_r} y_i + \sum_{s-s_r} y_i^*) / X$

$$\hat{\sigma}_*^2 = \frac{1}{n-1} \left(\sum_{s_r} \frac{1}{X_i} (y_i - \hat{\beta}^* x_i)^2 + \sum_{s-s_r} \frac{1}{X_i} (y_i^* - \hat{\beta}^* x_i)^2 \right).$$

به منظور تعیین k از رابطه‌ی (۵)، نیاز به بررسی اعتبار رابطه‌ی (۳) و استخراج $\text{var } E(\hat{\beta}^* | y_{\text{obs}})$ است. $\text{var } E(\hat{\beta}^* | y_{\text{obs}}) = \text{Evar}(\hat{\beta}^* | y_{\text{obs}})$ است. یادآور می‌شود که در پیوست ۱ نشان داده شده است که شرط (۳) برای مقادیر بزرگ و متوسط n_r برقرار است و این که

$$(6) \quad \text{var } E(\hat{\beta}^* | y_{\text{obs}}) = \frac{\sigma^2}{X_r} + \frac{(1-d_r)d_r n_{nr} X_{nr}}{X^2} \cdot \frac{\sigma^2}{n_r}$$

$$(7) \quad \text{Evar}(\hat{\beta}^* | y_{\text{obs}}) = \frac{X_{nr}}{X^2} \cdot \frac{\sigma^2}{n_r} (n_r + d_r - 2)$$

که در آن $d_r \leq d_s \leq 1$. با استفاده از روابط (۶) و (۷) و قرار دادن آن در رابطه (۵) خواهیم داشت

$$k = \frac{n_r X^r - n_r X X_r + (1-d_s)d_s n_{nr} X_{nr} X_r}{X_r X_{nr} (n_r + d_s - 2)} \approx \frac{X^r}{X_r} + (1-d_s)d_s \frac{n_{nr}}{n_r}$$

لازم به ذکر است که اگر همه‌ی $x_i = d_s$ باشد، بنا بر این $d_s = d_r$. اکنون با درنظر گرفتن $f_X = X_{nr}/X$ که نسبت x -کل در گروه بی‌پاسخی است و $f = n_{nr}/n$ که نرخ بی‌پاسخی است، چون معمولاً $(1-d_s)d_s \approx 0$ است، سرانجام رابطه‌ی زیر برای مقادیر معمول x و نرخ‌های بی‌پاسخی به دست می‌آید.

$$k \approx \frac{1}{1-f_X} + (1-d_s)d_s \frac{f}{1-f} \approx \frac{1}{1-f_X}$$

۳-۳-برآورد ضریب رگرسیون خطی ساده با جانه‌ی مانده‌ها

همانند بخش ۳-۱ و ۳-۲ مکانیسم بی‌پاسخی با $\Pr(R_i = 1) = P(R_i = 1)$ گم شدن کاملاً تصادفی در نظر گرفته شده است. مدل رگرسیون خطی ساده با $\text{var}(\varepsilon_i) = \sigma^2$ برای $i = 1, \dots, n$ به صورت $Y_i = \alpha + \beta x_i + \varepsilon_i$ فرض شده است. همه‌ی x_i ‌ها معلوم فرض شده‌اند. همچنین فرض شده است که در آن $\bar{x} = \sum_{i=1}^n x_i / n = 0$. بنا بر این برآوردهای بر اساس کل داده‌ها به وسیله‌ی $\hat{\beta} = \sum_{i=1}^n x_i y_i / SS_x$ مشخص شده است که در آن $\hat{\alpha} = \bar{y} = \sum_{i=1}^n y_i / n$ و $SS_x = \sum_{i=1}^n x_i^2$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)^2$$

فرض کنید $\hat{\alpha}_r$ ، $\hat{\beta}_r$ برآوردهای بر اساس نمونه‌ی پاسخ باشند که از رابطه‌ی $\hat{\beta}_r = \sum_{i \in S_r} (x_i - \bar{x}_r) y_i / SS_{x,r}$ و $\hat{\alpha}_r = \bar{y}_r - \hat{\beta}_r \bar{x}_r$ به دست می‌آیند. همچنین $SS_{x,r} = \sum_{i \in S_r} (x_i - \bar{x}_r)^2$ و $\bar{x}_r = \sum_{i \in S_r} x_i / n_r$ ، $\bar{y}_r = \sum_{i \in S_r} y_i / n_r$

جانه‌ی مانده‌های ساده به صورت زیر تعریف می‌شود؛ مانده‌های مشاهده شده برای e_i^* است. برای $e_j = (y_j - \hat{\alpha}_r - \hat{\beta}_r x_j)$ ، $j \in S_r$

جایگذاری از $(e_j, j \in s_r)$ قرعه کشی شده است. مقدار y جانهی شده به وسیله‌ی $y_i^* = \hat{\alpha}_r + \hat{\beta}_r \bar{x}_i + e_i^*$ مشخص می‌شود.

برآوردهای بر اساس جانهی عبارت‌اند از $\hat{\beta}^* = (\sum_{i \in s_r} x_i y_i + \sum_{i \in s - s_r} x_i y_i^*) / SS_x$

و $\bar{y}_{nr}^* = \sum_{s - s_r} y_i^* / (n - n_r)$ که در آن $\hat{\alpha}^* = (n_r \bar{y}_r + (n - n_r) \bar{y}_{nr}^*) / n$

می‌توان نشان داد $\hat{\sigma}_e^* = \frac{1}{n-2} \left\{ \sum_{s_r} (y_i - \hat{\alpha}^* - \hat{\beta}^* x_i)^2 + \sum_{s - s_r} (y_i^* - \hat{\alpha}^* - \hat{\beta}^* x_i)^2 \right\}$

که (برای خلاصه‌ی اثبات به پیوست ۲ نگاه کنید)

$$E(\hat{\sigma}_e^*) = \sigma^* E \left(\frac{n_r - 2}{n_r} \cdot \frac{n - 2f}{n - 2} \right) \approx \sigma^*$$

شد $f = (n - n_r) / n$ است. چون $\text{var}(\hat{\beta}) = \sigma^* / SS_x$ ، بنا بر این رابطه‌ی (۳) برقرار

است. به سهولت دیده می‌شود که $\text{var}(\hat{\beta}^* | y_{\text{obs}}) = s_e^* c_r / SS_x$ و $E(\hat{\beta}^* | y_{\text{obs}}) = \hat{\beta}_r$. می‌توان نشان داد که

$$E(c_r) = 1 - p_r \quad . \quad E(s_e^* | s_r) = \frac{n_r - 2}{n_r} \sigma^*$$

و $\text{var}(\hat{\beta}_r | s_r) = \sigma^* / SS_{x,r}$.

$$E(k) = \frac{E(1/SS_{x,r}) - 1/SS_x}{(1-p_r)E\{(n_r - 2)/n_r\}/SS_x} \approx \frac{1/E(SS_{x,r}) - 1/SS_x}{(1-p_r)/SS_x}$$

با استفاده از این حقیقت که به شرط s_r, n_r نمونه‌ای تصادفی ساده است به گونه‌ای که نماگرهای پاسخ آن با $\text{cov}(R_i, R_j) = -f(1-f)/(n-1)$ با $E(R_i, R_j) = 0$ همبسته‌اند، در می‌یابیم که

$$E(SS_{x,r}) = \left(p_r - \frac{1-p_r}{n-1} \right) SS_x$$

$$k = 1/(1-f), \quad E(k) = \frac{1}{p_r - \frac{1}{n}} \approx \frac{1}{p_r}.$$

۴- جانه‌ی چندگانه برای نمونه‌های طبقه‌بندی شده

۴-۱- ترکیب‌های جداگانه

یک روش ترکیب m مجموعه‌ی داده‌های کامل این است که برای هر طبقه این کار به طور جداگانه انجام شود یعنی یک مقدار k ای مجزا برای هر طبقه تعیین شود. پس ساختارکلی به صورت زیر است؛ نمونه‌ی s به صورت s_1, \dots, s_H تقسیم می‌شود. فرض کنید y_h کل داده‌های طرح ریزی شده از زیرنمونه‌ی s_h به اندازه‌ی n_h باشد. فرض شده است که y_{H+1}, \dots, y_n مستقل هستند. بخش مشاهده‌شده‌ی y_h با نماد $y_{h,obs}$ نمایش داده می‌شود که s_{hr} نمونه‌ی پاسخ از s_h به اندازه‌ی n_{hr} است. برآوردگر مبتنی بر مجموعه‌ی داده‌های کامل، جمع عبارات مستقل $\hat{\theta} = \sum_{h=1}^H \hat{\theta}_h$ است که در آن $\hat{\theta}_h$ بر اساس y_h به دست آمده است. $\hat{V}(\hat{\theta}) = \sum_{h=1}^H \hat{V}_h(y_h)$ از طریق $\text{var}(\hat{\theta}) = \sum_{h=1}^H \text{var}(\hat{\theta}_h)$ برآورد می‌شود که در آن $\hat{V}_h(y_h)$ برآورد واریانس $\hat{\theta}_h$ بر اساس y_h است. برای $y_i^*, i \in s_h - s_{hr}$ با استفاده از برخی روش‌ها بر اساس $y_{h,obs}$ y جانه‌ی می‌شود و y_h^* را به عنوان داده‌های کامل در نظر می‌گیریم ($y_h^* = \hat{\theta}_h(y_h)$). بر اساس y_h^* ، داریم $\hat{\theta}_h^* = \hat{\theta}_h(y_h^*)$ و $\hat{\theta}^* = \sum_{h=1}^H \hat{\theta}_h^* = \hat{V}_h(y_h^*)$. سپس برآوردگر مبتنی بر جانه‌ی به وسیله‌ی m مشخص می‌شوند. جانه‌ی چندگانه‌ی m جانه‌ی تکرار شده به $\hat{\theta}_{h,i} = \hat{\theta}_h(y_{h,obs})$ و $\hat{V}_{h,i}^* = \sum_{h=1}^H \hat{V}_h^*(y_{h,obs})$ مجموعه‌ی داده‌های کامل شده با m برآورد برای هر طبقه‌ی h ، به صورت $\hat{\theta}_{h,i}$ و برآورد واریانس مربوط $\hat{V}_{h,i}^*$ برای $i = 1, \dots, m$ منجر می‌شود. برآوردهای کل و واریانس‌های مربوط به صورت $\hat{\theta}_i^* = \sum_{h=1}^H \hat{\theta}_{h,i}^*$ برای $i = 1, \dots, m$ است. برآورد ترکیب‌شده برای طبقه‌ی h ام به وسیله‌ی $\bar{\theta}_h^* = \sum_{i=1}^m \hat{\theta}_{h,i}^* / m$ است از $\bar{V}_h^* = \sum_{i=1}^m \hat{V}_{h,i}^* / m$ و مؤلفه‌ی بین جانه‌ی به وسیله‌ی رابطه‌ی $B_h^* = \sum_{i=1}^m (\hat{\theta}_{h,i}^* - \bar{\theta}_h^*)^2 / (m-1)$ مشخص می‌شود. با دنبال کردن همان ایده‌ی بخش ۲،

رابطه‌ی (۱)، واریانس برآورده شده کل $\bar{\theta}_h^*$ ، به صورت $W_h = \bar{V}_h^* + (k_h + \frac{1}{m})B_h^*$ پیشنهاد می‌شود. برآورد کل ترکیب شده نیز به وسیله‌ی $\bar{\theta}_h^* = \sum_{i=1}^m \hat{\theta}_i^* / m = \sum_{h=1}^H \bar{\theta}_h^*$ مشخص می‌شود. نتیجه‌ی می‌گیریم که واریانس برآورده شده کل $\bar{\theta}^*$ نیز می‌تواند به صورت رابطه‌ی زیر بیان شود.

$$(8) \quad W_{\text{sep}} = \sum_{h=1}^H W_h = \bar{V}^* + \sum_{h=1}^H (k_h + \frac{1}{m})B_h^*$$

که در آن $\bar{V}^* = \sum_{i=1}^m \hat{V}_i^* / m = \sum_{h=1}^H \bar{V}_h^*$ برقرار باشد

$$(9) \quad E(\bar{V}_h^*) \approx \text{var}(\hat{\theta}_h)$$

و با استفاده از رابطه‌ی (۵)، باید در رابطه‌ی زیر صدق کند

$$(10) \quad E(k_h) = \frac{\text{var} E(\hat{\theta}_h^* | Y_{h,\text{obs}}) - \text{var}(\hat{\theta}_h)}{E \text{ var}(\hat{\theta}_h^* | Y_{h,\text{obs}})}$$

فرمول ترکیبی (۸) یک شق دیگر برای فرمول ترکیبی معمول (۱) است که به خصوص زمانی سودمند است که به عبارت‌هایی ساده برای k_h و نه برای k دست یافته باشیم. بخش بعدی یک عبارت برای k در این حالت بسط داده است.

