

محاسبه واریانس برای داده‌های حاصل از نمونه‌گیری چند مرحله‌ای و برآورد نسبی

آرمان بیداربخت‌نیا

مرکز آمار ایران

چکیده. در این مقاله روش‌هایی برای برآورد واریانس مورد بررسی قرار می‌گیرد که اثر روش‌های پیچیده نمونه‌گیری و همچنین استفاده از برآوردگرهای نسبی و بساطقه‌بندی را منعکس می‌کنند. نخست برآورد واریانس در طرح‌های مختلف نمونه‌گیری بررسی می‌شود و سپس با استفاده از روش خطی‌سازی سری تیلور، برآورد واریانس برای برآوردگر نسبی محاسبه می‌شود. در نهایت با استفاده از یک مثال، برآورد واریانس یک برآوردگر پیچیده به‌گونه‌ای بیان می‌شود که اثر روش نمونه‌گیری، روش برآورد و همچنین وزن‌دهی را نیز منعکس می‌کند.

۱- مقدمه

ممکن است نمونه‌ای با استفاده از روش‌های نمونه‌گیری چند مرحله‌ای انتخاب شده و سپس عناصر آن به‌منظور برآورد مقادیر کل جامعه وزن‌دهی شده باشند. به این ترتیب برآورد نسبی پس‌طبقه‌بندی شده تورمی ایجاد می‌کند که نه تنها اندازه بلکه ترکیب آن مثل سن-جنس-نژاد و مناطق سکونت را نیز منعکس می‌کند. سپس این داده‌های وزن‌دهی شده برای برآورد میانگین، نسبت و دیگر پارامترهای مورد علاقه به کار برده می‌شوند. هدف اصلی مقاله این است که نشان دهد چگونه برآوردگرهای واریانس، مراحل واقعی روش‌های نمونه‌گیری و برآورد نسبی را منعکس می‌کنند.

بعضی از واریانس‌ها را می‌توان با ترکیب دو روش شناخته‌شده خطی‌سازی برآوردهای نسبتی (وودراف، ۱۹۷۱) و بروردگر واریانس تعمیم‌یافته برای نمونه‌گیری چند مرحله‌ای (کندال و استورات، ۱۹۶۸) تقریب زد. با این ترکیب، می‌توان ویژگی‌های این دو روش را در برآوردهای واریانس انعکاس داد.

روش‌های نخست، برآورد نسبتی را با استفاده از بسط سری‌های تیلور خطی نموده و سهم متغیر از بسط خطی برای برآوردهای نسبتی را حفظ می‌کند. سپس واریانس نسبت با واریانس سهم متغیر از نسبت‌های خطی شده، تقریب زده می‌شود.

روش‌های دوم، حالت تعمیم‌یافته واریانس برای داده‌های جمع‌بندی‌شده را برای هر طرحی که در انتخاب نمونه به کار رفته باشد، شامل می‌شوند.

بعد از این روش‌ها دو دسته علائم مجموع می‌تواند وجود داشته باشد که دسته اول از برآوردهای نسبتی و دسته دوم از طرح‌های نمونه‌گیری ناشی می‌شود. واریانس تعمیم‌یافته می‌تواند با تعویض علائم مجموع، انتقال آن‌ها به جلو برای نمونه‌گیری و جمع بستن آن‌ها برای برآورد نسبتی به دست آید. به این ترتیب فقط علائم مجموع ناشی از روش‌های نمونه‌گیری باقی می‌ماند. سپس می‌توانیم فرمول واریانسی که قبلاً تهیه شده است را برای این نتیجه نهایی به کار ببریم.

این روش‌ها می‌تواند برای برآوردهای واریانس داده‌هایی به کار رود که توسط مراکز آمارهای بهداشتی و دیگر مراکز دولتی گردآوری شده باشد که معمولاً از برآوردها و طرح‌های نمونه‌گیری پیچیده بهره می‌گیرند و ممکن است که هنوز به فکر هر دو ویژگی برآوردهای واریانس نباشند. هیدبروگلو و رائو (۱۹۸۳) و شاه (۱۹۸۱) این روش‌ها را به ترتیب برای تحلیل داده‌های آمارگیری بهداشتی کانادا و برنامه خطاهای استاندارد برای داده‌های آمارگیری به کار برده‌اند؛ فرمول‌های واریانس در مورد نخست برای نمونه‌گیری با جایگذاری با احتمال‌های برابر به کار رفته است، در حالی که برای مورد اخیر نمونه‌گیری دو مرحله‌ای بدون جایگذاری با احتمال‌های برابر استفاده شده است. در عمل برای نمونه‌گیری چند مرحله‌ای به ندرت در همه مراحل از یک روش نمونه‌گیری استفاده می‌شود. برای مثال، ممکن است نمونه‌ها در مرحله اول با جایگذاری و با احتمال‌های متناسب با اندازه جامعه (pps) و در مرحله دوم بدون جایگذاری و با احتمال‌های برابر

انتخاب شود.

