

برآورد واریانس به روش جک‌نایف در آمارگیری‌های دوچارچوبی

مرجان نورینی

مرکز آمار ایران

چکیده: روش‌های معمول برآورد واریانس از قبیل روش خطی‌سازی سری تیلور (روش دلتا)^۱ در آمارگیری‌های چندچارچوبی^۲ عموماً مستلزم محاسبه‌ی مشتق‌های جزئی بوده و این محاسبات با افزایش تعداد چارچوب‌ها پیچیده‌تر می‌شود. برآورد واریانس به روش جک‌نایف^۳ روش دیگری است که ضمن سهولت در محاسبه، موجب کاهش چشم‌گیری در آریبی برآوردگر می‌شود. در این مقاله ابتدا به معرفی برآوردگرهای چندچارچوبی مجموع جامعه و سپس استفاده از روش جک‌نایف در برآورد واریانس برآوردگرهای دوچارچوبی و مقایسه‌ی آن با روش خطی‌سازی سری تیلور طی یک مطالعه‌ی شبیه‌سازی می‌پردازیم. واژگان کلیدی: برآورد واریانس جک‌نایف؛ برآورد واریانس خطی‌سازی؛ چارچوب‌های چندگانه؛ روش‌های باز نمونه‌گیری.

۱- مقدمه

چارچوب آماری یا فهرست واحدهای آمارگیری، اساس و مبنای یک طرح آمارگیری نمونه‌ای را تشکیل می‌دهد. گاهی ممکن است چارچوبی که تمامی واحدهای جامعه‌ی مورد مطالعه را پوشش دهد در دسترس نباشد، اما امکان دستیابی به پوشش کامل، با تلفیقی از دو یا چند چارچوب فراهم شود. در چنین حالتی به منظور دسترسی به پوشش مناسب، از دو چارچوب یا بیش‌تر به‌طور هم‌زمان استفاده می‌شود. گاهی نیز ممکن است یک چارچوب، پوشش کامل را برای جامعه‌ی مورد مطالعه فراهم کند، اما چارچوب ناقص دیگری موجود باشد که هزینه‌ی آمارگیری از آن کم‌تر از هزینه‌ی آمارگیری از چارچوب

کامل باشد. در این شرایط به دلیل پایین تر بودن هزینه‌ی آمارگیری از این چارچوب، می‌توان با هزینه‌ای مشخص و ثابت، از دو چارچوب استفاده کرده و اندازه‌ی نمونه را بزرگ‌تر و کارایی را افزایش داد. گاهی نیز ممکن است یک چارچوب فهرستی کامل، در دسترس باشد اما عملاً با گذشت زمانی نسبتاً طولانی به دلیل بروز تغییرات فراوان در آن، منبعی برای بروز خطاهای غیر نمونه‌گیری شود. از آن‌جا که یک فهرست ناحیه‌ای، کم‌تر در معرض تغییرات می‌باشد، ترکیب آن با یک چارچوب از اعضای جامعه که احتمالاً ناقص باشد، می‌تواند نتایج مفیدی را حاصل نماید. چنین آمارگیری‌هایی تحت عنوان آمارگیری‌های چندچارچوبی به کار می‌روند.

شاید بتوان گفت اولین شالوده‌ی آمارگیری‌های چندچارچوبی در سال ۱۹۴۹ با آمارگیری از فروشگاه‌های خرده‌فروشی که مجری آن دفتر سرشماری آمریکا بود گذاشته شد. سپس هارتلی [۵] نظریه‌ی مقدماتی چارچوب‌های چندگانه را توسعه داد. او با فرض این‌که اجتماع چارچوب‌ها، جامعه را پوشش می‌دهد، واحدهای جامعه را به زیرمجموعه‌های دو به دو ناسازگار شامل اجتماع‌ها و اشتراک‌های چارچوب‌های مختلف تقسیم کرد. بعد از وی لاند [۹]، فولر و بورمیستر [۴]، وگل [۱۹]، فوردو بوسکر [۳]، بانکیر [۲]، کالتون و اندرسون [۷]، اسکینر [۱۶] و اسکینر و رائو [۱۷] مقالاتی را در این زمینه ارائه نموده‌اند.

در آمارگیری‌های چندچارچوبی، برآوردهای مختلفی برای برآورد مجموع جامعه پیش‌نهاد شده است. این برآوردها عموماً توابعی ناخطی از مقادیر نمونه‌ای می‌باشند. برای برآورد واریانس این‌گونه برآوردها، روش‌های مختلفی از جمله خطی‌سازی سری تیلور پیش‌نهاد شده است.

اسکینر و رائو [۱۷] روش خطی‌سازی سری تیلور را پیش‌نهاد دادند. استفاده از این روش مستلزم محاسبه‌ی مشتق‌های جزئی بوده و این محاسبات با افزایش تعداد چارچوب‌ها پیچیده‌تر می‌شود. لیکن برآورد واریانس به روش جک‌نایف ضمن سهولت در محاسبه، موجب کاهش چشمگیری در آریبی برآوردها می‌شود. در این مقاله، ضمن استفاده از روش جک‌نایف در برآورد واریانس، به مقایسه‌ی آن با روش بسط سری تیلور می‌پردازیم. به این منظور در بخش ۲، برآوردهای چندچارچوبی معرفی می‌شوند. در بخش ۳، ابتدا شرح مختصری از روش جک‌نایف ارائه شده و سپس به کاربرد آن در آمارگیری‌های چندچارچوبی می‌پردازیم. در بخش ۴ نتایج حاصل از مقایسه و بررسی برآوردهای

چندچارچوبی و برآورد واریانس آن‌ها با استفاده از شبیه‌سازی آرایه می‌شود و در نهایت در بخش ۵ نتیجه‌گیری بیان می‌شود.

۲- برآوردهای مجموع جامعه

برای سادگی فرض کنید دو چارچوب A و B موجود است که هر دو ناقص و دارای واحدهای مشترک با یکدیگر هستند به طوری که مجموع آن‌ها روی هم، کل جامعه‌ی مورد مطالعه را پوشش می‌دهد. از چارچوب‌های A و B ، سه $(2^3 - 1)$ حوزه‌ی دو به دو ناسازگار به دست می‌آید.

حوزه‌ی a : شامل واحدهایی است که فقط در چارچوب A می‌باشند. $a = A \cap B^c$

حوزه‌ی b : شامل واحدهایی است که فقط در چارچوب B می‌باشند. $b = A^c \cap B$

حوزه‌ی ab : شامل واحدهایی است که در هر دو چارچوب می‌باشند. $ab = A \cap B$

(c نشان‌دهنده‌ی مکمل مجموعه می‌باشد).

نمادهای $N_a, N_b, N_{ab}, N_A, N_B$ به ترتیب تعداد واحدها در چارچوب‌های A و B و حوزه‌های a, b, ab می‌باشند.

دو نمونه‌ی مستقل S_A و S_B بر اساس طرح‌های نمونه‌گیری احتمالاتی $p_A(S_A)$ و $p_B(S_B)$ به اندازه‌ی n_A و n_B از دو چارچوب فوق گرفته می‌شود، به طوری که احتمال شمول نمونه‌ی حاصل از چارچوب A برابر است با $\pi_i^A = p\{i \in S_A\}$ و احتمال شمول نمونه‌ی حاصل از چارچوب B برابر است با $\pi_i^B = p\{i \in S_B\}$. ([۸])

بر اساس نمونه‌های مستقل فوق، n_a و n_{ab}^A ، تعداد واحدهای نمونه‌گیری حاصل از چارچوب A می‌باشند که به ترتیب در حوزه‌های a و ab قرار دارند. به همین ترتیب n_b و n_{ab}^B نیز تعداد واحدهای نمونه‌گیری حاصل از چارچوب B می‌باشند که به ترتیب در حوزه‌های b و ab هستند.

با فرض این که Y_a, Y_b, Y_{ab} به ترتیب مجموع جامعه در حوزه‌های a, b و ab باشند، داریم:

$$(۱) \quad Y = Y_a + Y_{ab} + Y_b$$

هدف برآورد Y است.

چندین برآوردهای نقطه‌ای تحت عنوان برآوردهای دو چارچوبی برای برآورد Y پیشنهاد

شده است که همه‌ی آن‌ها به فرم $\hat{Y} = \hat{Y}_a + \hat{Y}_{ab} + \hat{Y}_b$ می‌باشند. هر یک از این برآوردها، بسته به این که اطلاعات حاصل از دو نمونه برای برآورد Y چگونه با هم ترکیب می‌شوند با هم تفاوت‌هایی دارند.