۴-۲ - یک فرمول ترکیبی کلی

اکنون فرض کنید W با استفاده از رابطه‌ی (۱) مشخص شده باشد. هدف تعیین فاکتور بین جانهی k است. چون $E(W) = E(W_{\text{sep}})$ است، بنا بر این داریم

$$(11) \quad E \left\{ \sum_{h=1}^H \left(k_h + \frac{1}{m} \right) B_h^* \right\} = E \left(k + \frac{1}{m} \right) B^*$$

$$\text{که در آن، } B^* = \frac{1}{m-1} \sum_{i=1}^m (\hat{\theta}_i^* - \bar{\theta}^*)^* = \frac{1}{m-1} \sum_{i=1}^m \left\{ \sum_h (\hat{\theta}_{h,i}^* - \bar{\theta}_h^*) \right\}^*$$

داشته باشید که چون $E(B^* | y_{obs}) = \text{var}(\hat{\theta}^* | y_{obs}) = \sum_{h=1}^H \text{var}(\hat{\theta}_h^* | y_{obs})$

$E(B_h^* | y_{obs}) = \text{var}(\hat{\theta}_h^* | y_{obs})$ همچنین $E(B^* | y_{obs}) = E(\sum_{h=1}^H B_h^* | y_{obs}).$

بنا بر این اتحاد (۱۱)، به صورت

$$E\left\{\sum_{h=1}^H k_h E(B_h^* | Y_{obs})\right\} = E\{k E(B^* | Y_{obs})\}$$

می‌شود. در صورتی که بخواهیم از فرمول ترکیبی معمول (۱) استفاده کنیم، به حل دست می‌یابیم و لذا

$$(12) \quad k = \frac{\sum_{h=1}^H k_h \text{var}(\hat{\theta}_h^* | y_{obs})}{\text{var}(\hat{\theta}^* | y_{obs})} = \sum_{h=1}^H k_h \cdot \frac{\text{var}(\hat{\theta}_h^* | y_{obs})}{\text{var}(\hat{\theta}^* | y_{obs})}$$

که این رابطه یک میانگین موزون از k_h است. زمانی که همه‌ی k_h ‌ها برابر باشند، مثلاً $k_h = k$ ، به عبارت ساده‌ای برای $k = k$ دست می‌یابیم.

۵- چهار کاربرد برای نمونه‌های طبقه‌بندی شده و بی‌پاسخی تصادفی درون طبقات

۱-۵- براورد میانگین جامعه از نمونه‌ی طبقه‌بندی شده با جانه‌ی بی‌درنگ طبقه‌بندی شده

نمونه‌های تصادفی ساده‌ی طبقه‌بندی شده از یک جامعه‌ی متناهی به اندازه‌ی N ، با H طبقه‌ی بهاندازه‌ی N_h برای $h = 1, \dots, H$ را در نظر بگیرید. هدف براورد میانگین جامعه μ از متغیر y است. فرض کنید بی‌پاسخی کاملاً تصادفی، که توسط روین [۲] و لیتل و روین [۱] با MAR (گمشدگی تصادفی) نشان داده شده است، در هر طبقه وجود

داشته باشد. این مطلب به این معنا است که نماگرهای پاسخ در طبقه‌ی h ، $R_{h,N_h}, \dots, R_{h,1}$ با احتمال‌های $p_{hr} = P(R_{h,i} = 1)$ از هم مستقل هستند. روش جانهی، روش بی‌درنگ طبقه‌بندی شده است. فرض کنید $y_{h,obs} = (y_i : i \in s_{hr})$ باشد، $y_{h,obs}^* = (y_i^* : i \in s_{hr})$ نمونه‌ی پاسخ به اندازه‌ی n_{hr} از طبقه‌ی h باشد، $y_{h,obs}^*$ مقدار جانهی شده‌ی y_i^* در طبقه‌ی h ام به تصادف از قرعه‌کشی می‌شود. برآورده‌گر مبتنی بر داده‌های کل نمونه، میانگین موزون طبقه‌بندی شده‌ی معمول $\bar{Y}_{strat} = \sum_{h=1}^H N_h \bar{y}_h / N = \sum_{h=1}^H v_h \bar{y}_h$ و $v_h = N_h / n_h$ است که در آن s_h نمونه‌ی طبقه‌ی h ام و $n_h = |s_h|$ است. پس $\sigma_h^* = \sqrt{\sum_{i \in U_h} (y_i - \mu_h)^2 / (N_h - 1)}$ با $var(\bar{Y}_{strat}) = \sum_{h=1}^H v_h \sigma_h^* \left(\frac{1}{n_h} - \frac{1}{N_h} \right)$ واریانس جامعه در طبقه‌ی h ام می‌باشد. در اینجا U_h جمعیت طبقه‌ی h ام و μ_h میانگین در U_h است.

فرض کنید \bar{y}_{hr} میانگین نمونه‌ی مشاهده شده از طبقه‌ی h و $\sigma_{hr}^* = \sqrt{\frac{1}{n_{hr}-1} \sum_{i \in s_{hr}} (y_i - \bar{y}_{hr})^2}$ مبتنی بر جانهی به وسیله‌ی $\bar{Y}_{strat}^* = \sum_{h=1}^H N_h \bar{y}_h^* / N$ تعیین می‌شود که در آن $\bar{y}_h^* = \left(\sum_{i \in s_{hr}} y_i + \sum_{i \in s_{hr}} y_i^* \right) / n_h = \left(n_{hr} \bar{y}_{hr} + \sum_{i \in s_{hr}} y_i^* \right) / n_h$ است. فرض کنید m تکرار جانهی $\bar{Y}_{strat,i}^*$ به وسیله‌ی $i = 1, \dots, m$ برای مشخص می‌شود. برآورده‌گر ترکیب شده به وسیله‌ی $\bar{\bar{Y}}_{strat}^* = \sum_{i=1}^m \bar{Y}_{strat,i}^* / m$ مشخص می‌شود.

۱-۱-۵- ترکیبات طبقات مجزا

از بخش ۱-۳ نتیجه شد که $k_h = 1/(1-f_h)$ ، که در آن $f_h = (n_h - n_{hr}) / n_h$ نرخ بی‌پاسخی در طبقه‌ی h ام است. با استفاده از رابطه‌ی (۸)، فرمول ترکیبی برآورد واریانس $\bar{\bar{Y}}_{strat}^*$ به صورت

$$W_{\text{sep}} = \bar{V}^* + \sum_{h=1}^H \left(\frac{1}{1-f_h} + \frac{1}{m} \right) B_h^*$$

است. در اینجا، \bar{V}_h^* و \bar{V}^* میانگین m مقدار از برآورد واریانس مبتنی بر جانبه‌ی $V_h^* = v_h \hat{\sigma}_{h*} (\frac{1}{n_h} - \frac{1}{N_h})$

$$\hat{\sigma}_{h*} = \frac{1}{n_h - 1} \left(\sum_{s_{hr}} (y_i - \bar{y}_h^*)^2 + \sum_{s_h - s_{hr}} (y_i^* - \bar{y}_h^*)^2 \right)$$

۱-۵-۲- فرمول ترکیبی کلی. تعیین k در رابطه‌ی (۱)

از رابطه‌ی (۱۲)، لازم است که $\text{var}(v_h \bar{Y}_h^* | y_{\text{obs}})$ و $\text{var}(\bar{Y}_{\text{strat}}^* | y_{\text{obs}}) = \sum_{h=1}^H \text{var}(v_h \bar{Y}_h^* | y_{\text{obs}})$ تعیین شوند. پس از آن

$$k = \sum_{h=1}^H \frac{1}{1-f_h} \cdot \frac{\text{var}(v_h \bar{Y}_h^* | y_{\text{obs}})}{\text{var}(\bar{Y}_{\text{strat}}^* | y_{\text{obs}})}$$

اکنون،

$$E(\bar{Y}_h^* | y_{h,\text{obs}}) = \bar{y}_{hr}$$

و

$$\text{var}(\bar{Y}_h^* | y_{h,\text{obs}}) = \left\{ \frac{n_h - n_{hr}}{n_h} \right\} \cdot \left\{ \frac{(n_{hr} - 1)}{n_{hr}} \right\} \hat{\sigma}_{hr}^2 \approx \frac{f_h \hat{\sigma}_{hr}^2}{n_h}.$$

بنابراین می‌توانیم مقدار k را به صورت زیر تعیین کنیم.

$$k = \sum_{h=1}^H \frac{1}{1-f_h} \cdot \frac{f_h v_h \hat{\sigma}_{hr}^2 / n_h}{\sum_{k=1}^H f_k v_k \hat{\sigma}_{kr}^2 / n_h}$$

اگر اندازه‌های طبقه‌ی N_h بزرگ باشد، $\hat{V}(v_h \bar{Y}_h) = v_h \hat{\sigma}_{hr}^2 / n_h$ درنظر گرفته می‌شود.
همچنین اگر فرض کنیم $b_h = f_h \hat{V}(v_h \bar{Y}_h) / \sum_{k=1}^H f_k \hat{V}(v_k \bar{Y}_k)$ باشد آن‌گاه

$$(13) \quad k = \frac{\sum_{h=1}^H \hat{V}(v_h \bar{Y}_h) f_h \frac{1}{1-f_h}}{\sum_{h=1}^H \hat{V}(v_h \bar{Y}_h) f_h} = \sum_{h=1}^H b_h \cdot \frac{1}{1-f_h}$$

چون $\sum_{h=1}^H b_h = 1$ است، نتیجه می‌شود که k میانگین موزون عکس نرخ پاسخ است.

اگر همه‌ی $f_h = f$ ، نرخ بی‌پاسخی کلی باشد، آن‌گاه مقدار k مانند نمونه‌ی تصادفی ساده به صورت $k = 1/(1-f)$ درنظر گرفته می‌شود. در غیر این صورت، اگر یا نرخ بی‌پاسخی زیاد باشد و یا واریانس بروآورده‌ی $v_h \bar{Y}_h$ زیاد باشد، نرخ پاسخ طبقه‌ی $1-f_h$ وزن زیادی دارد.