بخش ۲ به معرفی برخی نمادهایی که در بخش‌های بعد به کار می‌رود، پرداخته است. بخش ۳ برآوردهای واریانس تعمیم‌یافته برای سرجمع‌ها، با استفاده از نمونه‌گیری چند مرحله‌ای را ارائه می‌کند. در بخش ۴، برآوردهای نسبتی پس طبقه‌بندی شده، خطی شده‌اند و فقط سهم متغیر حفظ شده است. سپس واریانس نسبت‌ها را با آن سهم متغیر نسبت‌های خطی شده تقریب می‌زنیم. نهایتاً یک مثال و چند دستورات عمل در بخش ۵ آمده است.

۲- علائم

فرض کنید که جامعه به L طبقه مستقل طبقه‌بندی شده باشد که با $s = 1, \dots, L$ نمادگذاری می‌شود و اعضای طبقه s در N_{1s} واحد نمونه‌گیری مقدماتی (PSU) گروه‌بندی شده باشند که با $i = 1, \dots, N_{1s}$ نمادگذاری می‌شود و i امین PSU شامل N_{Tsi} عضو است که با $z = 1, \dots, N_{Tsi}$ نشان داده می‌شود. علائم متناظر برای نمونه با حرف کوچک n و با زیرنویس‌های مشابه، همان‌گونه که در جدول (۱) مشاهده می‌کنیم، نشان داده می‌شود. از آن‌جا که واریانس برای سرجمع‌های این طبقات جمع‌پذیر است، برای نشان دادن واریانس حاصل از یک طبقه، از زیرنویس s برای طبقات صرف نظر می‌کنیم. جدول (۲) واریانس‌های سرجمع‌ها را به همراه نوع سرجمع و طرح نمونه‌گیری نشان می‌دهد.

برخی از فرمول‌ها برای این واریانس‌ها در بخش ۳ مورد بحث قرار می‌گیرد. واریانس نسبت‌ها در جدول (۱) می‌تواند خطی شده و به این ترتیب با سرجمع‌ها در دسته‌های مشابه قرار گیرد.

جدول ۱- علائم برای داده‌های نمونه‌ای، در حالی که خوشه‌بندی در دو مرحله انجام شده و در هر طبقه نمونه‌گیری سه مرحله‌ای اجرا شده باشد

نمونه	جامعه	
n_1	N_1	واحدهای مرحله اول
n_{r_i}	N_{r_i}	واحدهای مرحله دوم
$n_{r_{ij}}$	$N_{r_{ij}}$	واحدهای مرحله سوم
$i = 1, \dots, n_1$	$i = 1, \dots, N_1$	شاخص مرحله اول
$j = 1, \dots, n_{r_i}$	$j = 1, \dots, N_{r_i}$	شاخص مرحله دوم
$k = 1, \dots, n_{r_{ij}}$	$k = 1, \dots, N_{r_{ij}}$	شاخص مرحله سوم
$h = 1, \dots, q$	$h = 1, \dots, q$	شاخص سلول‌ها
$n = \sum_i \sum_j n_{r_i}$	$N = \sum_i \sum_j N_{r_{ij}}$	مقادیر کل:
$y_{hr} = \sum_i \sum_j \sum_k y_{hijk}$	$Y_{hr} = \sum_i \sum_j \sum_k y_{hijk}$	سرمجموع‌ها برای سلول
$x_{hr} = \sum_i \sum_j \sum_k x_{hijk}$	$X_{hr} = \sum_i \sum_j \sum_k x_{hijk}$	
$\frac{y_{hr}}{n}$	$\frac{Y_{hr}}{N}$	نسبت در سلول:
$r_{hr} = x_{hr} / y_{hr}$	$R_{hr} = X_{hr} / Y_{hr}$	نسبت:

x_{hijk} و y_{hijk} متغیرهایی برای صفات x و y به ترتیب هستند.