با فرض این که N_B و N_A معلوم بوده و $N_a > 0$ و $N_b > 0$ ، وزن‌های نمونه‌گیری w_i^A و w_i^B عبارت‌اند از:

$$w_i^A = N_A [\pi_i^A \sum_{j \in S_A} (\frac{1}{\pi_j^A})]^{-1}, \quad w_i^B = N_B [\pi_i^B \sum_{j \in S_B} (\frac{1}{\pi_j^B})]^{-1}$$

دو متغیر نشانگر زیر را برای هر یک از چارچوب‌های A و B در نظر می‌گیریم:

$$\delta_i^A = \begin{cases} 1 & \text{واحد } i \text{ متعلق به چارچوب } A \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

$$\delta_i^B = \begin{cases} 1 & \text{واحد } i \text{ متعلق به چارچوب } B \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

در نتیجه برآوردها در سه حوزه‌ی a و b و ab به صورت زیر است:

$$\begin{aligned} \hat{Y}_a^A &= \sum_{i \in S_A} w_i^A (1 - \delta_i^B) y_i, & \hat{Y}_{ab}^A &= \sum_{i \in S_A} w_i^A \delta_i^B y_i \\ \hat{Y}_b^B &= \sum_{i \in S_B} w_i^B (1 - \delta_i^A) y_i, & \hat{Y}_{ab}^B &= \sum_{i \in S_B} w_i^B \delta_i^A y_i \end{aligned}$$

برآوردها اندازه‌ی هر یک از حوزه‌ها نیز به صورتی مشابه با قرارداد $y_i = 1$ در تعاریف \hat{Y}_a^A ، \hat{Y}_b^B ، \hat{Y}_{ab}^A ، \hat{Y}_{ab}^B به دست می‌آید. همچنین دو برآوردها دیگر را نیز به صورت زیر تعریف می‌کنیم:

$$\begin{aligned} \hat{Y}_{ab}(\theta) &= \theta \hat{Y}_{ab}^A + (1 - \theta) \hat{Y}_{ab}^B \\ \hat{N}_{ab}(\theta) &= \theta \hat{N}_{ab}^A + (1 - \theta) \hat{N}_{ab}^B \end{aligned}$$

که مقدار بهینه‌ی θ واریانس $\hat{Y}_{ab}(\theta)$ را مینیمم می‌کند.

۱-۲- برآوردهای هارتلی و فولر- بورمیستر

هارتلی [۵] برآوردگر دوچارچوبی زیر را برای برآورد مجموع جامعه پیش‌نهاد کرد:

$$(۲) \quad \hat{Y}_H(\theta) = \hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}(\theta)$$

فولر و بورمیستر [۴] هم برآوردگر زیر را پیش‌نهاد دادند:

$$(۳) \quad \hat{Y}_{FB}(\beta_1, \beta_2) = \hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}(\beta_1) + \beta_2(\hat{N}_{ab}^A - \hat{N}_{ab}^B)$$

مقادیر بهینه‌ی پارامترهای θ ، β_1 و β_2 به ترتیب واریانس برآوردهای $\hat{Y}_H(\theta)$ و $\hat{Y}_{FB}(\beta_1, \beta_2)$ را مینیمم کرده و بنا بر این به کوواریانس‌های \hat{Y}_{ab}^A و \hat{Y}_{ab}^B بستگی دارند. اما در عمل این کوواریانس‌ها مجهول بوده و باید از روی نمونه برآورد شوند. در نتیجه $\hat{Y}_H(\hat{\theta}_H)$ و $\hat{Y}_{FB}(\hat{\beta}_{FB})$ به‌طور کلی توابعی خطی از y نخواهند بود و برای هر متغیر پاسخ باید به‌طور جداگانه محاسبه شوند. برای مثال اگر \hat{Y}_{H_1} ، \hat{Y}_{H_2} و \hat{Y}_{H_3} برآوردهای هارتلی برای تعداد کل بیمارانی که دچار تنگی نفس در گروه‌های سنی ۰-۱۶ و ۱۷-۴۵ و ۴۵ به بالا باشد، آن‌گاه $\hat{Y}_{H_1} + \hat{Y}_{H_2} + \hat{Y}_{H_3}$ لزوماً مساوی برآورد هارتلی تعداد کل بیماران دچار تنگی نفس در جامعه نخواهد بود. در نتیجه فولرو بورمیستر [۴] بر آن شدند برآوردگری ارائه دهند که از یک مجموعه وزن‌های یکسان برای همه‌ی متغیرهای پاسخ استفاده کند. آن‌ها $\hat{Y}_{FB,SRS}$ را برای برآورد Y معرفی کردند که این برآوردگر، تحت طرح نمونه‌گیری تصادفی ساده طراحی شده بود و نمی‌توانست مستقیماً با طرح‌های نمونه‌گیری پیچیده به کار رود. چون عموماً به‌صورتی ناسازگار با این نوع نمونه‌گیری‌ها طراحی شده بود.

۲-۲- برآوردگر ماکسیمم درست‌نمایی نما^۴

اسکینز [۱۶] از $\hat{Y}_{FB,SRS}$ تفسیری به‌صورت یک برآوردگر ماکسیمم درست‌نمایی داشت. چون این برآوردگر از مجموعه‌ی یکسانی از وزن‌ها استفاده می‌کرد، به برآوردگر فولر-بورمیستر ارجحیت داشت. ولی این برآوردگر، مستقیماً نمی‌توانست با طرح‌های نمونه‌گیری پیچیده به کار رود. چون به‌طور کلی به‌صورتی ناسازگار با این نوع طرح‌های نمونه‌گیری طراحی شده بود.

اسکینر و راتو [۱۸] تصمیم گرفتند این برآوردگر را طوری اصلاح کنند که به برآوردگری سازگار تحت طرح‌های نمونه‌گیری پیچیده نائل آیند. آن‌ها در این حالت، برآوردگر ماکسیمم درستنمایی‌نما را برای برآورد مجموع پیشنهاد کردند. این برآوردگر، اصلاح‌شده‌ی یک برآوردگر ماکسیمم درستنمایی تحت طرح نمونه‌گیری تصادفی ساده، برای رسیدن به برآوردگری سازگار با طرح‌های نمونه‌گیری پیچیده بود. آن‌ها برآوردگر ماکسیمم درستنمایی‌نما را به صورت زیر ارائه دادند:

$$\hat{Y}_{PML}(\theta) = \frac{N_A - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_a^A} \hat{Y}_a^A + \hat{N}_{ab}^{PML}(\theta) \frac{\hat{Y}_{ab}(\theta)}{\hat{N}_{ab}(\theta)} + \frac{N_B - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_b^B} \hat{Y}_b^B \quad (4)$$

به طوری که $\hat{N}_{ab}^{PML}(\theta)$ تابعی از \hat{N}_{ab}^A ، \hat{N}_{ab}^B و θ بوده و کوچک‌ترین ریشه‌ی معادله‌ی درجه‌ی دو، به صورت زیر است:

$$\left(\frac{\theta}{N_B} + \frac{1-\theta}{N_A}\right)x^2 - \left(1 + \frac{\theta}{N_B}\hat{N}_{ab}^A + \frac{1-\theta}{N_A}\hat{N}_{ab}^B\right)x + \hat{N}_{ab}(\theta) = 0$$

این برآوردگر برخلاف برآوردگرهای هارتلی و فولر-بورمیستر، تابعی خطی از y است. اسکینر و راتو [۱۷] انتخاب $\theta = \theta_p$ را برای مینیمم کردن واریانس مجانبی با $\hat{N}_{ab}^{PML}(\theta)$

$$\theta_p = \frac{N_a N_{BV}(\hat{N}_{ab}^B)}{N_a N_{BV}(\hat{N}_{ab}^B) + N_b N_{AV}(\hat{N}_{ab}^A)} \quad (5)$$

پیشنهاد کردند. در عمل N_a و N_b و واریانس‌ها در رابطه‌ی (۵) مجهول بوده و باید از روی داده‌ها برآورد شوند.