۳-۱-۵- یک عبارت دیگر برای k در رابطه‌ی (۱)

با استفاده از رابطه‌ی (۵) به طور مستقیم، می‌توان فرمول دیگری برای k به دست آورد. به شرط y_{obs} ، میانگین‌های نمونه‌ی جانهی شده‌ی \bar{Y}_h^* مستقل هستند، که این مطلب دلالت بر این دارد که $E(\bar{Y}_{strat}^* | y_{obs}) = \sum_{h=1}^H N_h \bar{y}_{hr} / N = \bar{y}_{strat,r}$ و $E(\bar{Y}_{strat}^* | y_{obs}) \approx \sum_{h=1}^H v_h f_h \hat{\sigma}_{hr}^2 / n_h$. درست مانند بخش ۳-۱، رابطه‌ی (۳) برقرار است و با استفاده از رابطه‌ی (۵)، رابطه‌ی زیر به دست می‌آید.

$$\begin{aligned} E(k) &\approx \frac{\text{var}(\bar{Y}_{strat,r}) - \text{var}(\bar{Y}_{strat})}{E\left(\sum_h v_h f_h \hat{\sigma}_{hr}^2 / n_h\right)} \\ &= \frac{\sum_{h=1}^H v_h \sigma_h^2 \left(E\left(\frac{1}{n_{hr}}\right) - \frac{1}{N_h} \right) - \sum_{h=1}^H v_h \sigma_h^2 \left(\frac{1}{n_h} - \frac{1}{N_h} \right)}{\sum_{h=1}^H v_h E\left\{ \frac{f_h}{n_h} E\left(\hat{\sigma}_{hr}^2 | n_{hr}\right) \right\}} \end{aligned}$$

$$(14) \quad \approx \frac{\sum_{h=1}^H v_h^\top \sigma_h^\top \frac{1-p_{hr}}{n_h} \cdot \frac{1}{p_{hr}}}{\sum_{h=1}^H v_h^\top \sigma_h^\top \frac{1-p_{hr}}{n_h}} = \frac{\sum_{h=1}^H v_h^\top \frac{\sigma_h^\top}{n_{hr}} E(f_h) \frac{1-f_h}{E(1-f_h)}}{\sum_{h=1}^H v_h^\top \frac{\sigma_h^\top}{n_{hr}} E(f_h)(1-f_h)}$$

اکنون، با فرض $\text{var}(\bar{Y}_{hr}) = E \text{ var}(\bar{Y}_{hr} | n_{hr}) = \sigma_h^\top E(1/n_{hr})$. می‌توان نتیجه گرفت که اگر اندازه‌های طبقه‌ی N_h بزرگ باشد، با درنظر گرفتن

$$(15) \quad \frac{1}{k} = \frac{\sum_{h=1}^H (1-f_h) f_h \hat{V}(v_h^\top \bar{Y}_{hr})}{\sum_{h=1}^H f_h \hat{V}(v_h^\top \bar{Y}_{hr})} = \sum_{h=1}^H a_h (1-f_h)$$

که در آن وزن‌ها $a_h = f_h \hat{V}(v_h^\top \bar{Y}_{hr}) / \sum_{k=1}^H f_k \hat{V}(v_k^\top \bar{Y}_{kr})$ است، عبارت مربوط به $E(k)$ به‌طور تقریبی صدق می‌کند. چون $\sum_{h=1}^H a_h = 1$ است، در می‌یابیم که $1/k$ میانگین موزون نرخ‌های پاسخ است. اگر همه‌ی $f_h = f$ ، نرخ بی‌پاسخی کلی باشد، همان‌طور که در بخش ۵-۱-۲ نشان داده شد، $(1-f)/f = k$ خواهد بود. همان‌گونه که در بخش ۵-۱-۲ دیدیم، به این نکته نیز در رابطه‌ی (15) اشاره کردیم که نرخ پاسخ طبقه f_h ، حتی اگر نرخ بی‌پاسخی زیاد باشد و یا واریانس برآورد شده‌ی $v_h^\top \bar{Y}_{hr}$ زیاد باشد، وزن زیادی اخذ می‌کند. برآورد مجموع براساس نمونه‌ی پاسخ به‌وسیله‌ی مشخص $\bar{Y}_{\text{strat},r} = \sum_h v_h^\top \bar{Y}_{hr}$ می‌شود. با درنظر گرفتن این حقیقت از رابطه‌ی (14) که فرمول $E(k) \approx \sum_{h=1}^H \text{var}(v_h^\top \bar{Y}_h) E(f_h) \frac{1}{E(1-f_h)} / \sum_{h=1}^H \text{var}(v_h^\top \bar{Y}_h) E(f_h)$ برای k به دست می‌آید. بنا بر این می‌توان نتیجه گرفت که اگر اندازه‌های طبقه‌ی N_h بزرگ باشد، با قرار دادن مقدار k ای که در رابطه‌ی (13) مشخص شده است، فرمول $E(k)$ به‌طور تقریبی برقرار است.

۳-۵-رگرسیون لوژستیک با متغیرهای تبیینی دودویی. برآورد لگاریتم (نسبت بختها)

متغیرهای Y_1, \dots, Y_n متغیرهای مستقلی هستند که مقادیر صفر و ۱ را اخذ می‌کنند و متغیر تبیینی x با مقادیر ثابت معلوم x_1, \dots, x_n که مقادیر صفر و ۱ را اخذ می‌کند، درنظر گرفته شده است. احتمال ردها با $P(Y_i = 1|x_i = 1) = \pi_i$ و $P(Y_i = 0|x_i = 0) = \pi_{0i}$ مشخص شده است. یک مدل گم‌شدن تصادفی برای متغیرهای پاسخ R_1, \dots, R_n با احتمال‌های $P(R_i = 1|x_i = 1) = p_{1r}$ و $P(R_i = 0|x_i = 0) = p_{0r}$ فرض می‌شود. می‌توان مدل را در قالب مدل لوجیت بهصورت $\log\left\{P(Y=1|x)/P(Y=0|x)\right\} = \alpha + \beta x$ بازپارامتری کرد که در آن برآورد $\beta = \log \frac{\pi_i/(1-\pi_i)}{\pi_{0i}/(1-\pi_{0i})}$ و $\alpha = \log\{\pi_i/(1-\pi_i)\}$ است. فرض کنید $s = (1, \dots, n)$ معرف کل نمونه‌ی با طبقات $S = \{i \in s : x_i = 0\}$ و $S_1 = \{i \in s : x_i = 1\}$ باشد. اندازه‌های n_1 و n_0 با n_s و n_r مشخص شده‌اند که $n_s = n - X$ و $n_r = \sum_{i=1}^n x_i = X$ است. نمونه‌های پاسخ در طبقات، $S_r = \{i \in S : R_i = 1\}$ و $S_{0r} = \{i \in S : R_i = 0\}$ با مجموع نمونه‌ی پاسخ $n_{0r} = |S_{0r}|$ و $n_r = |S_r|$ باشد. می‌بینیم که $n_r = n_r - X_r$ و $n_{0r} = \sum_{s_r} x_i = X_r$. می‌توان داده‌های s_r را که در آن n_{ijr} معرف تعداد مشاهدات با $i = x$ و $j = y$ است، بهصورت زیر نمایش داد (جدول ۱ را ببینید).

پس برآوردهای ماکسیمم درست‌نمایی برای π_{1r} و π_{0r} به ترتیب برابر $\hat{\pi}_{1r} = n_{1r}/n_r$ و $\hat{\pi}_{0r} = n_{0r}/n_r$ است و در نتیجه برآورد ماکسیمم درست‌نمایی β بهصورت $\hat{\beta}_r = \log \frac{\hat{\pi}_{1r}/(1-\hat{\pi}_{1r})}{\hat{\pi}_{0r}/(1-\hat{\pi}_{0r})} = \log(n_{1r}n_{0r}/n_{0r}n_{1r})$

مبتنی بر کل نمونه به وسیله‌ی $\hat{\beta} = \log \frac{\hat{\pi}_1/(1-\hat{\pi}_1)}{\hat{\pi}_0/(1-\hat{\pi}_0)} = \log(n_{1s}n_{0s}/n_{0s}n_{1s})$ با

نمادهایی آشکارا متناظر نمایش داده می‌شود. می‌توان این برآورد را بهصورت

جدول ۱ - داده‌های مشاهده شده و مجموع بی‌پاسخی برای دو رد

بی‌پاسخی	مجموع	$y = 1$	$y = 0$	y / x
$n_{\circ} - n_{\circ r}$	$n_{\circ r}$	$n_{\circ \backslash r}$	$n_{\circ \circ r}$	$x = 0$
$n_{\backslash} - n_{\backslash r}$	$n_{\backslash r}$	$n_{\backslash \backslash r}$	$n_{\backslash \circ r}$	$x = 1$

به همان شکل بخش ۱-۴، $\hat{\beta} = \log\{\hat{\pi}_{\backslash}/(1-\hat{\pi}_{\backslash})\} - \log\{\hat{\pi}_{\circ}/(1-\hat{\pi}_{\circ})\} = \hat{\beta}_{\backslash} - \hat{\beta}_{\circ}$ بیان کرد. همچنین $\hat{\beta}_{\backslash}$ و $\hat{\beta}_{\circ}$ که از طبقاتی با نمونه‌ی مجزای s_{\backslash} و s_{\circ} به دست آمده‌اند، مستقل هستند. برای اندازه‌های نمونه‌ی n_{\backslash} و n_{\circ} بزرگ، $\hat{\beta}$ به طور تقریبی دارای توزیع $N(\beta, \sigma_{\hat{\beta}}^2)$ است که در آن $\sigma_{\hat{\beta}}^2 = \{n_{\backslash}\pi_{\backslash}(1-\pi_{\backslash})\}^{-1} + \{n_{\circ}\pi_{\circ}(1-\pi_{\circ})\}^{-1}$ است. همچنان‌یعنی به طور تقریبی $\text{var}(\hat{\beta}_{\backslash}) = 1/\{n_{\backslash}\pi_{\backslash}(1-\pi_{\backslash})\}$ و $\text{var}(\hat{\beta}_{\circ}) = 1/\{n_{\circ}\pi_{\circ}(1-\pi_{\circ})\}$ می‌شود:

$$\hat{V}(\hat{\beta}) = \frac{1}{n_{\backslash}\hat{\pi}_{\backslash}(1-\hat{\pi}_{\backslash})} + \frac{1}{n_{\circ}\hat{\pi}_{\circ}(1-\hat{\pi}_{\circ})} = \left(\frac{1}{n_{\backslash\backslash}} + \frac{1}{n_{\circ\circ}} \right) + \left(\frac{1}{n_{\backslash\circ}} + \frac{1}{n_{\circ\backslash}} \right)$$

به گونه‌ای که در آن $\hat{V}_{\circ} = \left(\frac{1}{n_{\circ\backslash}} + \frac{1}{n_{\circ\circ}} \right)$ و $\hat{V}_{\backslash} = \left(\frac{1}{n_{\backslash\circ}} + \frac{1}{n_{\backslash\backslash}} \right)$ ، $\hat{V}(\hat{\beta}) = \hat{V}_{\backslash} + \hat{V}_{\circ}$ ترتیب براورد واریانس $\hat{\beta}_{\backslash}$ و $\hat{\beta}_{\circ}$ هستند.

اکنون روش جانه‌ی زیر د رنظر گرفته خواهد شد؛ برای هر مقدار گم شده در $s_{\backslash} - s_{\backslash r}$ ، مقدار جانه‌ی شده‌ی y^* به تصادف از توزیع براورده شده Y به شرط $x = 1$ به صورت زیر استخراج شده است:

$$\text{با احتمال } p_{\circ r} = n_{\circ r} / n_{\circ}, \quad y^* = 1 \quad \text{و با احتمال } p_{\backslash r} = n_{\backslash r} / n_{\backslash}, \quad y^* = 0.$$

همین روش جانه‌ی برای $s_{\circ r} - s_{\circ \circ r}$ استفاده می‌شود به گونه‌ای که y^* به تصادف از توزیع براورده شده Y به شرط $x = 0$ استخراج می‌شود. این روش همان روش جانه‌ی بی‌درنگ طبقه‌بندی شده است به گونه‌ای که مقادیر جانه‌ی شده به تصادف و با جایگذاری از $y_{\circ, \text{obs}} = (y_i : i \in s_{\circ r})$ و $y_{\backslash, \text{obs}} = (y_i : i \in s_{\backslash r})$ استخراج می‌شوند.