۳- واریانس

نمونه‌گیری می‌تواند با احتمال برابر یا نابرابر و یا با احتمال‌های متناسب با اندازه (pps)، با جایگذاری و یا بدون جایگذاری و با طرح‌های متقارن یا نامتقارن انجام شود. در هر مرحله، ممکن است که هر ترکیبی از این گزینه‌ها را در نظر داشته باشیم.

جدول ۲- واریانس‌ها با انواع طرح و سرجمع

سرجمع			انواع طرح
نمونه‌گیری سه مرحله‌ای ^۶	نمونه‌گیری دو مرحله‌ای ^۷	نمونه‌گیری یک مرحله‌ای ^۱	
$\text{var}(y_{h\tau})$	$\text{var}(y_{h\tau})$	$\text{var}(y_{h1})$	${}^{\tau}WR$ احتمال‌های نابرابر
$\text{var}(y_{h\tau})$	$\text{var}(y_{h\tau})$	$\text{var}(y_{h1})$	${}^{\tau}W0$
$\text{var}(y_{h\tau})$	$\text{var}(y_{h\tau})$	$\text{var}(y_{h1})$	${}^{\tau}WR$ احتمال‌های برابر
$\text{var}(y_{h\tau})$	$\text{var}(y_{h\tau})$	$\text{var}(y_{h1})$	${}^{\tau}W0$
$\text{var}(y_{h\tau})$	$\text{var}(y_{h\tau})$		ترکیب‌ها ^۵

$${}^1 y_{h1} = \sum_i y_{hi} \quad ; \quad {}^{\tau} y_{h\tau} = \sum_i \sum_j y_{hij} \quad ; \quad \text{بدون جایگذاری؛} \quad {}^{\tau} \text{ یا جایگذاری؛}$$

$${}^5 \text{ ترکیبی از نمونه‌گیری احتمالاتی با احتمال‌های برابر و نابرابر؛} \quad {}^6 y_{h\tau} = \sum_i \sum_j \sum_k y_{hijk}$$

یک فرمول واریانس تعمیم‌یافته برای هر برآورد $\hat{\theta}_h$ در h امین سلول، بر اساس احتمال‌های انتخاب کاملاً دلخواه ارائه می‌کنیم. به این ترتیب واریانس کل عبارت از مجموع واریانس‌ها برای تمام طبقات است.

علامت E برای عملگر مقدار مورد انتظار، var برای واریانس و vâr برای برآورد ناریب var به کار رفته است. می‌توان نوشت:

$$(1) \quad \text{var}(\hat{\theta}_h) = \text{var}(E(\hat{\theta}_h)) + E(\text{var}(\hat{\theta}_h))$$

که « > 1 » علامتی برای ارائه تمام مراحل نمونه‌گیری پس از مرحله نخست است.

اگر $\hat{\theta}_h = y_{h1}$ ، تعریف شده در جدول (۱) باشد، برآوردگر ناریب آن می‌تواند به شکل

زیر نوشته شود:

$$(۲) \quad \text{vâr}(\hat{\theta}_h) = \text{vâr}(\hat{\theta}_h) + \sum_i^{n_i} \pi_i^{(1)} \text{vâr}(y_{hi})$$

که $\pi_i^{(1)}$ عبارت است از احتمال این که واحد i ام در n_i واحد نمونه‌گیری مقدماتی (PSU) باشد.

رابطه (۱) ممکن است به سه مولفه تقسیم شود:

$$(۳) \quad \text{var}(\hat{\theta}_h) = \text{var} E E(\hat{\theta}_h) + E \text{var} E(\hat{\theta}_h) + E E \text{var}(\hat{\theta}_h).$$

اگر $y_{hi} = \sum_{j=1}^{n_{ij}} y_{hij}$ ، در رابطه (۲) جایگزین شود، برآورد نارایب (۳) را می‌توان به شکل زیر

نوشت:

$$(۴) \quad \text{vâr}(\hat{\theta}_h) = \text{vâr}(\hat{\theta}_h) + \sum_i^{n_i} \pi_i^{(1)} \text{vâr}(y_{hi}) + \sum_i^{n_i} \pi_i^{(1)} \sum_j^{n_{ij}} \pi_{ij}^{(2)} \text{vâr}(y_{hij})$$

که $\pi_{ij}^{(2)}$ عبارت است از احتمال این که واحد j ام از مرحله دوم در i امین واحد انتخاب شده در مرحله اول باشد. این رابطه برای مراحل بعدی نمونه‌گیری نیز قابل تعمیم است.