۳-۲- برآوردگرهای تک‌چارچوبی^۵

بانکیر [۲]، کالتون و اندرسون [۷] و اسکینر [۱۶] برای برآورد مجموع جامعه با در نظر گرفتن مشاهدات به‌گونه‌ای که آن‌ها از یک چارچوب با وزن‌های اصلاح‌شده برای مشاهدات در حوزه‌ی متداخل ab نمونه‌گیری شده‌اند، برآورگری ارائه دادند. وزن‌های اصلاح‌شده برای برآوردگرهای تک‌چارچوبی اسکینر [۱۶] و کالتون و اندرسون [۷] به

شناسایی واحدهای مشابه در نمونه‌ها نیازی ندارد. این وزن‌ها عبارت‌اند از:

$$w_i = \begin{cases} (\pi_i^A + \pi_i^B)^{-1} & i \in ab \\ (\pi_i^A)^{-1} & i \in a \\ (\pi_i^B)^{-1} & i \in b \end{cases}$$

آن‌ها برآوردگر زیر را برای برآورد مجموع پیش‌نهادهای کردند:

$$(۶) \quad \hat{Y}_{SF} = \sum_{i \in s_A} w_i y_i + \sum_{i \in s_B} w_i y_i$$

این برآوردگر قابل تعمیم به بیش از دو چارچوب می‌باشد و از هیچ‌گونه اطلاعات کمکی در ارتباط با N_A و N_B استفاده نمی‌کند. از آن جایی‌که استفاده از اطلاعات کمکی، باعث کاهش واریانس و افزایش دقت برآوردگر می‌شود [۲]، اسکینر و راثو [۱۷] دو روش برای تصحیح این برآوردگر ارائه دادند که عبارت‌اند از:

(۱) روش تعدیل نسبتی [۲]

(۲) روش رگرسیونی

اسکینر و راثو [۱۷] نشان دادند که فرایند تعدیل به برآوردگر تجربی زیر همگرا است:

$$(۷) \quad \hat{Y}_{SFrake} = \frac{N_A - \hat{N}_{ab}^{rake}}{\hat{N}_a} \hat{Y}_a + \frac{\hat{N}_{ab}^{rake}}{\hat{N}_{abs}} \hat{Y}_{abs} + \frac{N_B - \hat{N}_{ab}^{rake}}{\hat{N}_b} \hat{Y}_b$$

که در آن:

$$\hat{Y}_{abs} = \sum_{s_A} w_i \delta_i^B y_i + \sum_{s_B} w_i \delta_i^A y_i$$

$$\hat{N}_{abs} = \sum_{s_A} w_i \delta_i^B + \sum_{s_B} w_i \delta_i^A$$

و \hat{N}_{ab}^{rake} کوچک‌ترین ریشه‌ی معادله‌ی درجه‌ی دو به‌صورت زیر است:

$$(۸) \quad \hat{N}_{abs} x^2 - [\hat{N}_{abs}(N_A + N_B) + \hat{N}_{as}^A \hat{N}_{bs}^B] x + \hat{N}_{abs} N_A N_B = 0$$

$$\text{که } \hat{N}_{bs}^B = \sum_{s_B} w_i \delta_i^B \text{ و } \hat{N}_{as}^A = \sum_{s_A} w_i \delta_i^A$$

جزئیات محاسبه‌ی این برآوردگر در [۱] آورده شده است.

تصحیح با N_A و N_B معلوم به عنوان دو متغیر کمکی مورد استفاده که به ترتیب با \hat{N}_{AS} و \hat{N}_{BS} برآورد می‌شوند، از طریق روش رگرسیونی، برآوردگر زیر را نتیجه می‌دهد:

$$(۹) \quad \hat{Y}_{SFree} = \hat{Y}_S + \hat{\beta}_{S\Upsilon}^T [N_A - \hat{N}_{AS}, N_B - \hat{N}_{BS}] = 0$$

به طوری که مقدار بهینه $\hat{\beta}_{S\Upsilon}$ برابر است با:

$$\hat{\beta}_{S\Upsilon} = -Cov\left\{\left[\frac{\hat{N}_{AS}}{V(\hat{N}_{AS})}, \frac{\hat{N}_{BS}}{V(\hat{N}_{BS})}\right], \hat{Y}_S\right\}$$

۳- برآورد واریانس

روش جک‌نایف برآوردهایی از پارامتر مورد مطالعه را از هر یک از زیرنمونه‌های نمونه اصلی به دست آورده و سپس برآورد واریانس برآوردگر نمونه اصلی را از تغییرپذیری بین برآوردهای این زیر نمونه‌ها نتیجه می‌دهد.

نخستین بار کنولی [۱۰]، جک‌نایف را به عنوان روشی بر اساس حذف هر بار یک مشاهده از مجموعه داده‌های اولیه و محاسبه مجدد برآوردگر با استفاده از بقیه داده‌ها برای کاهش ارببی برآوردگر ضریب همبستگی پیاپی معرفی نمود. در مقاله‌ای دیگر [۱۱] این روش را تعمیم داده و خواص کلی کاهش ارببی آن را در یک جامعه متنهای بیان کرد. سپس توکی [۱۸] پیش‌نهاد کرد که این روش علاوه بر کاهش ارببی، می‌تواند برای برآورد واریانس نیز استفاده شود.

فرض می‌کنیم X_1, \dots, X_n متغیرهای تصادفی مستقل و هم توزیع با مقادیر حقیقی از یک توزیع نامعلوم F باشند و آماره $\hat{\theta}(X_1, \dots, X_n)$ ، برآوردگر برای پارامتر نامعلوم θ بوده که از روی نمونه اصلی محاسبه می‌شود. $\hat{\theta}_i$ برآوردگر برای $\hat{\theta}$ بر اساس نمونه‌ای به اندازه $n-1$ است که بعد از حذف i امین واحد نمونه به دست می‌آید. آنگاه برآوردگر جک‌نایف $\hat{\theta}$ این نمونه عبارت است از:

$$\hat{\theta}_J = (n-1)^{-1} \sum_{i=1}^n [n\hat{\theta} - (n-1)\hat{\theta}_i]$$

برآورد واریانس جک‌نایف برآوردگر بالا عبارت است از:

$$v_J(\hat{\theta}) = n^{-1} (n-1) \sum_{i=1}^n (\hat{\theta}_i - \hat{\theta})^2$$

همان‌طور که ملاحظه می‌شود روش جک‌نایف به فرضیات مدل بستگی کمتری داشته و نیازی به فرمول‌های نظری که در روش سنتی می‌باشد، ندارد و در مجموع به‌عنوان یک روش ناپارامتری در تحلیل‌های آماری شناخته می‌شود. با این حال، جک‌نایف نیاز به محاسبه‌ی مکرر آماره به تعداد n مرتبه دارد که واقعاً این کار در قدیم امکان‌پذیر نبود.

در این بخش، برآورد جک‌نایف واریانس برآوردهای دوچارچوبی ارایه شده و سپس هم‌ارزی مجانبی آن با برآوردهای خطی‌سازی واریانس (برآوردهای که اسکینر و رائو با استفاده از روش خطی‌سازی سری تیلور به دست آوردند [۱۷]) اثبات می‌شود.

فرض کنید چارچوب‌های A و B به‌ترتیب دارای H و L طبقه باشند به‌طوری‌که طبقات h و l در هر یک از آن‌ها به‌ترتیب شامل N_h^A و N_l^B واحد بوده و در مجموع \tilde{N}_l^B و \tilde{N}_h^A واحد نمونه‌گیری اولیه (psu) دارند که به‌ترتیب \tilde{n}_l^B و \tilde{n}_h^A واحد از آن‌ها نمونه‌گیری می‌شود. به‌طوری‌که $\tilde{n}_B = \sum_{l=1}^L \tilde{n}_l^B$ و $\tilde{n}_A = \sum_{h=1}^H \tilde{n}_h^A$ کل نمونه‌ای است که

به‌ترتیب از چارچوب‌های A و B انتخاب می‌شود. وزن طبقات h و l در چارچوب‌های A و B به‌ترتیب با $W_h^A = N_h^A / N_A$ و $W_l^B = N_l^B / N_B$ نشان داده می‌شود [۸]. لازم به یادآوری است که اغلب در نمونه‌گیری‌های چندمرحله‌ای، نمونه‌ی بدون جایگذاری از واحدهای نمونه‌گیری اولیه با احتمال متناسب با اندازه (به‌دلیل کارایی بالای نمونه‌گیری بدون جایگذاری نسبت به نمونه‌گیری با جایگذاری) به عمل آمده و سپس در مرحله‌ی برآورد واریانس، برای صرفه‌جویی در زمان و هزینه و اجتناب از محاسبات بالقوه پرهزینه احتمال‌های شمول توأم، نمونه‌گیری را به‌گونه‌ای در نظر می‌گیرند که گویی با جایگذاری انتخاب شده است [۱۵].