مقادیر براورد شده در s_r را می‌توان به همان شکل داده‌های اصلی نشان داد که در آن n_{ij}^* معرف تعداد مقادیر جانهی شده با $x = i$ و $y = j$ است (جدول ۲ را ببینید).
براورد π_i براساس جانهی به وسیله‌ی $\hat{\pi}_i^* = (n_{11r} + n_{11}^*) / n_1$ مشخص می‌شود
به‌گونه‌ای که براورد مبتنی بر جانهی

$$\hat{\beta}_i^* = \log \left\{ \hat{\pi}_i^* / (1 - \hat{\pi}_i^*) \right\} = \log \left\{ (n_{11r} + n_{11}^*) / (n_1 - n_{11r} - n_{11}^*) \right\}$$

است. به همین ترتیب، براوردهای β_i^* و $\beta_{\circ i}^*$ براساس جانهی به صورت $\hat{\beta}_i^* = \hat{\beta}_i^* - \hat{\beta}_{\circ i}^*$ و $\hat{\beta}_{\circ i}^* = \log \left\{ (n_{\circ 1r} + n_{\circ 1}^*) / (n_{\circ} - n_{\circ 1r} - n_{\circ 1}^*) \right\}$ تعیین می‌شوند.
جانهی تکرار شده به m براورد $\hat{\beta}_i^*$ ، $\hat{\beta}_{\circ i}^*$ برای $i = 1, \dots, m$ منجر می‌شود.
براورد ترکیب شده به وسیله‌ی

$$\bar{\beta}^* = \sum_{i=1}^m \hat{\beta}_i^* / m = \sum_{i=1}^m \hat{\beta}_{\circ i}^* / m - \sum_{i=1}^m \hat{\beta}_{\circ, i}^* / m = \bar{\beta}_i^* - \bar{\beta}_{\circ i}^*$$

بیان می‌شود. براورد واریانس جانهی شدهی \hat{V}^* برای $\hat{\beta}^*$ به صورت زیر مشخص می‌شود؛

$$(16) \quad \hat{V}^* = \frac{1}{n_{11r} + n_{11}^*} + \frac{1}{n_{\circ 1r} + n_{\circ 1}^*} + \frac{1}{n_{\circ \circ r} + n_{\circ \circ}^*} + \frac{1}{n_{\circ \circ \circ r} + n_{\circ \circ \circ}^*}$$

می‌بینیم که $E(\hat{V}^* | y_{\text{obs}}) \approx \frac{1}{n_1 \hat{\pi}_{1r} (1 - \hat{\pi}_{1r})} + \frac{1}{n_{\circ} \hat{\pi}_{\circ r} (1 - \hat{\pi}_{\circ r})}$ و رابطه‌ی (۳) برقرار است. همچنین توجه می‌کنیم که رابطه‌ی (۹) نیز به طور مجزا برای هر رده برقرار است.

۱-۲-۵- ترکیب رددهای مجزا

ابتدا از رهیافت بخش ۱-۴ استفاده می‌کنیم و مقادیر مجزای k_1 و k_{\circ} را برای دو رده تعیین می‌کنیم. طبقه‌ی اول $\{i \in s : x_i = 1\}$ را در نظر می‌گیریم. در پیوست ۳ نشان داده شده است که $\text{var}(\hat{\beta}_r^* | y_{1,\text{obs}}) \approx f_1(1 - f_1) \hat{V}(\hat{\beta}_r)$ و $E(\hat{\beta}_r^* | y_{1,\text{obs}}) \approx \hat{\beta}_r$. با استفاده از رابطه‌ی (۱۰)، به طور تقریبی داریم

جدول ۲- مجموعهای جانه‌ی شده برای دو رد

مجموع	$y = 1$	$y = 0$	y / x
$n_0 - n_{0r}$	n_{01}^*	n_{00}^*	$x = 0$
$n_1 - n_{1r}$	n_{11}^*	n_{10}^*	$x = 1$

$$\begin{aligned} E(k_1) &= \frac{\text{var}(\hat{\beta}_{1r}) - \text{var}(\hat{\beta}_1)}{E\{f_1(1-f_1)\hat{V}(\hat{\beta}_{1r})\}} = \frac{E \text{var}(\hat{\beta}_{1r}|n_{1r}) - \text{var}(\hat{\beta}_1)}{E\{f_1(1-f_1)E[\hat{V}(\hat{\beta}_{1r})|n_{1r}]\}} \\ &\approx \frac{\frac{1}{\pi_1(1-\pi_1)} \left(E\left(\frac{1}{n_{1r}}\right) - \frac{1}{n_1} \right)}{\frac{(1-p_{1r})/p_{1r}}{1-p_{1r}}} = \frac{1}{p_{1r}} \end{aligned}$$

که با قرار دادن $k_1 = 1/(1-f_1)$ به طور تقریبی در رابطه صدق می‌کند. درست به همین ترتیب پی می‌بریم که $k_0 = 1/(1-f_0)$ است که در آن $f_0 = (n_0 - n_{0r})/n_0$ نسبت بی‌پاسخی در طبقه‌ی s_0 است. مؤلفه‌ی بین جانه‌ی برای $\hat{\beta}_0^*$ به وسیله‌ی $B_0^* = \frac{1}{m-1} \sum_{i=1}^m (\hat{\beta}_{0,i}^* - \bar{\beta}_0^*)$ جانه‌ی برای $\hat{\beta}_0^*$ است. بنا بر این یک واریانس براوردشده براورد ترکیبی $\bar{\beta}^*$ برای β بر اساس جانه‌ی، با استفاده از رابطه‌ی (۸) به صورت زیر مشخص می‌شود

$$W_{\text{sep}} = \bar{V}^* + \sum_{x=0}^1 \left(\frac{1}{1-f_x} + \frac{1}{m} \right) B_x^*$$

که در آن \bar{V}^* میانگین m تکرار براورد واریانس جانه‌ی شده‌ی \hat{V}^* است که از طریق رابطه‌ی (۱۶) مشخص شده است.

۲-۵- فرمول ترکیبی کلی. تعیین k در رابطه‌ی (۱)

$$\text{var}(\hat{\beta}_o^* | y_{o,\text{obs}}) = f_o(1-f_o)\hat{V}(\hat{\beta}_{or}) \quad \text{و} \quad \text{var}(\hat{\beta}_r^* | y_{r,\text{obs}}) = f_r(1-f_r)\hat{V}(\hat{\beta}_{rr})$$

است، از رابطه‌ی (۱۲) داریم:

$$(۱۷) \quad k = \frac{1}{1-f_r} \cdot \frac{f_r(1-f_r)\hat{V}(\hat{\beta}_{rr})}{\sum_{x=0}^1 f_x(1-f_x)\hat{V}(\hat{\beta}_{xr})} + \frac{1}{1-f_o} \cdot \frac{f_o(1-f_o)\hat{V}(\hat{\beta}_{or})}{\sum_{x=0}^1 f_x(1-f_x)\hat{V}(\hat{\beta}_{xr})}$$

به همین ترتیب، $\text{var}(\hat{\beta}_r) \approx (n_{vr}/n_v)\text{var}(\hat{\beta}_{vr}|n_{vr}) = (1-f_r)\text{var}(\hat{\beta}_{vr}|n_{vr})$ است. بنا بر این می‌توان واریانس برآوردهای کل نمونه‌ی $\hat{\beta}$ و $\hat{\beta}_r$ را به ترتیب از طریق $\hat{V}(\hat{\beta}) = (1-f_r)\hat{V}(\hat{\beta}_{rr})$ و $\hat{V}(\hat{\beta}_r) = (1-f_o)\hat{V}(\hat{\beta}_{or})$ براورد کرد. بنا بر این

$$k = \frac{1}{1-f_r} \cdot \frac{f_r\hat{V}(\hat{\beta}_r)}{\sum_{x=0}^1 f_x\hat{V}(\hat{\beta}_x)} + \frac{1}{1-f_o} \cdot \frac{f_o\hat{V}(\hat{\beta}_o)}{\sum_{x=0}^1 f_x\hat{V}(\hat{\beta}_{xr})} = \frac{1}{1-f_r}b_r + \frac{1}{1-f_o}(1-b_r)$$

درست مانند بخش ۱-۵ می‌بینیم که k میانگین موزون نرخ‌های پاسخ است. اگر همه‌ی $f_h = f$ ، نرخ بی‌پاسخی کلی باشد، آن‌گاه $(1-f)(1-f_r)k = 1/f$. در غیر این صورت، نرخ پاسخ طبقه‌ی f_x زمانی که یا نرخ بی‌پاسخی زیاد باشد و یا واریانس برآوردهای $\hat{\beta}_x$ زیاد باشد، وزن زیادی خواهد داشت. یا این که از رابطه‌ی (۱۷)

$$\frac{1}{k} = \frac{\sum_{x=0}^1 (1-f_x)f_x\hat{V}(\hat{\beta}_{xr})}{\sum_{x=0}^1 f_x\hat{V}(\hat{\beta}_{xr})} = \sum_{x=0}^1 a_x(1-f_x)$$

که در آن وزن‌ها به صورت $a_x = f_x\hat{V}(\hat{\beta}_{xr})/(f_r\hat{V}(\hat{\beta}_{rr})+f_o\hat{V}(\hat{\beta}_{or}))$ می‌باشد. بنا بر این می‌توان $1/k$ را به عنوان میانگین موزون از نرخ‌های پاسخ بیان کرد.

اگر هدف براورد π_1 و π_0 باشد، مقدار $k = 1/(1-f_1)$ و $k = 1/(1-f_0)$ به ترتیب برای π_1 و π_0 به دست می‌آید.

۳-۵- رگرسیون لوزسیتیک با متغیرهای تبیینی رسته‌ای گسسته. براورد لگاریتم (نسبت بخت‌ها)

اگر متغیر تبیینی x به صورت رسته‌ای است که مثلاً H رده را تعریف می‌کند، می‌توان نتایج را به صورت زیر تعمیم داد.

فرض کنید برای $H-1$ این داده‌های $x_h = P(Y=1|x=h)$ ، $h=0, \dots, H-1$ است. رگرسیون لوزسیتیکی که برای این داده‌های تعریف می‌شود، از طریق معرفی $H-1$ متغیر تبیینی دودویی x_{H-1}, \dots, x_1, x_0 صورت می‌گیرد که در آن اگر مشاهدات به رده‌ی h تعلق داشته باشد، برای $x_h = 1$ ، $h=0, \dots, H-1$ و در غیر این صورت صفر خواهد بود. پس اگر $x_0 = x_1 = \dots = x_{H-1} = 0$ باشد، مشاهده به رده‌ی صفر تعلق دارد. مدل لوجیت با $(x_0, x_1, \dots, x_{H-1})$ به صورت

$$\log \left\{ P(Y=1|x) / P(Y=0|x) \right\} = \alpha + \beta_0 x_0 + \beta_1 x_1 + \dots + x_{H-1} \beta_{H-1}$$

نمایش داده می‌شود و می‌توان نتیجه گرفت که برای رده‌ی h در برابر رده‌ی صفر طریق جانه‌ی چندگانه دقیقاً مشابه متغیر دودویی x با جایگزینی رده‌ی h با رده‌ی ۱ انجام می‌شود.