فرمول‌های بالا را می‌توان به این ترتیب خلاصه نمود: یک برآوردگر نارایب برای واریانس در نمونه‌گیری چند مرحله‌ای، وقتی که مرحله اول نمونه‌گیری بدون جایگذاری باشد، از مجموع دو مولفه حاصل شده است. مولفه اول واریانس را در حالتی که فقط نمونه‌گیری مرحله اول اجرا شده برآورد می‌کند. مولفه دوم عبارت از مجموع موزون برآوردها درون واحدهای انتخاب شده در مرحله اول، برای واریانس حاصل از مراحل دیگر نمونه‌گیری است (واحدهای مرحله اول ثابت فرض می‌شوند)؛ وزن‌ها عبارت از احتمال‌های انتخاب این واحدهای مرحله اول می‌باشد. (دوربین، ۱۹۵۳).

اگر نمونه‌گیری در مرحله اول به صورت جایگذاری انجام شده باشد، با توجه به این که $\pi_i^{(1)} \rightarrow 0$ ، فقط جمله اول در (۴) باقی می‌ماند. در این مورد، برای یک نمونه‌گیری چند مرحله‌ای با هر تعداد مرحله نمونه‌گیری، وقتی که در مرحله اول برای هر انتخاب احتمال‌های نامساوی به کار برده می‌شود، در حالی که مراحل دیگر دلخواه بوده و در واحدهای متمایز انتخاب شده در مرحله اول استخراج به صورت مستقل انجام می‌شود،

برآورد واریانس‌ها ساده است.

واریانس‌های مربوط به چند وضعیت متفاوت برای نمونه‌گیری را ملاحظه کنید.

الف) $\text{var}(y_{hi})$ برای حالت بدون جایگذاری با احتمال نابرابر در یک مرحله:

فرض کنید p_i احتمال انتخاب شدن i امین فرد در r امین انتخاب باشد و

$$\sum_i^{N_i} p_i = 1, \quad \pi_i = \sum_r p_i, \quad \pi_{i' i''} = \sum_{r \neq s} p_i p_{i''}$$

کندال و استوارت (۱۹۶۸ صفحه ۱۷۲) نشان داده‌اند که

$$(5) \quad \text{var}(y_{hi}) = \sum_{i=1}^{N_i} \pi_i (1 - \pi_i) y_{hi}^2 + \sum_{i \neq i'} \sum_{i''} (\pi_{i' i''} - \pi_i \pi_{i'}) y_{hi} y_{hi''}$$

با استفاده از:

$$E\left(\sum_{i \neq i'} \sum_{i''} g(y_i, y_{i'})\right) = \sum_{i \neq i'} \sum_{i''} \pi_{i' i''} g(y_i, y_{i'}) \quad \text{و} \quad E\left(\sum_i g(y_i)\right) = \sum_i \pi_i g(y_i)$$

برای هر تابع g از مشاهدات، برآورد ناریب (۵) به صورت زیر داده شده است:

$$(6) \quad \hat{\text{var}}(y_{hi}) = \frac{1}{2} \sum_{i \neq i'} \sum_{i''} \frac{(\pi_i \pi_{i'} - \pi_{i' i''})}{\pi_{i' i''}} (y_{hi} - y_{hi''})^2,$$

برای نمونه‌گیری یک مرحله‌ای، در فرمول عمومی (۶)،

$$\hat{\text{var}}(y_{hi}) = \hat{\text{var}}(y_{hi}), \quad (7)$$

ب) $\text{var}(y_{h\tau})$ برای حالت بدون جایگذاری با احتمال نابرابر در نمونه‌گیری دو مرحله‌ای:

$$\hat{\text{var}}(y_{h\tau}) = \hat{\text{var}}(y_{h\tau}) + \sum_i \pi_i^{(1)} \hat{\text{var}}(y_{hi}),$$

که جمله اول با استفاده از (۶) معلوم است و جمله دوم عبارت از مجموع موزون واریانس‌ها برای مراحل دوم در واحدهای انتخاب شده در مرحله اول است.

پ) $\text{var}(y_{h1})$ برای حالت بدون جایگذاری با احتمال برابر در نمونه‌گیری یک مرحله‌ای:

داریم $\pi_i = \frac{n_i}{N_1} = F_1$ و $\pi_{i' } = \frac{n_i}{N_1} \frac{(n_i - 1)}{(N_1 - 1)}$. با استفاده از این عبارتها، می‌توان (۶) را به صورت زیر نوشت:

$$(۷) \quad \hat{\text{var}}(y_{h1}) = (1 - F_1) \frac{n_1}{n_1 - 1} \sum_i^{n_1} (y_{hi} - \bar{y}_h)^2$$

که \bar{y} میانگین y_i ها است.