این تقریب عموماً منجر به بیش برآورد شدن واریانس برآوردهای مجموع می‌شود اما در صورتی که کسرهای نمونه‌گیری مرحله‌ی اول بزرگ نباشد، اریبی نسبی بزرگ نخواهد بود. برآوردهای جک‌نایف و دیگر برآوردهای واریانس بازنمونه‌گیری از این روش در برآورد واریانس استفاده می‌کنند.

حال فرض کنید که نمونه‌هایی به اندازه‌ی (≥ 2) \tilde{n}_l^B و \tilde{n}_h^A به‌ترتیب با احتمال‌های

شامل متناسب با اندازه، $\tilde{\pi}_{hi}^A = \tilde{n}_h^A p_{hi}^A$ و $\tilde{\pi}_{lj}^B = \tilde{n}_l^B p_{lj}^B$ از چارچوب‌های A و B انتخاب می‌شود، به طوری که p_{hi}^A و p_{lj}^B به ترتیب برابر با احتمال انتخاب واحد i ام و j ام در طبقات h و l متناسب با اندازه‌ی psu ها بوده و $\sum_i p_{hi}^A = 1$ و $\sum_j p_{lj}^B = 1$ ($p_{hi}^A = u_i / \sum_A u_i$ و $p_{lj}^B = u_j / \sum_B u_j$ که u_i اندازه‌ی تقریبی i امین psu است).

با در نظر گرفتن دو بردار $A = (Y_a, Y_{ab}, N_{ab}, N_A)^T$ و $B = (Y_b, Y_{ab}, N_{ab}, N_B)^T$ ، آن‌گاه برآورد A برابر است با:

$$\hat{A} = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} \frac{\hat{A}_{hi}}{\tilde{\pi}_{hi}^A} = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} \frac{\hat{A}_{hi}}{\tilde{n}_h^A p_{hi}^A} = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} \frac{N_h^A a_{hi}}{\tilde{n}_h^A} = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} N_h^A \bar{a}_h$$

که \hat{A}_{hi} برآوردگر نارایب بردار مجموع A_{hi} بر اساس نمونه‌گیری در مرحله‌ی دوم و مراحل بعدی بوده و $a_{hi} = \hat{A}_{hi} / (N_h^A p_{hi}^A)$ می‌باشد. برآوردگر $\hat{B} = \sum_{l=1}^L \sum_{j=1}^{\tilde{n}_l^B} N_l^B b_{lj} / \tilde{n}_l^B$ نیز به صورتی مشابه تعریف می‌شود. تحت فرضیات نمونه‌گیری با جایگذاری، a_{hi} ها، برآوردگرهای نارایب مستقل و هم‌توزیع میانگین جامعه در طبقه‌ی h ام چارچوب A با $E(a_{hi}) = \bar{A}_h = A_h / N_h^A$ می‌باشند در حالی که برای $h \neq h'$ ، a_{hi} و $a_{h'i'}$ مستقل بوده اما لزوماً هم‌توزیع نیستند. به طور مشابه b_{lj} ها هم برآوردگرهای نارایب مستقل و هم‌توزیع بردار میانگین جامعه در طبقه‌ی l ام چارچوب B با $E(b_{lj}) = \bar{B}_l = B_l / N_l^B$ می‌باشند.

پارامتری به صورت $\tau = g(\bar{A}, \bar{B})$ را که تابعی از میانگین‌های جامعه، $\bar{A} = A/N$ و

$\bar{B} = B/N$ است در نظر می‌گیریم. میانگین‌های بالا به ترتیب با $\hat{A} = \sum_{h=1}^H w_h^A \bar{a}_h$ و

$$\hat{B} = \sum_{l=1}^L w_l^B \bar{b}_l$$

برآورد شده و $\hat{\tau} = g(\hat{A}, \hat{B})$ با $\hat{\tau}$ برآورد می‌شود [۸].

برآوردگرهای دوچارچوبی می‌توانند در قالب $\hat{\tau}$ با در نظر گرفتن مجموع جامعه به صورت $Y = N\bar{Y} = N g(\bar{A}, \bar{B})$ بیان شوند. در این صورت میانگین \bar{Y} برابر است با:

$$(۱۰) \quad g(\bar{A}, \bar{B}) = \bar{Y} = \left(\frac{N_A}{N}\right) (\bar{A}_\vee + \theta \bar{A}_\vee) + \left(\frac{N_B}{N}\right) \{\bar{B}_\vee + (1 - \theta) \bar{B}_\vee\}$$

در این صورت برآوردگر هارتلی برابر با $\hat{Y}_H(\theta_H) = N g(\hat{A}, \hat{B})$ خواهد بود [۱].

ماتریس واریانس بردار \hat{A} برابر با $\Sigma_A = \sum_{h=1}^H (W_h^A)^\vee \Sigma_h^A / \tilde{n}_h^A$ بوده که با

$S_A = \sum_{h=1}^H (W_h^A)^\vee S_h^A / \tilde{n}_h^A$ برآورد می‌شود. به‌طور مشابه، ماتریس واریانس \hat{B} یعنی

Σ_B نیز با $S_B = \sum_{l=1}^L (W_l^B)^\vee S_l^B / \tilde{n}_l^B$ برآورد می‌شود که S_l^B و S_h^A به‌ترتیب برآوردگرهای

ماتریس واریانس $\sqrt{\tilde{n}_l^B} \bar{b}_l$ و $\sqrt{\tilde{n}_h^A} \bar{a}_h$ می‌باشند و عبارت‌اند از:

$$S_h^A = (\tilde{n}_h^A - 1) \sum_{i=1}^{\tilde{n}_h^A} (a_{hi} - \bar{a}_h)(a_{hi} - \bar{a}_h)^T$$

$$S_l^B = (\tilde{n}_l^B - 1) \sum_{j=1}^{\tilde{n}_l^B} (b_{lj} - \bar{b}_l)(b_{lj} - \bar{b}_l)^T$$

برای اثبات سازگاری برآوردگر جک‌نایف واریانس و همچنین هم‌ارزی مجانبی آن با برآوردگر خطی‌سازی واریانس، نیاز به شرط‌های زیر است:

(آ) فرض کنیم که وقتی $\tilde{n}_A, \tilde{n}_B \rightarrow \infty$ آن‌گاه: $W_h^A \tilde{n}_h^A / \tilde{n}_h^A = O(1)$ و $W_l^B \tilde{n}_l^B / \tilde{n}_l^B = O(1)$

(ب) فرض می‌کنیم $g_A(\tilde{a}, \tilde{b})$ برداری q بعدی از مشتق‌های اول تابع $g(\cdot)$ نسبت به مؤلفه‌های a ، $g_B(\tilde{a}, \tilde{b})$ برداری r بعدی از مشتق‌های اول تابع $g(\cdot)$ نسبت به مؤلفه‌های b باشد که در \tilde{a} و \tilde{b} مقدار دهی می‌شوند. همچنین فرض می‌کنیم $g_A''(\tilde{a}, \tilde{b})$ یک ماتریس q از مشتق‌های دوم، $\partial^2 g(\partial a_j \partial a_k)$ و $g_B''(\tilde{a}, \tilde{b})$ یک ماتریس r بعدی از مشتق‌های دوم $\partial^2 g(\partial b_j \partial b_k)$ در \tilde{a} و \tilde{b} باشند و g_A'' و g_B'' در یک همسایگی از (\bar{A}, \bar{B}) پیوسته و کراندار باشند.

(ج) $\tilde{n} = \tilde{n}_A + \tilde{n}_B$ به‌طوری که $\tilde{n}_A / \tilde{n} \rightarrow k \in (0, 1)$

شرط‌های بالا توسط رائو و وو [۱۲] برای تحقیق و بررسی خواص برآوردگرهای

واریانس در طرح نمونه‌گیری طبقه‌بندی چندمرحله‌ای مورد استفاده قرار گرفته است. شرط آذعان می‌دارد که تعداد طبقات در چارچوب‌ها می‌تواند کراندار یا بی‌کران باشد. اگر اندازه‌ی نمونه‌ی حاصل از طبقات در چارچوب‌های A و B کراندار باشد یعنی $\tilde{n}_h^A = O(1)$ و $\tilde{n}_l^B = O(1)$ آن‌گاه شرط آهم‌ارز است با $W_h^A = O(\tilde{n}_A^{-1})$ و $W_h^B = O(\tilde{n}_B^{-1})$.