۴-۵- رگرسیون لوزسیتیک با مقادیر گم شده در یک متغیر تبیینی دودویی

این وضعیت همانند بخش ۲-۵ است با این تفاوت که y به طور کامل در s مشاهده شده است، $(y_1, \dots, y_n) = y$ ، و مقادیر گم شده برای متغیر تبیینی x وجود دارد. Y_n, \dots, Y_1 متغیرهای مستقل با مقادیر صفر و ۱ هستند و متغیرهای تبیینی صفر و ۱ با

جدول ۳- داده‌های مشاهده شده و مجموعه‌ای بی‌پاسخی برای طبقه‌ی y

$y = 1$	$y = 0$	y / x
$n_{\circ \mid r}$	$n_{\circ \circ r}$	$x = 0$
$n_{\circ \mid r}$	$n_{\circ \circ r}$	$x = 1$
$n_{\circ r}^{\circ}$	$n_{\circ \circ r}^{\circ}$	مجموع
$n_1^{\circ} - n_{\circ \mid r}^{\circ}$	$n_0^{\circ} - n_{\circ \circ r}^{\circ}$	بی‌پاسخی

مقادیر ثابت x_1, x_2, \dots, x_n که برخی از آن‌ها گم شده است نیز وجود دارد. متغیرهای پاسخ نشان‌دهنده‌ی گم شدگی x_i ‌ها است اما اکنون با مدل گم شدن تصادفی،

$$P(R_i = 1 | y_i = 0) = q_{\circ r} \quad P(R_i = 1 | y_i = 1) = q_{\circ r}$$

در غیر این صورت، مدل مشابه بخش ۵-۲ با احتمال‌های ردیافی $\pi_i = P(Y_i = 1 | x_i = 1)$ و

$$\log\{P(Y = 1 | x) / P(Y = 0 | x)\} = \alpha + \beta x \quad \pi_0 = P(Y_i = 1 | x_i = 0)$$

$$\text{با } \beta = \log \frac{\pi_1 / (1 - \pi_1)}{\pi_0 / (1 - \pi_0)}$$

اکنون فرض کنید $s^{\circ} = \{i \in S : y_i = 0\}$ و $s^1 = \{i \in S : y_i = 1\}$ با

اندازه‌های n° و n^1 باشد. نمونه‌های پاسخ در طبقات، $\{i \in s^1 : R_i = 1\}$ و

$s_r^1 = \{i \in s^1 : R_i = 1\} = s_r^{\circ} \cup s_r^1$ با نمونه‌ی پاسخ کل $s_r = \{i \in S : R_i = 1\}$ است.

حال، داده‌ها می‌توانند همانند قبل معرفی شوند با این تفاوت که مجموعه‌ای بی‌پاسخی به هر طبقه‌ی y تعلق دارد (جدول ۳ را بینید).

برآورد ماکسیمم درستنمایی $\hat{\pi}_{\circ r}, \hat{\pi}_{\circ r}^1$ ، براساس s_r مشابه قبل است که برای

برآورد کل نمونه‌ی $\hat{\beta}$ بوده است. روش جانهی، روش بی‌درنگ طبقه‌بندی شده برای

طبقه‌های y است. برای هر مقدار گم شده‌ی x در $s_r^1 - s_r^{\circ}$ ، مقدار جانهی شده‌ی x^*

به تصادف از $(x_i : i \in s_r^1)$ قرعه‌کشی می‌شود. به همین ترتیب، مقادیر جانهی شده در $s_r^{\circ} - s^{\circ}$ نیز به تصادف از $(x_i : i \in s_r^{\circ})$ قرعه‌کشی می‌شود. مقادیر

جدول ۴- مجموعهای جانه‌ی شده برای طبقه‌های y

$y = 1$	$y = 0$	y / x
$n_{\circ 1}^*$	$n_{\circ \circ}^*$	$x = 0$
n_{11}^*	$n_{\circ 0}^*$	$x = 1$
$n_1^{\circ} - n_{1r}^{\circ}$	$n_{\circ}^{\circ} - n_{\circ r}^{\circ}$	مجموع

جانه‌ی شده در $s_r - s$ می‌تواند به همان صورت داده‌های اصلی بیان شود به‌طوری که اکنون n_{ij}^* معروف تعداد مقادیر جانه‌ی شده با $x = i$ و $y = j$ است (جدول ۴ را بینید). لازم است یک عبارت تقریبی برای امید ریاضی و واریانس $\hat{\beta}_*$ ، که اکنون با نمایش داده می‌شود، به شرط داده‌های مشاهده شده، پیدا کنیم. این مطلب در پیوست ۴ نشان داده شده است که

$$\text{var}(\hat{\beta}_* | y, x_{\text{obs}}) \approx f'(1-f') \left(\frac{1}{n_{1r}} + \frac{1}{n_{\circ 1r}} \right) + f''(1-f'') \left(\frac{1}{n_{\circ r}} + \frac{1}{n_{\circ \circ r}} \right)$$

۹

$$E(\hat{\beta}_* | y, x_{\text{obs}}) \approx \hat{\beta}_r.$$

در این حالت $f'(1-f') = (n_1^{\circ} - n_{1r}^{\circ})/n_1^{\circ}$ نرخ بی‌پاسخی در طبقه‌ی ۱ و $f''(1-f'') = (n_{\circ}^{\circ} - n_{\circ r}^{\circ})/n_{\circ}^{\circ}$ نرخ بی‌پاسخی در طبقه‌ی ۰ است. لازم به ذکر است که $\hat{q}_{\circ r} = n_{\circ r}^{\circ}/n_{\circ}^{\circ}$ و $\hat{q}_{1r} = n_{1r}^{\circ}/n_1^{\circ} = 1-f'$ است. بنابراین مخرج رابطه‌ی (۵) به‌صورت زیر در می‌آید.

$$(18) \quad E \left\{ f'(1-f') \left(\frac{1}{n_{1r}} + \frac{1}{n_{\circ 1r}} \right) + f''(1-f'') \left(\frac{1}{n_{\circ r}} + \frac{1}{n_{\circ \circ r}} \right) \right\}$$

صورت کسر رابطه‌ی (۵) نیز همانند قبل معادل، $\text{var}(\hat{\beta}_r) - \text{var}(\hat{\beta})$ ، است و می‌توان به طور تقریبی به‌صورت زیر بیان کرد

$$(۱۹) \quad \text{var}(\hat{\beta}_r) - \text{var}(\hat{\beta}) = \frac{1}{n_1\pi_1(1-\pi_1)} \cdot \frac{1-p_{1r}}{p_{1r}} + \frac{1}{n_0\pi_0(1-\pi_0)} \cdot \frac{1-p_{0r}}{p_{0r}}$$

$p_{0r} = P(R_i = 1 | x_i = 0)$ و $p_{1r} = P(R_i = 1 | x_i = 1)$ می‌باشد. به برآوردهای دیگری از p_{0r} و p_{1r} نیز نیاز داریم. چون $\hat{p}_{1r} = \hat{\pi}_1(1-f^1) + (1-\hat{\pi}_1)(1-f^0)$, $p_{1r} = \pi_1 q_{1r} + (1-\pi_1)q_{0r}$ همین ترتیب، $\hat{p}_{0r} = \hat{\pi}_0(1-f^0) + (1-\hat{\pi}_0)(1-f^1)$ همچنین می‌توانیم از برآوردهای دیگری چون $\hat{p}_{1r} \approx n_{1r}$ و $n_0\hat{p}_{0r} \approx n_{0r}$ نیز استفاده کنیم. با استفاده از روابط (۱۸) و (۱۹) می‌توان نتیجه گرفت که رابطه‌ی زیر برقرار است.

$$\begin{aligned} k &= \frac{\left(\frac{1}{n_{1r}} + \frac{1}{n_{0r}}\right)\left(\hat{\pi}_{1r}f^1 + (1-\hat{\pi}_{1r})f^0\right) + \left(\frac{1}{n_{0r}} + \frac{1}{n_{1r}}\right)\left(\hat{\pi}_{0r}f^0 + (1-\hat{\pi}_{0r})f^1\right)}{f^1(1-f^1)\left(\frac{1}{n_{1r}} + \frac{1}{n_{0r}}\right) + f^0(1-f^0)\left(\frac{1}{n_{0r}} + \frac{1}{n_{1r}}\right)} \\ &= \frac{f^1\left(\frac{1}{n_{1r}} + \frac{1}{n_{0r}}\right) + f^0\left(\frac{1}{n_{1r}} + \frac{1}{n_{0r}}\right)}{f^1(1-f^1)\left(\frac{1}{n_{1r}} + \frac{1}{n_{0r}}\right) + f^0(1-f^0)\left(\frac{1}{n_{0r}} + \frac{1}{n_{1r}}\right)} \end{aligned}$$

لازم به ذکر است که اگر $f^1 = f^0 = f$ باشد، آن‌گاه $k = 1/(1-f)$ است. در غیر این صورت، می‌توان $k/1$ را به عنوان یک ترکیب خطی از نرخ‌های پاسخ $w_0 = \frac{1}{n_{0r}} + \frac{1}{n_{1r}}$ و $w_1 = \frac{1}{n_{1r}} + \frac{1}{n_{0r}}$ بیان کرد. فرض کنید $(1-f^1, 1-f^0)$ باشد. بنا بر این خواهیم داشت

$$\frac{1}{k} = a_1(1-f^1) + a_0(1-f^0)$$

که در آن $a_{\circ} = f^{\circ}w_{\circ} / (f^{\circ}w_{\circ} + f^{\circ}w_1)$ و $a_i = f^{\circ}w_i / (f^{\circ}w_{\circ} + f^{\circ}w_i)$ است. لازم به ذکر است که در حالت کلی $a_i + a_{\circ} \neq 1$.

۶- سؤال: آیا می‌توان از فرمول ترکیبی یکسان برای یک وضعیت مشخص و یک روش جانه‌ی برای همه‌ی براوردهای علمی استفاده کرد؟

در این بخش سعی می‌شود یک رهیافت کلی برای این مسئله ارائه شود. فرض کنید s معرف کل نمونه و y داده‌های کل نمونه باشد. سه حالت ممکن است وجود داشته باشد:

- ۱- s یک نمونه از یک جامعه‌ی متناهی و $(y_i : i \in s) = y$ با مدل طرح باشد. بنا بر این متغیرهای تصادفی مشاهده شده، (s, s_r) هستند و y_{obs} معادل (s, s_r) است.

- ۲- وضعیت مشابهی همانند حالت ۱، ولی با یک مدل جامعه به جای مدل طرح وجود داشته باشد. در این حالت، متغیرهای تصادفی مشاهده شده

$$y_{\text{obs}} = \{(y_i : i \in s_r), s_r, s\}$$

- ۳- یک مطالعه‌ی مشاهده‌ای که در آن $(1, \dots, n) = s$ باشد و $y = (y_1, \dots, y_n)$ باشد. در این حالت متغیرهای تصادفی مشاهده شده به صورت

$$y_{\text{obs}} = \{(y_i : i \in s_r), s_r\}$$

به عنوان یک مثال، حالتی را با وجود بی‌پاسخی با گم‌شدن گی کاملاً تصادفی (که متغیرهای پاسخ R_i با احتمال $p_r = P(R_i = 1)$ مستقل هستند) و جانه‌ی بی‌درنگ در نظر می‌گیریم. حالت معرفی شده در بخش ۳-۱ با نمونه‌ی تصادفی ساده، یک حالت خاص شماره‌ی ۱ است و دریافتیم که برای براورد میانگین جامعه با استفاده از میانگین نمونه،

$$(20) \quad k = \frac{1}{1-f}$$

که در آن $\hat{p}_r = (n - n_r) / n = 1 - f$ نرخ بی‌پاسخی است.