ت) $\text{var}(y_{h1})$ برای حالت با جایگذاری با احتمال نابرابر:

حال باید حالت $i = i'$ را برای $\pi_{i'}$ در نظر بگیریم، اما جمله π_i و $\pi_{i'}$ در مجموع دوگانه هنوز هم باید پسوند‌های متفاوتی داشته باشد. رابطه (۶) برای این نمونه‌گیری نیز برقرار است.

ث) $\text{var}(y_{h1})$ برای حالت با جایگذاری با احتمال برابر:

در این مورد تئوری ساده‌تر می‌شود، یعنی $\pi_i = np_i$ و $\pi_{i' } = n(n-1) p_i p_{i'}$ که عبارت است از احتمال این که هر کدام از انتخاب‌های نمونه‌گیری با جایگذاری شامل عضو i ام باشد و تحت این تعاریف می‌توان (۶) را به شکل زیر نوشت:

$$(۸) \quad \text{var}(y_{h1}) = \frac{n_1}{2} \sum_i^{n_1} \sum_j^j p_i p_j (y_i - y_j)^2,$$

از آن‌جا که شرایطی مشابه با (۶) برقرار است، با قرار دادن $i = i'$ ، برآوردگر نارایب (۸) را می‌توان به شکل زیر نوشت:

$$(۹) \quad \hat{\text{var}}(y_{h1}) = \frac{n_1}{(n_1 - 1)} \sum_i^{n_1} (y_{hi} - \bar{y}_h)^2,$$

که تفاوت آن با (۷) فقط در $(1 - F_1)$ ، یعنی فاکتور مربوط به بدون جایگذاری بودن طرح

است. بخشی از هدف طرح نمونه (یعنی انتخاب π_{ii} و سپس π_i) کاهش واریانس برآوردگر تا حد ممکن است. ما می‌توانیم مجموعه‌هایی از π_{ii} بیابیم که در تولید واریانس کوچک برای تمام برآوردهایی که ممکن است به کار ببریم، مؤثر باشند. برور (۱۹۶۳) مقادیر π_{ii} و π_r را در حالتی که دو واحد نمونه داشته باشیم ($n=2$) معرفی کرد که ویژگی‌های مطلوب واریانس کوچک در (ح) و $(\pi_i \pi_r - \pi_{ii}^2) > 0$ که در (۶) نشان داده شد را دارا است.

(ج) $\text{var}(y_{hr})$ برای حالت بدون جایگذاری با احتمال برابر:

$$(10) \quad y_{hr} = \sum_i^n \sum_j^{n_{ri}} \sum_k^{n_{rj}} y_{hijk}$$

فرض کنید که نمونه‌گیری با احتمال برابر و بدون جایگذاری در هر یک از سه مرحله نمونه‌گیری در طبقه انجام می‌شود. همان‌گونه که در قسمت (پ) انجام شد، چنین احتمال‌هایی را در فرمول عمومی (۴) جایگزین می‌کنیم. می‌توان نشان داد که

$$(11) \quad \text{var}(y_{hr}) = n_1 \sigma_{y_n}^2 (1 - F_1) + F_1 \sum_i^{N_1} n_{r1} \sigma_{y_{hi}}^2 (1 - F_{r1}) + F_1 \sum_i^{N_1} F_{r1} \sum_j^{N_{r1}} n_{rj} \sigma_{y_{hij}}^2 (1 - F_{rj})$$

که $F_1 = n_1 / N_1$ ، $F_{r1} = n_{r1} / N_{r1}$ و $F_{rj} = n_{rj} / N_{rj}$ کسرهای نمونه‌گیری هستند.