در مجموع فرض می‌شود $\sum_l W_l^B \Sigma_l^B = O(1)$ و $\sum_h W_h^A \Sigma_h^A = O(1)$ یعنی متوسط وزنی کواریانس درون طبقات کراندار باشد. شرط ج نیز تضمین می‌کند که نمونه‌ی حاصل از یک چارچوب، به‌طور مجانبی روی نمونه‌ی چارچوب دیگر هیچ تأثیری نمی‌گذارد.

قضیه‌ی ۱. فرض می‌کنیم شروط آ تا ج برقرار باشند. آن‌گاه:

$$\text{Var}(\hat{\tau}) = g_A^T(\bar{A}, \bar{B}) \Sigma_A g_A(\bar{A}, \bar{B}) + g_B^T(\bar{A}, \bar{B}) \Sigma_B g_B(\bar{A}, \bar{B}) + o(\tilde{n}^{-1}) \quad (11)$$

($o(\tilde{n}^{-1})$ عبارت از جملاتی از مرتبه‌ی کوچک‌تری از \tilde{n}^{-1} است)

برآوردگر خطی‌سازی واریانس بالا عبارت است از:

$$v_L(\hat{\tau}) = g_A^T(\hat{A}, \hat{B}) S_A g_A(\hat{A}, \hat{B}) + g_B^T(\hat{A}, \hat{B}) S_B g_B(\hat{A}, \hat{B}). \quad (12)$$

اثبات قضیه و همچنین اثبات سازگاری برآوردگر $v_L(\hat{\tau})$ در [۱] آمده است. فرض می‌کنیم $\hat{\tau}_{(hi)}^A$ برآوردگری به‌صورت $\hat{\tau}$ باشد که بعد از حذف مشاهده‌ی i امین psu نمونه‌ی طبقه‌ی h ام به دست آمده و $\hat{\tau}_{(hi)}^A = g(\hat{A}_{(hi)}, \hat{B})$ که $\hat{A}_{(hi)}$ نیز برآوردگر \bar{A} بعد از حذف i امین psu نمونه‌ی طبقه‌ی h ام در چارچوب A می‌باشد. به‌طور مشابه $\hat{\tau}_{(lj)}^B = g(\hat{A}, \hat{B}_{(lj)})$ که $\hat{B}_{(lj)}$ نیز برآوردگر \bar{B} بعد از حذف j امین psu نمونه‌ی طبقه‌ی l ام در چارچوب B است. برآوردگر جک‌نایف واریانس $\hat{\tau}$ (با توجه به استقلال نمونه‌ها) عبارت است از (جزئیات در [۱] دیده شود):

$$(۱۳) \quad v_J(\hat{t}) = \sum_{h=1}^H \frac{\hat{n}_h^A - 1}{\hat{n}_h^A} \sum_{i=1}^{\hat{n}_h^A} (\hat{t}_{(hi)}^A - \hat{t})^2 + \sum_{l=1}^L \frac{\hat{n}_l^B - 1}{\hat{n}_l^B} \sum_{j=1}^{\hat{n}_l^B} (\hat{t}_{(lj)}^B - \hat{t})^2$$

رائو و لوهر [۱۴] در قضیه‌ی ۲ (قضیه‌ی زیر) ثابت کردند که $v_J(\hat{t})$ و $v_L(\hat{t})$ به‌طور مجانبی هم‌ارز هستند.

قضیه‌ی ۲. با فرض این که شروط آ-ج برقرار باشند. آن‌گاه:

$$(۱۴) \quad v_J(\hat{t}) = v_L(\hat{t}) + o_p(\tilde{n}^{-1})$$

$o_p(\tilde{n}^{-1})$ عبارت است از جملاتی از مرتبه‌ی کوچک‌تری از \tilde{n}^{-1} در احتمال^۷.

برآوردگر جک‌نایف واریانس بالا در حالتی که \hat{t} تابع همواری از \hat{A} و \hat{B} باشد، مورد بررسی قرار گرفته و سازگاری و نیز هم‌ارزیش با برآوردگر خطی‌سازی واریانس، طی قضیه‌های ۱ و ۲ ثابت می‌شود [۱]. در بین برآوردگرهای دوچارچوبی، برآوردگرهای SF_{rake} و به‌صورت توابع همواری از میانگین جوامع، قابل بیان هستند. دیگر برآوردگرها نیز مادامی که پارامترهای θ_p ، θ_H ، β_{FB} و β_{S_2} ثابت بوده و از روی داده‌ها برآورد نشوند، نیز توابع همواری از میانگین جوامع خواهند بود. بنا بر این قضیه‌ی ۲ بیان می‌دارد که برآوردگر جک‌نایف واریانس برای شکل بهینه‌ی هر برآوردگر سازگار است [۱۴].

۱-۳- جک‌نایف کامل و اصلاح‌شده

برآوردگرهای $\hat{\theta}_p$ ، $\hat{\theta}_H$ ، $\hat{\beta}_{FB}$ و $\hat{\beta}_{S_2}$ همگی توابعی از ماتریس‌های واریانس S_B و S_A هستند و به‌طور کلی نمی‌توانند به‌صورت توابع همواری از میانگین در نمونه‌گیری طبقه‌بندی بیان شوند (واریانس‌ها معمولاً نامعلوم بوده و باید از روی نمونه برآورد شوند). بنا بر این قضیه‌ی ۲ همیشه نمی‌تواند به‌طور مستقیم با برآوردگرهایی که توابعی از S_A و S_B هستند به کار رود. مشکل دیگر که در استفاده از روش جک‌نایف در خیلی از نمونه‌های طبقه‌بندی رخ می‌دهد، هنگام محاسبه‌ی $S_{(hi)}^A$ و $S_{(lj)}^B$ است، چون

$$(۱۵) \quad S_{(hi)}^A = S_A + \frac{(W_h^A)^\gamma}{\tilde{n}_h^{A-\gamma}} \left[\frac{\gamma S_h^A}{\tilde{n}_h^A} - \frac{\tilde{n}_h^A}{(\tilde{n}_h^A - 1)^\gamma} (a_{hi} - \bar{a}_h)(a_{hi} - \bar{a}_h)^T \right]$$

در صورتی که $\tilde{n}_h^A = 2$ باشد رابطه‌ی (۱۵) قابل محاسبه نیست. به‌طور مشابه $S_{(ij)}^B$ نیز وقتی که $\tilde{n}_i^B = 2$ قابل محاسبه نیست. برای رفع این مشکل راثو و لوهر [۱۴] یک برآورد جک‌نایف اصلاح‌شده در حالتی که $\tilde{n}_h^A = 2$ یا $\tilde{n}_i^B = 2$ ارائه دادند. جک‌نایف (جک‌نایف اصلاح‌شده برای طرح نمونه‌گیری با ۲ واحد نمونه‌گیری اولیه در هر طبقه که بعداً به معرفی آن می‌پردازیم) حتی زمانی که برآوردگرها به S_A و S_B بستگی دارند، برآورد سازگاری از واریانس می‌دهد [۱۴].

در \hat{t} فرض بر این است که پارامترهای θ_H ، θ_p ، β_{S_Y} و β_{F_B} ثابت بوده و از روی نمونه برآورد می‌شوند. اما اگر این پارامتر ثابت نبوده و به S_A و S_B بستگی داشته باشند باید اثر این برآوردگرها در \hat{t} منظور شود. فرض می‌کنیم که

$$\hat{t} = g(\hat{A}, \hat{B}) = f(\hat{A}, \hat{B}, \beta)$$

$$\beta = [E_A + E_B]^{-1} [e_A + e_B]$$

که هر عنصر ماتریس p بعدی E_A و هر عنصر بردار r بعدی e_A ترکیب خطی از عناصر \sum_A هستند و این ویژگی برای E_B و e_B برقرار است. آنگاه برآوردگرهای $\hat{Y}_{SFreg}(\hat{\beta}_{S_Y})$ و $\hat{Y}_{FB}(\hat{\beta}_{F_B})$ ، $\hat{Y}_H(\hat{\theta}_H)$ ، $\hat{Y}_{PML}(\hat{\theta}_p)$ همگامی به‌صورت $N\hat{\zeta} = Nf(\hat{A}, \hat{B}, \hat{\beta})$ می‌باشند به‌طوری که $\hat{\beta}$ برآورد S_A و S_B را در β جایگزین می‌کند. برای مثال اگر بردارهای A و B را به‌صورت زیر در نظر بگیریم:

$$A = (Y_a, Y_{ab}, N_{ab}, N_A)^T \quad B = (Y_a, Y_{ab}, N_{ab}, N_B)^T$$

آنگاه برآوردگر هارتلی $\hat{Y}_H(\hat{\theta}_H)$ را می‌توان به‌صورت زیر نوشت:

$$\hat{Y}_H(\hat{\theta}_H) = \hat{Y}_H(\theta_H) + (\hat{\beta}_H - \beta_H) \left[N_A(\hat{A}_\gamma - \bar{A}_\gamma) - N_B(\hat{B}_\gamma - \bar{B}_\gamma) \right]$$