با محدود کردن توجه به براوردهای خطی که در آن براوردگر جانه‌ی شده $\hat{\theta}^*$ ، پارامتر مشابهی مانند $\hat{\theta}$ را براورد می‌کند، نشان خواهیم داد که رابطه‌ی (۲۰) در حالت کلی برای

هر سه حالت معرفی شدهی بالا، زمانی که مکانیسم بی‌پاسخی، گمشدگی کاملاً تصادفی است و از روش جانهی بی‌درنگ استفاده شده باشد، برقرار است. اولین سؤالی که بررسی می‌کنیم این است که آیا روش جانهی بی‌درنگ به برآوردهای معتبر جانهی- مبنا منجر می‌شود به‌گونه‌ای که این مقدار k بتواند استفاده شود. جواب در حالت کلی خیر است. نیازی آشکار برای یک روش جانهی، حداقل به طور تقریبی، آن است که

$$(21) \quad E(\hat{\theta}^* | y, s) = \hat{\theta}$$

یعنی برآوردهای جانهی شده باید پارامتر مشابهی مانند $\hat{\theta}$ را برآورد کند. می‌توان گفت که مقدار مورد انتظار برآوردهای جانهی شده به شرط کل داده‌های نمونه‌گیری شده، باید معادل برآورد کل نمونه باشد. در حالت ۱، زمانی که s مشخص باشد، y اضافی است و رابطه‌ی (21) بیان می‌کند که $E(\hat{\theta}^* | s) = \hat{\theta}$ است. در حالت ۳، s تصادفی نیست و بنا بر این غیر ضروری است، در حالی که در حالت ۲ هر دوی s و y مورد نیاز است. در این مقاله برآوردهایی را درنظر گرفته‌ایم که خطی در $(y_i : i \in s)$ هستند. نتایج زیر که در پیوست ۵ اثبات شده‌اند، برآوردهای خطی را مشخص می‌کنند که با استفاده از روش جانهی بی‌درنگ در رابطه‌ی (21) صدق می‌کند و نشان می‌دهند که برای چنین برآوردهایی، $k = 1/(1-f)$ است.

لم. فرض کنید $E(\hat{\theta}^* | y, s) = \hat{\theta} = \sum_{i \in s} a_i(s) y_i$ اگر و تنها اگر برای همه‌ی i ، $a_i(s) = a(s)$. یعنی $a_i(s) = a(s)$ ، $i \in s$. تذکر. در حالت ۳، s هیچ اطلاعی نمی‌دهد و $a_i(s) = a_i(s)$.

قضیه. فرض کنید $E(\hat{\theta}^* | y, s) = \hat{\theta} = \sum_{i \in s} a_i(s) y_i$ باشد. همچنین فرض کنید که رابطه‌ی (3) برقرار است. بنا بر این $E(k) = \frac{E(1/\hat{p}_r) - 1}{1 - p_r - \frac{1}{n}[E(1/\hat{p}_r) - 1]} \approx \frac{1}{p_r}$ و رابطه‌ی $k = 1/(1-f)$ می‌تواند مورد استفاده قرار گیرد.

اکنون به برخی از حالت‌های خاص می‌پردازیم؛

۱- با درنظر گرفتن $a(s) = 1/n$ ، مشابه بخش ۱-۳، می‌بینیم که رابطه‌ی (۲۱)

برقرار است.

۲- ضریب رگرسیون برای رگرسیونی که از مبدأ می‌گذرد، به صورت $\hat{\beta} = \sum_{i=1}^n y_i / \sum_{i=1}^n x_i$ است. در این حالت، با درنظر گرفتن $a = 1 / \sum_{i=1}^n x_i$ رابطه‌ی (۲۱) برقرار است و بنا بر این $k = 1/(1-f)$ است.

۳- براورد ضریب رگرسیون در رگرسیون خطی معمول، حالتی است که در رابطه‌ی (۲۱) صدق نمی‌کند و در آن $\hat{\beta} = \sum (x_i - \bar{x}) y_i / \sum (x_i - \bar{x})^2$ است. در این حالت، $a_i = (x_i - \bar{x}) / \sum_{j=1}^n (x_j - \bar{x})^2$ است که مستقل از i نیست. می‌توان نشان داد که به طور تقریبی $E(\hat{\beta}^* | y) \approx p_r \hat{\beta}$ و به طور دقیق به صورت $\frac{np_r - 1}{n - 1} \hat{\beta}$ است. بنا بر این برای مسائل رگرسیون عادی، روش جانه‌ی بی‌درنگ نمی‌تواند مؤثر واقع شود. لازم به ذکر است که با استفاده از مطالب بخش ۳-۳ می‌توان با استفاده از جانه‌ی مانده‌ها از $k = 1/(1-f)$ استفاده کرد.

بدیهی است که زمانی که y با مقدار معلوم x در یک گروه بی‌پاسخی همبسته باشد، باید صرف نظر از مشکلات براورد تحت بررسی، این همبستگی را در جانه‌ها مورد استفاده قرار داد.

۷- بحث و نتیجه‌گیری

در این مقاله نشان داده شد که امکان بسط یک قضیه‌ی کلی برای جانه‌ی چندگانه، بدون نیاز به قرعه‌کشی تصادفی جانه‌ها از یک توزیع پسین بیزی وجود دارد. برای نمونه‌های طبقه‌بندی شده با براوردهای طبقه‌بندی شده، نیاز به مطالعات بیشتر برای این که چه

برآورد واریانسی به کار گرفته شود، وجود دارد. یعنی آیا رابطه‌ی (۸) با استفاده از ترکیب طبقات مجزا که در رابطه‌ی (۱۰) مشخص شده است، استفاده شود و یا ترکیب کلی رابطه‌ی (۱) با مقدار k مشخص شده در رابطه‌ی (۱۲) به کار گرفته شود.

از مواردی که در این مقاله ارائه شده است می‌توان نتیجه گرفت که فرمول روش جانهی غیر بیزی نوعاً به معیاری از سهم اطلاع گم شده در نمونه‌ی پاسخ در مقایسه با کل نمونه بستگی دارد. در ساده‌ترین حالت در بخش ۳-۱، اطلاعات گم شده به‌وسیله‌ی $(1-f)$ که عکس ترخ پاسخ است، اندازه‌گیری شده است. هر چه این عامل بیش‌تر باشد، وزن مؤلفه‌ی بین جانهی نیز بیش‌تر است. در مدل نسبت بخش ۳-۲ با روش جانهی رگرسیونی بی‌درنگ مانده‌ها، میزان اطلاعات گم شده عکس نسبت x - کل در نمونه‌ی پاسخ در مقایسه با کل نمونه است. نشان داده شده است که در رگرسیون خطی ساده با عبارت واریانسی که مستقل از متغیر تبیینی است، مجدداً اطلاع گم شده به‌وسیله‌ی $(1-f)$ اندازه‌گیری می‌شود. یک پیشنهاد برای مطالعات بعدی، بررسی امکان تعمیم این نتیجه برای تعیین k ، از طریق تعریف معیارهای مرتبط با اطلاعات گم شده، با استفاده از فرمول تعریف پایه‌ای (۵) می‌باشد.

مطالعه‌ی عملکرد بازه‌های اطمینان مربوط نیز باقی می‌ماند. برخی مطالعات شبیه‌سازی مقدماتی که در این مقاله گنجانده نشده‌اند، نشان می‌دهند که برای رگرسیون خطی ساده با استفاده از روش جانهی مانده‌ها و $(1-f)$ ، بازه‌ی اطمینان به صورت $\bar{\beta} \pm z_{\alpha/2} \sqrt{W}$ (که در آن $z_{\alpha/2}$ نقطه‌ی بالایی توزیع نرمال استاندارد می‌باشد) است که به طور تقریبی به سطح اسمی $(1-\alpha)$ دست می‌یابد.

مرجع‌ها

- [1] Little, R.J.A.; Rubin, D.B. (2002). *Statistical Analysis with Missing Data*. Wiley, New York.
- [2] Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. Wiley, New York.
- [3] Rubin, D.B. (1996). Multiple imputation after 18 + years (with discussion). *Journal of the American Statistical Association*, **91**, 473–489.

- [4] Schafer, J.L. (1997). *Analysis of Incomplete Multivariate Data*. Chapman and Hall, London.

پیوست

۱- جانه‌ی چندگانه برای مدل نسبت در بخش ۳-۲

ابتدا شرط (۳) که معادل $E(\hat{\sigma}_*^*) \approx \sigma^*$ می‌باشد را درنظر بگیرید. فرض کنید $\hat{\sigma}_{nr}^* = \sum_{s=s_r}^1 (y_i^* - \hat{\beta}_{nr} x_i)^*/(n_{nr}-1)$ و $\hat{\beta}_{nr} = \sum_{s=s_r} y_i^*/X_{nr}$ در این حالت، $n_{nr} = n - n_r$. سپس عبارت $\hat{\sigma}_*^*$ می‌تواند به صورت زیر بیان شود:

$$\hat{\sigma}_*^* = \frac{1}{n-1} \left((n_r-1)\hat{\sigma}_r^* + (n_{nr}-1)\hat{\sigma}_{nr}^* + \frac{X_r X_{nr}}{X} (\hat{\beta}_r - \hat{\beta}_{nr}) \right)$$

در این حالت، $E(Y_i^* | y_{\text{obs}}) = \hat{\beta}_r x_i + \bar{e} \sqrt{x_i}$ و $\bar{e} = \sum_{s_r} e_i / n_r$ که در آن $\text{var}(Y_i^* | y_{\text{obs}}) = x_i s_e^*$ است که در آن $s_e^* = \frac{1}{n_r} \sum_{s_r} (e_i - \bar{e})^2$ است. با استفاده از این اطلاعات، می‌توان نشان داد که

$$E(\hat{\sigma}_*^*) = \sigma^* \left(1 - \frac{c_1}{n-1} - \frac{4c_r}{(n-1)n_r} - c_n f \frac{n-1}{n \cdot n_r} \right)$$

که در آن c_1, c_r, c_n در فاصله‌ی (۱۰) قرار می‌گیرند. بنا بر این، $\hat{\sigma}_*^* \approx \sigma^*$ و رابطه‌ی (۳) برای مقادیر زیاد و متوسط n_r برقرار است.