$$\begin{aligned} \sigma_{y_n}^2 &= \frac{1}{N_1 - 1} \sum_i^{N_1} (y_{hi} - \frac{1}{N_1} \sum_i^{N_1} y_{hi})^2 \\ \sigma_{y_{hi}}^2 &= \frac{1}{N_{r1} - 1} \sum_j^{N_{r1}} (y_{hij}^* - \frac{1}{N_{r1}} \sum_j^{N_{r1}} y_{hij}^*)^2 \\ \sigma_{y_{hij}}^2 &= \frac{1}{N_{rj} - 1} \sum_k^{N_{rj}} (y_{hijk} - \frac{1}{N_{rj}} \sum_k^{N_{rj}} y_{hijk})^2 \\ y_{hi} &= \frac{n_{ri}}{N_{r1}} \sum_j^{N_{r1}} y_{hij}^* , \quad y_{hij}^* = \frac{n_{rj}}{N_{rj}} \sum_k^{N_{rj}} y_{hijk} \end{aligned}$$

برای داده‌های متقارن، یعنی $n_1 n_{r1} n_{rj} = n_1 n_r n_r = n$ ، برآورد ناریسب (۱۱) به

صورت زیر داده شده است:

$$(۱۲) \text{var}(y_{h_1}) = n^2 \left(\frac{s_i^2}{n_1} (1 - F_1) + \frac{n_1}{N_1} \frac{s_r^2}{n_1 n_r} (1 - F_r) + \frac{n_1}{N_1} \frac{n_r}{N_r} \frac{s_r^2}{n_1 n_r n_r} (1 - F_r) \right)$$

در (۱۲،۳) همه عبارات‌ها بعد از عبارت اول در کسرهای نمونه‌گیری مرحله قبل، $(n_1 / N_1)(n_r / N_r) \dots$ ضرب شده و s^2 نیز جایگزین σ^2 شده است. توجه کنید که اگر n_1 / N_1 قابل چشم‌پوشی باشد، تمام عبارات‌های دیگر پس از عبارت نخست نیز قابل چشم‌پوشی هستند.

(چ) $\text{var}(y_{h_1})$ برای حالت بدون جایگذاری با احتمال برابر در نمونه‌گیری دو مرحله‌ای:

نتایج دو مرحله‌ای، با قرار دادن $N_1 = n_1 = 1$ در (۱۱) پس از تغییرات مناسب در زیرنویس‌ها حاصل می‌شود.

(ح) $\text{var}(y_{h_2})$ برای حالت با جایگذاری با احتمال متناسب با اندازه جامعه (pps) برای دو مرحله نخست و بدون جایگذاری با احتمال برابر برای مرحله سوم:

استفاده از نمونه‌گیری با احتمال‌های برابر در نمونه‌گیری چند مرحله‌ای به دلیل بزرگ شدن واریانس، به ندرت صورت می‌گیرد. وقتی که واحدها از نظر اندازه به‌طور قابل ملاحظه‌ای متفاوت باشند، نمونه‌گیری با احتمال برابر منجر به بزرگ شدن واریانس‌ها می‌شود. در مقابل، در حالت متقارن وقتی که تمام واحدها در هر مرحله‌ای اندازه یکسان داشته باشند، این مشکل پیش نمی‌آید. بنا بر این ناچار هستیم که برای کاهش واریانس نمونه‌گیری، طرح نمونه‌گیری دیگری را جستجو کنیم.

با تنوع دادن به احتمال‌ها در هر مرحله، ممکن است به هدف بالا دست یابیم. اگر

احتمال کلی انتخاب یک عضو منفرد در یک نمونه‌گیری چند مرحله‌ای $\frac{n}{N}$ باشد، آن

نمونه‌گیری را به دلیل این که اعضا نمونه به‌طور مساوی وزن‌دهی شده‌اند، خود وزن گویند. بنا بر این، واریانس نمونه می‌تواند برای برخی برآوردها کاهش یابد. یک راه ساده برای دست‌یابی به طرح نمونه‌گیری pps خود وزن عبارت است از انتخاب n_1 واحد نمونه‌گیری مقدماتی (PSU) با احتمال‌های $p_i^{(1)}$ در هر انتخاب، n_{pi} واحد مرحله دوم از

هر یک از n_1 واحد نمونه‌گیری مقدماتی (PSU) با احتمال‌های $p_{ij}^{(r)}$ و $n_{\tau ij}$ واحد مرحله سوم با احتمال‌های $p_{ijk}^{(r)}$ در هر انتخاب، که $P_i^{(1)} = \frac{N_{\tau i+}}{N}$ و $P_{ij}^{(r)} = \frac{N_{\tau ij}}{N_{\tau i+}}$

و $P_{ijk}^{(r)} = \frac{1}{N_{\tau ij}}$ زیرنویس «+» به مفهوم جمع‌بستن روی زیرنویس مورد نظر است.