به‌طوری که مقدار بهینه‌ی β_H برابر است با:

$$\beta_H = - \frac{k^\gamma \sum_A (1, \gamma) - \sum_B (1, \gamma) - \sum_B (\gamma, \gamma)}{k^\gamma \sum_A (\gamma, \gamma) + \sum_B (\gamma, \gamma)}, \quad k = \frac{N_A}{N_B}$$

$\Sigma_A(1,2)$ عبارت است از عنصر $(1,2)$ امین Σ_A . برای برآوردگر ماکسیمم درستنمایی‌نما، پارامتر مناسب عبارت است از $\beta_p = k^T \Sigma_A(3,3) / \Sigma_B(3,3)$. فرض کنید $\hat{\zeta}_{(hi)}^A$ ، برآوردگری به صورت $\hat{\zeta}$ باشد که بعد از حذف مشاهدات i امین psu نمونه‌ی طبقه‌ی h حاصل شده است. آنگاه $(\hat{A}_{(hi)}, \hat{B}, \hat{\beta}_{(hi)}^A)$ که $\hat{\zeta}_{(hi)}^A = f(\hat{A}_{(hi)}, \hat{B}, \hat{\beta}_{(hi)}^A)$ که $\hat{\beta}_{(hi)}^A$ نیز برآوردگر β با استفاده از S_B و $S_{(hi)}^A$ با فرض $\tilde{n}_h^A > 2$ می‌باشد. به طور مشابه $\hat{\zeta}_{(lj)}^B = f(\hat{A}, \hat{B}_{(lj)}, \hat{\beta}_{(lj)}^B)$ است به شرط این که $\tilde{n}_l^B > 2$ باشد. آنگاه برآوردگر جک‌نایف واریانس $\hat{\zeta}$ عبارت است از:

$$v_J(\hat{\zeta}) = \sum_{h=1}^H \frac{\tilde{n}_h^A - 1}{\tilde{n}_h^A} \sum_{i=1}^{\tilde{n}_h^A} (\hat{\zeta}_{(hi)} - \hat{\zeta})^2 + \sum_{l=1}^L \frac{\tilde{n}_l^B - 1}{\tilde{n}_l^B} \sum_{j=1}^{\tilde{n}_l^B} (\hat{\zeta}_{(lj)} - \hat{\zeta})^2$$

(۱۶)

برآوردگر جک‌نایف اصلاح شده‌ی واریانس، $v_{MJ}(\hat{\zeta})$ نیز به همان صورت $v_J(\hat{\zeta})$ است با این تفاوت که به جای $\hat{\beta}_{(hi)}^A$ یا $\hat{\beta}_{(lj)}^B$ از $\hat{\beta}$ استفاده کرده و قابل کاربرد با $\tilde{n}_h^A = 2$ یا $\tilde{n}_l^B = 2$ می‌باشد.

برای به دست آوردن برآوردگر خطی واریانس $\hat{\zeta}$ با استفاده از خطی‌سازی سری تیلور آن با فرض این که f تابعی مشتق‌پذیر با مشتقات جزئی مرتبه‌ی دوم پیوسته و کراندار در یک همسایگی از $f(\bar{A}, \bar{B}, \beta)$ است و $\partial f(\bar{A}, \bar{B}, \beta) / \partial \beta$ و همچنین تحت شرط ج داریم:

$$\text{Var}(\hat{\zeta}) = E(\hat{\zeta} - \zeta)^2 = g_A^T(\bar{A}, \bar{B}) \Sigma_A g_A(\bar{A}, \bar{B}) + g_B^T(\bar{A}, \bar{B}) \Sigma_B g_B(\bar{A}, \bar{B}) + o(\tilde{n}^{-1})$$

و برآوردگر خطی‌سازی واریانس $\hat{\zeta}$ عبارت است از:

$$(17) \quad v_L(\hat{\zeta}) = \hat{g}_A^T(\hat{A}, \hat{B}) S_A \hat{g}_A(\hat{A}, \hat{B}) + \hat{g}_B^T(\hat{A}, \hat{B}) S_B \hat{g}_B(\hat{A}, \hat{B})$$

که \hat{g}_A و \hat{g}_B از $\hat{\beta}$ به جای β استفاده می‌کنند.

قضیه‌ی ۳. هر سه برآوردگر واریانس $v_J(\hat{\zeta})$ ، $v_{MJ}(\hat{\zeta})$ و $v_L(\hat{\zeta})$ به طور مجانبی در حالتی که اختلاف بین هر جفت از آن‌ها از مرتبه‌ی $o_p(\tilde{n}^{-1})$ باشد هم‌ارزند. (برای اثبات

[۱۴] دیده شود).

برآوردهای جک‌نایف واریانس مذکور برای حالت نمونه‌گیری طبقه‌بندی بدون جایگذاری از واحدهای نمونه‌گیری اولیه با فرض ناچیز بودن کسرهای نمونه‌گیری در دو چارچوب، مطرح شده و قابل کاربرد با طرح‌های نمونه‌گیری طبقه‌بندی چندمرحله‌ای بدون جایگذاری نیز می‌باشند. در حالتی که نمونه‌گیری با جایگذاری از واحدهای نمونه‌گیری اولیه به عمل آید نیز همین فرمول‌های واریانس صدق می‌کنند. یک حالت خاص ممکن است زمانی رخ دهد که در یکی از چارچوب‌ها مثلاً چارچوب فهرستی B ، نمونه‌گیری تصادفی ساده‌ی طبقه‌بندی انجام شده و کسرهای نمونه‌گیری ناچیز باشد. در این حالت b_{ij} برداری از مقادیر مرتبط با j امین واحد در طبقه‌ی l ام چارچوب B و $\tilde{n}_l^B = n_l^B$ عبارت از تعداد واحدهای نمونه‌گیری شده از N_l^B واحد در طبقه‌ی l ام چارچوب B می‌باشد. در صورتی که کسر نمونه‌گیری n_l^B/N_l^B ناچیز نباشد، به‌طور مشابه با روش جک‌نایف در حالت تک چارچوبی ([۱۹])، $(n_l^B - 1)/n_l^B$ با $(n_l^B - 1)/n_l^B (1 - n_l^B/N_l^B)$ جایگزین می‌شود.

۲-۳- برآورد جک‌نایف واریانس برآوردهای دوچارچوبی

تمام برآوردهای دو چارچوبی می‌توانند به‌صورت زیر بیان شوند:

$$(۱۸) \quad \hat{Y} = \sum_{t \in S_A} \tilde{w}_t^A y_t + \sum_{t \in S_B} \tilde{w}_t^B y_t$$

که از وزن‌های اصلاح‌شده‌ی \tilde{w}_t^A و \tilde{w}_t^B استفاده می‌کند. برای مثال در $\hat{Y}_{PML}(\theta)$ از وزن‌های زیر استفاده می‌شود:

$$\tilde{w}_t^A = \begin{cases} w_t^A \frac{N_A - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_a^A} & t \in a \\ \frac{w_t^A \theta \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_{ab}(\theta)} & t \in ab \end{cases}$$

$$\tilde{w}_t^B = \begin{cases} w_t^B \frac{N_B - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_b^B} & t \in b \\ \frac{w_t^B (1 - \theta) \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_{ab}(\theta)} & t \in ab \end{cases}$$

برای محاسبه‌ی برآورد جک‌نایف $\hat{Y}_{(hi)}^A$ به‌سادگی وزن‌های w_t^A با $w_{t(hi)}^A$ جایگزین می‌شود. اگر واحد t در خوشه‌ی k ام طبقه‌ی g ام در چارچوب A باشد آن‌گاه

$$w_{t(hi)}^A = \begin{cases} 0 & hi = gk \quad \text{اگر} \\ \frac{\tilde{n}_h^A}{(\tilde{n}_h^A - 1)} w_t^A & h = g \quad i \neq k \quad \text{اگر} \\ w_t^A & h \neq g \quad \text{اگر} \end{cases}$$