در مرحله‌ی بعد، به $E(\hat{\beta}^*|y_{\text{obs}})$ و $\text{var}(\hat{\beta}^*|y_{\text{obs}})$ می‌بینیم که

$$E(\hat{\beta}_{nr}|y_{\text{obs}}) = \hat{\beta}_r + (\bar{e}/X_{nr}) \sum_{s=s_r} \sqrt{x_i}$$

$$\hat{\beta}^* = (\hat{\beta}_r X_r + \hat{\beta}_{nr} X_{nr})/X$$

می‌باشد. این اطلاعات بیان می‌کند که $\text{var}(\hat{\beta}_{nr}|y_{\text{obs}}) = s_e^r/X_{nr}$

$$\text{var}(\hat{\beta}^*|y_{\text{obs}}) = (X_{nr}/X^r) s_e^r$$

$$E(\hat{\beta}^*|y_{\text{obs}}) = \hat{\beta}_r + (\bar{e}/X^r) \sum_{s=s_r} \sqrt{x_i}$$

است. و این نتیجه می‌شود که

$$\text{var} E(\hat{\beta}^*|y_{\text{obs}}) = \text{var}(\hat{\beta}_r) + \frac{\left(\sum_{s=s_r} \sqrt{x_i} \right)^r}{X^r} \text{var}(\bar{e}) + 2 \frac{\sum_{s=s_r} \sqrt{x_i}}{X} \text{cov}(\hat{\beta}_r, \bar{e})$$

اکنون، $\text{cov}(\hat{\beta}_r, \bar{e}) = 0$ است. با استفاده از نابرابری کوشی شوارتس، $b_i = 1$ و $a_i = \sqrt{x_i}$ با $(\sum a_i b_i)^r \leq \sum a_i^r \sum b_i^r$ که این می‌شود که

$$\left(\sum_{s=s_r} \sqrt{x_i} \right)^r \leq nX$$

$$\text{var}(\bar{e}) = \left(\sigma^r/n_r \right) \left(1 - \left(\sum_{s=s_r} \sqrt{x_i} \right)^r / n_r X_r \right) = (1-d_r) \sigma^r/n_r, \quad 0 \leq d_r \leq 1$$

$$\left(\sum_{s=s_r} \sqrt{x_i} \right)^r / X^r = d_r n_{nr} X_{nr} / X^r, \quad 0 \leq d_r \leq 1$$

با بر این،

$$\text{var} E(\hat{\beta}^*|y_{\text{obs}}) = \frac{\sigma^r}{X_r} + \frac{(1-d_r)d_r n_{nr} X_{nr}}{X^r} \cdot \frac{\sigma^r}{n_r}$$

سپس درمی‌یابیم که از $E(s_e^r) = \sigma^r \left(1 - \frac{1}{n_r} \right) - \text{var}(\bar{e}) = \sigma^r (n_r + d_r - 2)/n_r$ این اطلاع خواهیم داشت

$$E \operatorname{var}(\hat{\beta}^* | y_{\text{obs}}) = \frac{X_{nr}}{X} \cdot \frac{\sigma^2}{n_r} (n_r + d - 2)$$

۲- جانه‌ی چندگانه در رگرسیون خطی ساده در بخش ۳-۳. خلاصه‌ای از اثبات:

$E(\hat{\sigma}_*^2) = \sigma^2 E\left(\frac{n_r - 2}{n_r} \cdot \frac{n - f}{n - 2}\right), f = (n - n_r)/n$ در اینجا $\hat{\sigma}_*^2 = \frac{1}{n-2}(SS_e^r + SS_e^{nr})$ و $SS_e^r = \sum_{s_r} (y_i - \hat{\alpha}^* - \hat{\beta}^* x_i)^2$ در آن $\hat{\sigma}_*^2 = \frac{1}{n-2}(SS_e^r + SS_e^{nr})$ است. $SS_e^{nr} = \sum_{s-s_r} (y_i^* - \hat{\alpha}^* - \hat{\beta}^* x_i)^2$ می‌توان دو مجموع مربعات مانده‌ها را به صورت زیر بیان کرد

$$SS_e^r = \sum_{s_r} (y_i - \hat{\alpha}_r - \hat{\beta}_r x_i)^2 + \sum_{s_r} \left[(\hat{\alpha}_r - \hat{\alpha}^*)^2 + (\hat{\beta}_r - \hat{\beta}^*) x_i \right]$$

$$SS_e^{nr} = \sum_{s-s_r} (y_i - \hat{\alpha}_{nr}^* - \hat{\beta}_{nr}^* x_i)^2 + \sum_{s-s_r} \left[(\hat{\alpha}_{nr}^* - \hat{\alpha}^*)^2 + (\hat{\beta}_{nr}^* - \hat{\beta}^*) x_i \right]$$

که در آن $\hat{\beta}_{nr}^*, \hat{\alpha}_{nr}^*$ برآوردهایی هستند که تنها بر اساس y_i^* جانه‌ی شده برای $i \in s - s_r$ به وجود آمده‌اند. از این مطلب می‌توان نتیجه گرفت که

$$E(SS_e^r | y_{\text{obs}}) = n_r s_e^r + n_r \operatorname{var}(\hat{\alpha}^* | y_{\text{obs}}) + \left(\sum_{s-s_r} x_i^r \right) \operatorname{var}(\hat{\beta}^* | y_{\text{obs}}) + 2n_r \bar{x}_r \operatorname{cov}(\hat{\alpha}^*, \hat{\beta}^* | y_{\text{obs}})$$

$$= n_r s_e^r + f(1-f) s_e^r + (1-c_r) SS_x c_r s_e^r / SS_x + \gamma n_r \bar{x}_{nr} f \bar{x}_{nr} s_e^r / SS_x$$

$$\left\{ \bar{x}_{nr} = \sum_{s=s_r} x_i / (n - n_r) \right\} \text{ با}$$

$$(22) \Rightarrow E(SS_e^r | y_{obs}) = s_e^r \{ n_r + f(1-f) + c_r (1-c_r) + \gamma n_r \bar{x}_r \bar{x}_{nr} f (1/SS_x) \}$$

پس از بعضی محاسبات جبری، در می‌یابیم که

$$E(SS_e^r | y_{obs}) = \sum_{s=s_r} \text{var}(Y_i^* - \bar{Y}_{nr}^* | y_{obs}) + \sum_{s=s_r} (x_i - \bar{x}_{nr})^* \text{var}(\hat{\beta}_{nr}^* | y_{obs})$$

$$- \gamma \sum_{s=s_r} (x_i - \bar{x}_{nr}) \text{cov}(Y_i^* - \bar{Y}_{nr}^*, \hat{\beta}_{nr}^* | y_{obs})$$

$$+ (n - n_r) \text{var}(\hat{\alpha}_{nr}^* - \hat{\alpha}^* | y_{obs}) + c_r SS_x \text{var}(\hat{\beta}_{nr}^* - \hat{\beta}^* | y_{obs})$$

$$+ \gamma (n - n_r) \bar{x}_{nr} \text{cov}(\hat{\alpha}_{nr}^* - \hat{\alpha}^*, \hat{\beta}_{nr}^* - \hat{\beta}^* | y_{obs})$$

$$= (n - n_r - 1) s_e^r + s_e^r - \gamma s_e^r + (n - n_r) \left\{ (1-f) s_e^r \frac{1}{n - n_r} + \bar{x}_{nr}^* s_e^r \frac{1}{SS_{x,nr}} \right\}$$

$$+ c_r SS_x s_e^r \left(\frac{1}{SS_{x,nr}} + \frac{c_r}{SS_x} - \gamma \frac{1}{SS_x} \right)$$

$$+ \gamma (n - n_r) \bar{x}_{nr} \left(f \bar{x}_{nr} s_e^r \frac{1}{SS_x} - \bar{x}_{nr} s_e^r \frac{1}{SS_{x,nr}} \right)$$

که در آن $SS_{x,nr} = \sum_{s=s_r} (x_i - \bar{x}_{nr})^*$ و $SS_{x,nr} = \sum_{s=s_r} x_i^r - (n - n_r) \bar{x}_{nr}^r = c_r SS_x - (n - n_r) \bar{x}_{nr}^r$ می‌باشد. از مطالب بالا نتیجه می‌شود که

$$E(SS_e^r | y_{obs}) = s_e^r ((n - n_r - 1) + (1-f)^r + c_r^r - \gamma c_r$$

$$(23) \quad + \gamma f \bar{x}_{nr}^r (n - n_r) / SS_x)$$

از روابط (22) و (23) در می‌یابیم که

$$\begin{aligned}
 (n - r)E(\hat{\sigma}_*^* | y_{\text{obs}}) &= s_e^* \left(n - f - c_r + f \frac{1}{SS_x} \bar{x}_{nr} n \bar{x} \right) = s_e^* (n - f - c_r) \\
 E(c_r | n_r) &= \frac{1}{SS_x} \sum_{i=1}^n E(-R_i | n_r) x_i^* = (-n_r/n) = f \quad \text{چون} \\
 (n - r)E(\hat{\sigma}_*^*) Es_e^*(n - f - c_r) &= E(n - f - c_r) E(s_e^* | s_r) \\
 &= E(n - f - c_r) \frac{n_r - r}{n_r} \sigma^* \\
 &= \sigma^* E \left\{ \frac{n_r - r}{n_r} (n - f - E(c_r | n_r)) \right\} \\
 &= \sigma^* E \left\{ \frac{n_r - r}{n_r} (n - rf) \right\}
 \end{aligned}$$

۳- رگرسیون لوژستیک با متغیرهای تبیینی دودویی. ترکیب رددهای مجزا

هدف تعیین $\text{var}(\hat{\beta}_1^* | y_{1,\text{obs}})$ و $E(\hat{\beta}_1^* | y_{1,\text{obs}})$ است. به شرط $n_{11}, y_{1,\text{obs}}$ دو جمله‌ای با $(n_{11} - n_{1r}, \hat{\pi}_{1r})$ است. بنابراین، $\text{var}(n_{11}^* | y_{1,\text{obs}}) = (n_{11} - n_{1r}) \hat{\pi}_{1r} (1 - \hat{\pi}_{1r})$ و $E(n_{11}^* | y_{1,\text{obs}}) = (n_{11} - n_{1r}) \hat{\pi}_{1r}$. به شرط $T = \log \{(a + Z)/(b - Z)\}$ است که در آن Z به صورت $y_{1,\text{obs}}$ دارای توزیع دوجمله‌ای با پارامترهای (n, p) و a و b مقادیر ثابت هستند. براساس خطی‌سازی تیلور در نزدیکی $E(Z) = np$ به رابطه‌ی $T \approx \log \{(a + np)/(b - np)\} + (Z - np)(a + b)/\{(a + z)(b - z)\}$ دست می‌یابیم و خواهیم داشت.

$$E(T) \approx \log \frac{a+np}{b-np} \quad (24)$$

$$\text{var}(T) \approx \left(\frac{a+b}{(a+np)(b-np)} \right)^2 np(1-p)$$

با درنظر گرفتن $b = n_r - n_{\text{ir}}$ و $a = n_{\text{ir}} - n_{\text{or}}$ به نتیجه‌های دست $E(\hat{\beta}_r^* | y_{\text{ir,obs}}) \approx \hat{\beta}_{\text{ir}}$ و $\text{var}(\hat{\beta}_r^* | y_{\text{ir,obs}}) \approx \left(n_r / \{n_r \hat{\pi}_{\text{ir}} n_{\text{ir}} (1 - \hat{\pi}_{\text{ir}})\} \right)^2 (n_r - n_{\text{ir}}) \hat{\pi}_{\text{ir}} (1 - \hat{\pi}_{\text{ir}})$ می‌یابیم. فرض کنید $f_r = (n_r - n_{\text{ir}}) / n_r$ نرخ بی‌پاسخی در طبقه‌ی s_1 باشد. بنا بر این

$$\begin{aligned} \text{var}(\hat{\beta}_r^* | y_{\text{ir,obs}}) &\approx \frac{f_r n_r}{n_r^2} \cdot \frac{1}{\hat{\pi}_{\text{ir}} (1 - \hat{\pi}_{\text{ir}})} = f_r (1 - f_r) \cdot \frac{1}{n_{\text{ir}} \hat{\pi}_{\text{ir}} (1 - \hat{\pi}_{\text{ir}})} \\ &= f_r (1 - f_r) \hat{V}(\hat{\beta}_{\text{ir}}) \end{aligned}$$