لذا طرح نمونه‌گیری pps یک شرط لازم برای کاهش واریانس نمونه فراهم می‌کند، اگر $n_1 n_{\tau i} n_{\tau ij} = n_1 n_{\tau i} n_{\tau j} = n$ و $(n_{\tau ij} p_{ijk}^{(r)}) (n_{\tau i} p_{ij}^{(r)}) (n_i p_i^{(1)})$ که احتمال انتخاب کلی برای هر واحد مقدماتی است.

برای مقدار کل y_{hr} در (۱۰) از نمونه‌گیری pps با جایگذاری برای دو مرحله اول و نمونه‌گیری با احتمال برابر بدون جایگذاری برای مرحله سوم، می‌توان نشان داد که

$$\begin{aligned} \text{var}(y_{hr}) = & \frac{n_1}{N} \sum_i^{N_i} N_{\tau i+} (T_i - \bar{T})^2 + \frac{n_1}{N} \sum_i^{N_i} n_{\tau i} \sum_j^{N_{\tau ij}} N_{\tau ij} (T_{ij} - \bar{T}_i)^2 \\ & + \frac{n_1}{N} \sum_{i=1}^{N_i} n_{\tau i} \sum_{j=1}^{N_{\tau ij}} n_{\tau ij} N_{\tau ij} \sigma_{hij}^2 (1 - F_{\tau ij}) \end{aligned} \quad (13)$$

که $T_{ij} = n_{\tau ij} \mu_{ij}$; $T_i = \sum_j^{N_{\tau ij}} T_{ij}$; $\mu_{ij} = E(m_{ij})$; $m_{ij} = \frac{1}{n_{\tau ij}} \sum_k^{n_{\tau ij}} y_{hijk}$

$$\begin{aligned} \bar{T} = n_1 \sum_i^{N_i} \left(\frac{N_{\tau i+}}{N} \right) \bar{T}_i \quad \text{که} \\ \sigma_{hij}^2 = \frac{1}{N_{\tau ij} - 1} \sum_{k=1}^{N_{\tau ij}} (y_{hijk} - \mu_{ij})^2, N_{\tau i+} = \sum_{j=1}^{N_{\tau ij}} N_{\tau ij}, \bar{T}_i = \sum_j^{N_{\tau ij}} \frac{(N_{\tau ij})}{N_{\tau i+}} T_{ij} \end{aligned}$$

ممکن است یک برآورد ناریب از (۱۳) یافت شود که نتیجه آن با بحث سابق در مورد نمونه‌گیری با جایگذاری در یک مرحله، سازگار باشد.

خ) $\text{var}(y_{hr})$ برای pps با جایگذاری در مرحله اول و حالت بدون جایگذاری با احتمال برابر در مرحله دوم:

نتیجه برای حالت دو مرحله‌ای، با قرار دادن $n_1 = N_1 = 1$ و ایجاد تغییرات مناسب در

علائم، از (۱۳) به دست می‌آید. به‌طور مشابه، می‌توانیم واریانس نسبت و کواریانس نسبت‌ها را نیز به دست آوریم.

۴- خطی‌سازی

برآوردهای نسبتی می‌تواند با یک بسط سری‌های تیلور تحت جمع، خطی شود. سپس واریانس سهم متغیر از این بسط، شبیه به واریانس نسبت اصلی است. هر نسبت از متغیرهای u_1, \dots, u_k را با تابع $f(u_1, u_2, \dots, u_k)$ نشان می‌دهیم. داریم:

$$(14) \quad \text{var}(f(u_1, \dots, u_k)) \approx \text{var}\left(\sum_i u_i \frac{\partial F(u_1, \dots, u_k)}{\partial u_i}\right)$$

که $E(u_i) = U_i$ ($i = 1, \dots, k$) و علامت « \approx » به مفهوم این است که هر دو طرف علامت به‌طور تقریبی مساوی هستند.

مثال ۱) فرض کنید x و y متغیرهای تصادفی با مقادیر مورد انتظار X و Y باشند. نسبت $\frac{x}{y}$ را در نظر بگیرید. واریانس نسبت با واریانس سهم متغیر از یک بسط خطی از آن، تقریب زده شده است:

$$(15) \quad \text{var}\left(\frac{x}{y}\right) \approx \text{var}\left(\frac{x - Ry}{Y}\right)$$

$$\text{که } R = \frac{X}{Y}$$

مثال ۲) نسبت $X' = \left(\frac{x}{y}\right) \hat{Y}$ اغلب برای اهداف برآورد استفاده شده است که x و y متغیر هستند در حالی که \hat{Y} یک عدد معلوم است.

$$(16) \quad \text{var}\left(\frac{x}{y} \hat{Y}\right) \approx \text{var}\left(\frac{\hat{Y}}{Y}(x - Ry)\right)$$

با استفاده از این دو روش که در بخش ۳ و ۴ ارائه شد، می‌توانیم یک واریانس برای

برآورد نسبی پیچیده به دست آوریم. ما این واریانس را با استفاده از یک مثال واقعی در بخش ۵ ارائه خواهیم کرد.