به‌طور مشابه $\hat{Y}_{(lj)}^B$ با استفاده از وزن‌های $w_{t(lj)}^B$ محاسبه شده و سپس برآورد جک‌نایف واریانس \hat{Y} از رابطه‌ی زیر به دست می‌آید.

$$v_j(\hat{Y}) = \sum_{h=1}^H \frac{\tilde{n}_h^A - 1}{\tilde{n}_h^A} \sum_{i=1}^{\tilde{n}_h^A} (\hat{Y}_{(hi)}^A - \hat{Y})^2 + \sum_{l=1}^L \frac{\tilde{n}_l^B - 1}{\tilde{n}_l^B} \sum_{j=1}^{\tilde{n}_l^B} (\hat{Y}_{(lj)}^B - \hat{Y})^2$$

اگر وزن‌های اصلاح شده‌ی \tilde{w}_t^A و \tilde{w}_t^B به عناصر S_A و S_B بستگی داشته باشند، برای محاسبه‌ی برآورد جک‌نایف $\hat{Y}_{(hi)}^A$ احتیاج به محاسبه‌ی $S_{(hi)}^A$ است. ماتریس $S_{(hi)}^A$ به ازای $\tilde{n}_h^A \geq 3$ می‌تواند با استفاده از یک جک‌نایف جداگانه درون هر تکرار جک‌نایف، این بار با مجموعه داده‌ای که مشاهدات i امین psu طبقه‌ی h ام چارچوب A حذف شده است، به دست آید.

۴- مطالعه‌ی شبیه‌سازی

در این بخش پس از معرفی نحوه‌ی اجرای شبیه‌سازی، به بررسی ویژگی‌های تجربی و مقایسه‌ی برآوردگرهای دوچارچوبی و همچنین مقایسه‌ی برآوردگرهای واریانس (خطی‌سازی، جک‌نایف و جک‌نایف اصلاح‌شده) با استفاده از شبیه‌سازی می‌پردازیم. یک نمونه‌ی خوشه‌ای دومرحله‌ای با \tilde{n}^A خوشه و m عنصر از هر خوشه به‌عنوان

نمونه‌ی حاصل از چارچوب A ($n_A = \tilde{n}_A \cdot m$) و یک نمونه‌ی تصادفی ساده با n_B مشاهده به‌عنوان نمونه‌ی حاصل از چارچوب B تولید می‌شود. با فرض نامتناهی بودن جامعه‌ی تحت مطالعه، N_a/N و N_b/N را با γ_a و γ_b جایگزین می‌کنیم. نمونه‌ی حاصل از چارچوب A شامل مقادیر

$$\{(y_{ij}, m_{ai}), i = 1, \dots, \tilde{n}_A \quad j = 1, \dots, m\}$$

می‌باشد که m_{ai} عبارت از تعداد عناصر نمونه‌گیری شده از i امین خوشه‌ی نمونه‌ای، متعلق به حوزه‌ی a و y_{ij} نیز مقدار مرتبط با j امین عنصر نمونه‌ای در i امین خوشه‌ی نمونه می‌باشد. نمونه‌ی حاصل از چارچوب B شامل مقادیر نمونه‌ای $\{(y_j, n_b), j = 1, \dots, n_B\}$ است که در آن n_b ، تعداد عناصر نمونه‌ای متعلق به حوزه‌ی b و y_j ، مقدار مرتبط با j امین عنصر نمونه‌ای است (جزئیات شبیه‌سازی در [۱] آمده است).

برای شبیه‌سازی، نمونه‌گیری را ۱۰۰۰۰ بار تکرار کرده و از هر نمونه، برآوردهای $\hat{Y}_{PML}(\hat{\theta}_P)$ و $\hat{Y}_{FB}(\hat{\beta}_{FB})$ ، $\hat{Y}_H(\hat{\theta}_H)$ و $\hat{Y}_{SF}(\hat{\beta}_{SY})$ و \hat{Y}_{SFrake} و \hat{Y}_{SF} چارچوبی تک برآوردهای تک چارچوبی $\hat{Y}_{PML}(\hat{\theta}_P)$ و $\hat{Y}_{FB}(\hat{\beta}_{FB})$ ، $\hat{Y}_H(\hat{\theta}_H)$ و $\hat{Y}_{SF}(\hat{\beta}_{SY})$ با استفاده از مقادیر بهینه‌ی $\hat{\theta}_P$ ، $\hat{\beta}_{FB}$ ، $\hat{\theta}_H$ و $\hat{\beta}_{SY}$ و به ازای مقادیر مختلف γ_a ، γ_b ، n_B و \tilde{n}_A محاسبه می‌کنیم و سپس میانگین توان دوم خطای تجربی هر برآوردها (متوسط توان دوم انحراف برآوردها از مقدار واقعی آن) را به‌صورت زیر محاسبه می‌کنیم:

$$EMSE = \frac{1}{R} \sum_{r=1}^R (\hat{Y}_r - \hat{Y})^2$$

که \hat{Y}_r عبارت است از برآورد \hat{Y} برای r امین تکرار شبیه‌سازی، $\hat{Y} = \frac{1}{R} \sum_{r=1}^R \hat{Y}_r$ و R ، تعداد تکرارهای شبیه‌سازی است.

برای هر ۶ برآوردها Y ، سه برآوردها خطی‌سازی واریانس (L)، جک‌نایف (J) و جک‌نایف اصلاح‌شده (MJ) را محاسبه می‌کنیم. در محاسبه‌ی واریانس جک‌نایف به‌دلیل این‌که از مقادیر برآوردهای پارامترها استفاده می‌شود و این مقادیر، خود به ماتریس‌های S_A و S_B وابسته‌اند، $S_{(hi)}^A$ و $S_{(ij)}^B$ با استفاده از یک جک‌نایف جداگانه درون هر تکرار

جک‌نایف محاسبه می‌شوند. درصد اریبی نسبی (RB) هر سه برآوردگر واریانس را با استفاده از رابطه‌ی زیر محاسبه می‌کنیم:

$$RB = \frac{100(EV - EMSE)}{EMSE}$$

که EV عبارت است از میانگین 10000 برآورد واریانس در هر روش. همچنین با استفاده از بازه‌ی اطمینان $(1 - \alpha)$ درصد توزیع نرمال:

$$\text{برآورد} \pm z_{(1-\frac{\alpha}{2})} SE \text{ (برآورد)}$$

و نیز صدک توزیع t با $(\tilde{n}_A - 1)$ درجه‌ی آزادی، احتمال پوشش تجربی برای روش‌های خطی‌سازی، جک‌نایف اصلاح‌شده و جک‌نایف محاسبه شده است. در این شبیه‌سازی، تمام محاسبات با استفاده از نرم افزار $\text{SPSS} 2000$ صورت گرفته است.

نتایج شبیه‌سازی در جدول‌های ۱ و ۲ (در جدول ۱، $\frac{n_A}{n_B} = 3$ و در جدول ۲، $\frac{n_A}{n_B} = 6$)

نشان می‌دهد که برآوردگر ماکسیمم درست‌نمایی‌نما، همواره کاهش‌ی کوچک در $EMSE$ نسبت به برآوردگرهای هارتلی و فولر-بورمیستر در اثر برآورد $\hat{\theta}_H$ و $\hat{\beta}_{FB}$ نشان می‌دهد. در مقایسه با برآوردگرهای دو چارچوبی، برآوردگر تک چارچوبی، افزایش قابل ملاحظه‌ای را در $EMSE$ وقتی $\gamma_a \leq \gamma_b$ یا n_A خیلی بزرگ‌تر از n_B باشد نشان می‌دهد. اریبی نسبی هر سه برآوردگر واریانس با افزایش اندازه‌ی نمونه میل به کاهش دارد. با این حال برای نمونه‌های کوچک‌تر، روش‌های خطی‌سازی و جک‌نایف اصلاح‌شده، $EMSE$ را کم برآورد می‌کنند. (نتایج به ازای مقادیر دیگر پارامترهای γ_a و γ_b و دیگر اندازه‌های نمونه از دو چارچوب به غیر از آنچه در جدول‌های ۱ و ۲ است در [۱] آمده است).