۴- رگرسیون لوزیستیک با مقادیر گم شده در یک متغیر تبیینی دودویی

به منظور تعیین $E(\hat{\beta}_* | y, x_{\text{obs}})$ و $\text{var}(\hat{\beta}_* | y, x_{\text{obs}})$ ، نیاز است که $\hat{\beta}_*$ به روشی متفاوت با روش بخش ۲-۵ بیان شود تا بتواند مجموع دو عبارت مستقل به شرط داده‌های مشاهده شده‌ی (y, x_{obs}) باشد:

$$\hat{\beta}_* = \log \frac{(n_{\text{ir}} + n_{\text{ir}}^*) (n_{\text{or}} + n_{\text{or}}^*)}{(n_{\text{or}} + n_{\text{or}}^*) (n_{\text{ir}} + n_{\text{ir}}^*)} = \log \frac{(n_{\text{ir}} + n_{\text{ir}}^*)}{(n_{\text{ir}} + n_{\text{ir}}^*)} - \log \frac{(n_{\text{or}} + n_{\text{or}}^*)}{(n_{\text{or}} + n_{\text{or}}^*)} = \hat{\beta}_r^* - \hat{\beta}_o^*$$

به شرط $\text{var}(\hat{\beta}_* | y, x_{\text{obs}}) = \text{var}(\hat{\beta}_r^* | y, x_{\text{obs}}) + \text{var}(\hat{\beta}_o^* | y, x_{\text{obs}})$ و n_{ir}^* دارای توزیع دوجمله‌ای با $(n_{\text{ir}}^* - n_{\text{ir}}, p^*)$ است که در آن $p^* = n_{\text{ir}}^* / n_{\text{ir}}$ است و n_{or}^* دارای توزیع دوجمله‌ای با $(n_{\text{or}}^* - n_{\text{or}}, p^*)$ است که در آن $p^* = n_{\text{or}}^* / n_{\text{or}}$ می‌باشد. پس با استفاده از رابطه‌ی (۲۴)، می‌توان به این نتیجه رسید که $E(\hat{\beta}_* | y, x_{\text{obs}}) = \log \{p^* / (1 - p^*)\}$

است. $\text{var}(\hat{\beta}_* | y, x_{\text{obs}}) = \left(n_r / \{n_r p^r n_r (1-p^r)\} \right)^r (n_r - n_r^r) p^r (1-p^r)$ بنابراین $\{ \text{var}(\hat{\beta}_* | y, x_{\text{obs}}) \approx f^r / \{n_r p^r (1-p^r)\} = f^r (1-f^r) / \{n_r p^r (1-p^r)\}$ به همین ترتیب $E(\hat{\beta}_* | y, x_{\text{obs}}) \approx \log \{p^r / (1-p^r)\}$. همچنین $E(\hat{\beta}_* | y, x_{\text{obs}}) \approx \hat{\beta}_r$

$$\text{var}(\hat{\beta}_*^r | y, x_{\text{obs}}) \approx f^r / \{n_r p^r (1-p^r)\} = f^r (1-f^r) / \{n_r p^r (1-p^r)\}$$

$$\frac{1}{n_r p^r (1-p^r)} = \frac{1}{n_{\text{or}}} + \frac{1}{n_{\text{oor}}} \quad \text{و} \quad \frac{1}{n_r^r p^r (1-p^r)} = \frac{n_r^r}{n_{\text{or}} n_{\text{oor}}} = \frac{1}{n_{\text{or}}} + \frac{1}{n_{\text{oor}}}$$

می‌دانیم و از این‌ها نتیجه می‌شود که

$$\text{var}(\hat{\beta}_* | y, x_{\text{obs}}) \approx f^r (1-f^r) \left(\frac{1}{n_{\text{or}}} + \frac{1}{n_{\text{oor}}} \right) + f^r (1-f^r) \left(\frac{1}{n_{\text{or}}} + \frac{1}{n_{\text{oor}}} \right).$$

۵- اثبات لم و قضیه‌ی مربوط به بخش ۶

به منظور اثبات لم و قضیه در بخش ۶ به دانستن بعضی واقعیت‌ها نیاز داریم. در همه‌ی سه حالت مطرح شده در بخش ۶

-۱ دارای توزیع دوجمله‌ای با (n_r, p_r) و مستقل از s می‌باشد.

-۲ به شرط n_r, s نمونه‌ای تصادفی ساده از s به اندازه‌ی n_r است.

$$P(R_i = 1, R_j = 1 | n_r) = \frac{n_r}{n_r} \cdot \frac{n_r - 1}{n_r - 1} \quad \text{و} \quad P(R_i = 1 | n_r) = n_r / n \quad -۳$$

(رابطه‌ی ۲ به دست می‌آید)

$$E(Y_i^* | y_{\text{obs}}) = \bar{y}_r \quad (\Rightarrow E(Y_i^* | y, s, n_r) = \bar{y}_s \Rightarrow E(Y_i^* | y, s) = \bar{y}_s) \quad -۴$$

$$\hat{\sigma}_r^2 = \frac{1}{n_r - 1} \sum_{i \in s_r} (y_i - \bar{y}_r)^2 \quad \text{که در آن} \quad \text{var}(Y_i^* | y_{\text{obs}}) = \frac{n_r - 1}{n_r} \hat{\sigma}_r^2 \quad -۵$$

$$\hat{\sigma}^* = \frac{1}{n-1} \sum_{i \in s} (y_i - \bar{y}_s)^* \quad \text{که در آن } E(\hat{\sigma}_r^* | y, s, n_r) = \hat{\sigma}^* - \hat{\sigma}$$

$$\left(\Rightarrow \text{var}(Y_i^* | y, s, n_r) = \frac{n-1}{n} \hat{\sigma}^* \approx \hat{\sigma}^* \right)$$

$$\text{var}(\bar{Y}_r | y, s, n_r) = f \hat{\sigma}^* / n_r = \hat{\sigma}^* \left(\frac{1}{n_r} - \frac{1}{n} \right) - \hat{\sigma}$$

اثبات لم:

$$\begin{aligned} E(\hat{\theta}^* | y, s) &= E \left\{ E \left(\sum_{i \in s_r} a_i(s) y_i + \sum_{i \in s - s_r} a_i(s) Y_i^* | Y_{\text{obs}} \right) | y, s \right\} \\ &=^{(*)} E \left(\sum_{i \in s_r} a_i(s) y_i | y, s \right) + E \left(\sum_{i \in s - s_r} a_i(s) \bar{Y}_r | y, s \right) \end{aligned}$$

عبارت اول:

$$\begin{aligned} E \left(\sum_{i \in s_r} a_i(s) y_i | y, s \right) &= E \left\{ E \left(\sum_{i \in s} a_i(s) y_i R_i | y, s, n_r \right) | y, s \right\} \\ &=^{(*)} E \left(\sum_{i \in s} a_i(s) y_i \frac{n_r}{n} | y, s \right) =^{(*)} p_r \hat{\theta} \end{aligned}$$

عبارت دوم:

$$\begin{aligned} E \left(\sum_{i \in s - s_r} a_i(s) \bar{Y}_r | y, s \right) &= E \left\{ E \left(\frac{1}{n_r} \sum_{i \in s} \sum_{j \in s} a_i(s) y_j (1 - R_i) R_j | y, s, n_r \right) | y, s \right\} \\ &=^{(*)} E \left(\frac{1}{n_r} \sum_{i \in s} \sum_{j \in s, j \neq i} a_i(s) y_j \left(\frac{n_r}{n} - \frac{n_r}{n} \cdot \frac{n_r - 1}{n - 1} \right) | y, s \right) \\ &=^{(*)} \frac{1 - p_r}{n - 1} \sum_{i \in s} \sum_{j \in s, j \neq i} a_i(s) y_j = \frac{1 - p_r}{n - 1} (n \bar{a}(s) n \bar{y}_s - \hat{\theta}) \end{aligned}$$

$$\text{که در آن } \bar{a}(s) = \sum_{i \in s} a_i(s) / n$$

(۲۱) این به آن معنی است که $\hat{\theta} = n\bar{a}(s)\bar{y}_s = \bar{a}(s)\sum_{i \in s} y_i \Leftrightarrow$

اثبات قضیه:

$\hat{\theta}^* = a(s) \left(\sum_{i \in s_r} y_i + \sum_{i \in s - s_r} y_i^* \right)$ و $\hat{\theta} = a(s) \sum_{i \in s} y_i = na(s)\bar{y}_s$ با استفاده از لامبرت است.

$$\begin{aligned} E(\hat{\theta}^* | Y_{\text{obs}}) &=^{(*)} a(s)(n_r \bar{y}_r + (n - n_r) \bar{y}_r) = na(s)\bar{y}_r \\ \text{var}(\hat{\theta}^* | Y_{\text{obs}}) &=^{(*)} \{a(s)\}^* (n - n_r) \frac{n_r - 1}{n_r} \hat{\sigma}_r^* \end{aligned}$$

با این،

$$\begin{aligned} \text{var} E(\hat{\theta}^* | Y_{\text{obs}}) &= \text{var}[na(s)\bar{Y}_r] = E[n^* \{a(s)\}^* \text{var}(\bar{Y}_r | Y, s)] \\ &\quad + \text{var}\{na(s)E(\bar{Y}_r | Y, s)\} \\ &= n^* E[a(s)]^* \{E[\text{var}(\bar{Y}_r | Y, s, n_r) | Y, s] + \text{var}\{E(\bar{Y}_r | Y, s, n_r) | Y, s\}\} \\ &\quad + \text{var}\{na(s)E\{E(\bar{Y}_r | Y, s, n_r) | Y, s\}\} \\ &=^{(*)} n^* E[a(s)]^* \{\hat{\sigma}_r^* E\{(\sqrt{n_r} - 1/n) | s\} + \text{var}(\bar{Y}_s | Y, s)\} + \text{var}\{na(s)\bar{Y}_s\} \\ &= nE[a(s)]^* \hat{\sigma}_r^* [E(\sqrt{\hat{p}_r} | s) - 1] + \text{var}\hat{\theta} \\ &+ \text{var}\hat{\theta} =^{(*)} n\{E(\sqrt{\hat{p}_r}) - 1\} E[a(s)]^* \hat{\sigma}_r^* + \text{var}\hat{\theta} \end{aligned}$$

سپس،

$$\begin{aligned}
E \operatorname{var}(\hat{\theta}^* | Y_{\text{obs}}) &= E \left([a(s)]' (n - n_r) \frac{n_r - 1}{n_r} \hat{\sigma}_r^2 \right) \\
&= E \left\{ (n - n_r) (1 - 1/n_r) [a(s)]' E(\hat{\sigma}_r^2 | Y, s, n_r) \right\} \\
&= E \left\{ (n - n_r) (1/n_r - 1) [a(s)]' \hat{\sigma}_r^2 \right\} \\
&= {}^{(1)} \left[n(1 - p_r) - \{E(1/\hat{p}_r) - 1\} \right] E([a(s)]' \hat{\sigma}_r^2)
\end{aligned}$$

از رابطه‌ی (۵) دریافتیم که

$$\begin{aligned}
E(k) &= \frac{\{E(1/\hat{p}_r) - 1\} E([a(s)]' \hat{\sigma}_r^2)}{\left[(1 - p_r) - \frac{1}{n} \{E(1/\hat{p}_r) - 1\} \right] E([a(s)]' \hat{\sigma}_r^2)} \\
&= \frac{\{E(1/\hat{p}_r) - 1\}}{1 - p_r - \frac{1}{n} \{E(1/\hat{p}_r) - 1\}} \\
&\approx \frac{(1/p_r) - 1}{1 - p_r} = \frac{1}{p_r}
\end{aligned}$$

زهرا رضایی قهرودی

دکتری آمار

تهران، خیابان دکتر فاطمی، خیابان باباطهر، خیابان شهید فکوری، شماره‌ی ۱۴۵، پژوهشکده‌ی آمار.

رایانشانی: z_rezaei@srtc.ac.ir