نتایج حاصل از بازه‌ی اطمینان نشان می‌دهد که روش جک‌نایف کامل به وضوح احتمال پوشش بالاتری نسبت به دو روش دیگر دارد در حالی که دو روش دیگر به وضوح احتمال پوشش کم‌تری دارند که البته با افزایش اندازه‌ی نمونه، احتمال پوشش افزایش می‌یابد.

۵- بحث و نتیجه‌گیری

نتایج حاصل از بررسی خواص برآوردگرهای دوچارچوبی نشان می‌دهد که برآوردگر ماکسیم درست‌نمایی‌نما به دلیل این که از وزن‌های یکسانی برای همه‌ی متغیرها استفاده می‌کند نسبت به برآوردگر هارتلی و فولر- بورمیستر در آمارگیری‌هایی با متغیرهای متعدد ارجحیت دارد. در صورتی که $\frac{N_a}{N} \cong \frac{N_b}{N}$ و $(\gamma_a \cong \gamma_b)$ و دو طرح نمونه‌گیری مشابه باشند، برآوردگر تک‌چارچوبی با تصحیح نسبتی، نسبت به دیگر برآوردگرها از کارایی بالاتری برخوردار است. برآوردگر جک‌نایف دارای اربیبی کوچک‌تری نسبت به برآوردگر واریانس خطی‌سازی است. اگر چه احتمال‌های پوشش تجربی برای هر سه روش تقریباً مشابه هستند، اما وقتی که صدک توزیع t در محاسبه‌ی بازه‌ی اطمینان استفاده می‌شود، روش جک‌نایف احتمال پوشش خیلی بالاتری را نسبت به دو برآوردگر دیگر نشان می‌دهد. جک‌نایف معمولاً با توابع ناخطی نظیر نسبت مجموع دو جامعه به کار می‌رود. در این حالت مشتق‌های جزئی که در محاسبه‌ی برآورد خطی‌سازی چنین مقادیر ناخطی استفاده می‌شود، در آمارگیری‌های دوچارچوبی پیچیده‌تر از برآوردگر تک‌چارچوبی است. در صورتی که با استفاده از جک‌نایف می‌توان از این محاسبات پیچیده اجتناب کرد. روش‌های بازنمونه‌گیری دیگر برآورد واریانس نظیر تکرار مکرر متعادل $(BRR)^{\wedge}$ و خودگردان^۹ نیز می‌توانند به موازات جک‌نایف در این‌گونه طرح‌ها مورد استفاده قرار گیرند. مزیت این روش‌ها این است که برخلاف جک‌نایف با توابع ناهموار نظیر میانه نیز قابل استفاده هستند.

جدول ۲- میانگین توأم دوم خطای براوردگرهای چادچوب دوکان، درصد آریبی نسبی و احتمال پوشش سه براوردگر واریانس برای $\frac{n_A}{n_B} = 6$

احتمال پوشش										
احتمال پوشش					درصد آریبی نسبی					
نرمال					نرمال					
J	MJ	L	J	MJ	L	J	MJ	L	EMSE	برآوردگر
$\tilde{n}_A = 2, m = 3, (n = \tilde{n}_A \times m), n_B = 10, \gamma_a = 0.1, \gamma_b = 0.2$										
0.952	0.937	0.937	0.85	0.922	0.922	-0.63	-11.06	-11.36	2.37	H
0.948	0.920	0.902	0.918	0.908	0.885	1.8	-16.66	-17.66	2.43	FB
0.947	0.949	0.946	0.943	0.939	0.930	-1.10	-7.81	-8.2	2.3	PML
0.972	0.937	0.938	0.958	0.936	0.934	-7.99	-7.99	-7.99	5.45	SF
0.952	0.932	0.914	0.947	0.921	0.921	0.08	-12.05	-12.69	2.28	SF _{reg}
0.957	0.948	0.940	0.945	0.905	0.95	-3.07	-0.7	-1.7	2.32	SF _{rate}

توضیحات

1. Delta Method
 2. Multiple frame surveys
 3. Jackknife
 4. Pseudo-maximum likelihood estimator
 5. Single frame estimator
۶. برای دو دنباله از اعداد $\{a_n\}$ و $\{b_n\}$ ، اگر $a_n = O(b_n)$ ، اگر به ازای هر n یک c ثابت $|a_n| \leq c|b_n|$ ؛ و $a_n = o(b_n)$ اگر وقتی $n \rightarrow \infty$ ، $\frac{a_n}{b_n} \rightarrow 0$.
۷. برای دو دنباله از متغیرهای تصادفی $\{X_n\}$ و $\{Y_n\}$ ، اگر برای هر $\varepsilon > 0$ ، یک ثابت $C_\varepsilon > 0$ وجود داشته باشد به طوری که $p\{|X_n| \geq C_\varepsilon|Y_n|\}$ برای هر n ؛ $X_n = o_p(Y_n)$ اگر وقتی $n \rightarrow \infty$ ، $\frac{X_n}{Y_n} \xrightarrow{p} 0$.
8. Balanced repeated replication (BRR)
 9. Bootstrap

مرجع‌ها

- [۱] نورینی، مرجان (۱۳۸۲). طرح‌های نمونه‌گیری چندچارچوبی. پایان‌نامه‌ی کارشناسی ارشد. دانشگاه صنعتی اصفهان، اصفهان.
- [2] Bankier, M.D. (1986). Estimators based on several stratified samples with application to multiple frame surveys, *Journal of the American Statistical Association*, **81**, 1074-1079.
- [3] Ford, B.L. and Bosecker, R.R. (1979). Multiple frame estimation with stratified overlap domain, in *Proceedings of the Social Statistics Section, American Statistical Association*, 219-224.
- [4] Fuller, W.A. and Burmeister, L.F. (1972). Estimators for samples selected from two overlapping frames, in *Proceedings of the Social Statistics Section, American Statistical Association*, 245-249.
- [5] Hartley, H.O. (1962). Multiple frame surveys, in *Proceedings of the Social Statistics Section, American Statistical Association*, 203-206.
- [6] Hartely, H.O. (1974). Multiple frame methodology and selected

- application, *Sankhya, Ser. C*, **36**, 203–206.
- [7] Kalton, G. and Anderson, D.W. (1986). Sampling rare populations, *Journal of the Royal statistical Society, Ser. A*, **149**, 65–82.
- [8] Lohr, S. and Rao, J.N.K. (1997). Jackknife variance estimation in multiple frame surveys. *Proceedings of the Section on Survey Research Methods, American Statistical Association*, 552–557
- [9] Lund, R.E. (1968). Estimators in multiple frame surveys, in *Proceedings of the Social Statistics Section, American Statistical Association*, 282–288.
- [10] Quenouille, M.H. (1949). Approximation tests of correlation in time Series, *Journal of Royal Stat. Soc.*, **B 11**, 68–84
- [11] Quenouille, M.H. (1956). Notes on bias in estimation, *Bimetrika*, **43**, 353–360,
- [12] Rao, J.N.K. and Wu, C.F J. (1985). Inference from stratified samples: Second-order analysis of three methods for nonlinear statistics, *Journal of the American Statistical Association*, **80**, 620–630.
- [13] Rao, J.N.K. and Lohr, S.L. (1997). *Jackknife Variance Estimation in Dual Frame Surveys*, Technical Report, Department of Mathematics and Statistics, Carleton University.
- [14] Rao, J.N.K. and Lohr, S.L. (2000). Inference from dual frame survey, *Journal of the American Statistical Association*, **95**, 271–280.
- [15] Sarndal, C.E., Swensson, B. and Wretman, J. (1992). *Model Assisted Survey Sampling*, Springer-Verlag, New York.
- [16] Skinner, C.J. (1991). On the efficiency of raking ratio estimators for multiple frame surveys, *Journal of the American Statistical Association*, **86**, 779–784.
- [17] Skinner, C.J. and Rao, j.N.K. (1996). Estimation in dual frame survey with complex design, *Journal of the American Statistical Association*, **91**, 349–356.
- [18] Tukey, J.W. (1958). Bias and Confidence in Note Quite Large Samples

(Abstract), *Annals of Mathematical Statistics*, **29**, 614.

- [19] Vogel, F.A. (1975). Survey with overlapping frame-problems in application, in Proceedings of the Social Statistics Section, *American Statistical Association*, 694-699.
- [20] Wolter, K.M. (1985). *Introduction to Variance Estimation*, Springer-Verlag, New York.

مرجان نورینی

فوق لیسانس آمار

تهران، خیابان فاطمی، نبش خیابان رهی معیری، مرکز آمار ایران.

رایانشانی: Mar_noorini@yahoo.com