۵- یک مثال و خلاصه بحث

برآورد x برای مشخصه X جامعه در آمارگیری جاری جمعیتی به شکل زیر است: (گزارش فنی ۴۰، اداره سرشماری، صفحه ۱۵۵)

$$(17) \quad X' = \sum_a \frac{x_{aSR} + \sum_{c=1}^C \frac{x_{acNS}}{z_c} Z_c}{y_{aSR} + \sum_{c=1}^C \frac{y_{acNS}}{z_c} Z_c} \hat{Y}_a$$

که زیرنویس $c = 1, \dots, C$ ، برابر با ۴۸ سلول از رسته رنگ-محل سکونت) برای طبقات غیر خود نماینده (NS) از تعدیل نسبی اول و $a = 1, \dots, A$ (برابر با ۶۰ سلول از رسته‌های سن-جنسیت-نژاد) از تعدیل نسبی دوم حاصل می‌شود. ما رابطه (۱۷) را به شکل زیر بیان می‌کنیم.

$$(18) \quad X' = \sum_a \frac{x'_a}{y'_a} \hat{Y}_a$$

که x'_a و y'_a در (۱۷) تعریف شده‌اند، عبارت‌های $x_{aSR}, y_{aSR}, x_{acNS}, y_{acNS}, z_c$ و \hat{Y}_a در (۱۷) به شکل زیر تعریف می‌شوند:

x_{aSR} = مقدار کل موزون نمونه، از آخرین واحدهای نمونه‌گیری (USU) در واحدهای نمونه‌گیری مقدماتی (PSUs) خود نماینده (SR)، برای جامعه با ویژگی مطلوب x در a امین رسته سن-جنسیت-نژاد. وزن‌ها عبارت از عکس احتمال انتخاب USUها هستند. در عمل، وزن‌ها که با w_{sij} نشان داده می‌شوند، شامل وزن‌دهی خاص و فاکتورهای تعدیل برای بی‌پاسخی نیز هستند.

y_{aSR} = همان‌طور که برای x_{aSR} گفته شد، اما برای مقدار کل جامعه.

$x_{acNS} =$ همان طور که برای x_{aSR} گفته شد، اما برای c امین رسته نژاد- محل سکونت در مورد جامعه غیر خود نماینده (NS)

$y_{acNS} =$ همان طور که برای x_{ac} گفته شد، اما برای مقدار کل جامعه.

$Z_c =$ مقدار کل برآورد شده در سرشماری سال ۱۹۸۰ برای جامعه NSR در c امین رسته ادغام شده نژاد- محل سکونت، بر اساس جمعیت سرشماری ۱۹۸۰ برای PSUهای نمونه NSR که روی تمام طبقات NS وزن دهی و جمع شده است.

$Z_{csij} = 1$ اگر (sij) امین فرد در PSU مربوط به نمونه NSR متعلق به c امین رسته رنگ- محل سکونت باشد و در دیگر موارد $Z_{csij} = 0$.

$Z_c =$ جمعیت سرشماری ۱۹۸۰ در طبقات NS در c امین رسته ادغام شده نژاد- محل سکونت.

$Y_a =$ مقدار کل مستقل جامعه برای ایالات متحده در a امین رسته سن- جنسیت- نژاد برای ماه جاری CPS.

برای واریانس X' ، نخست X' را تحت علامت مجموع که روی a جمع می‌بندد، خطی می‌کنیم و سپس فقط سهم متغیر از بسط خطی را در نظر می‌گیریم. در مرحله دوم علامت‌های مجموع را تعویض کرده، آن‌ها را برای نمونه‌گیری به جلو منتقل می‌کنیم و برای برآوردهای نسبت جمع می‌بندیم. در مرحله سوم واریانس سرجمع حاصل می‌تواند با استفاده از فرمولی که در بخش ۳ ارائه شد، به دست آید.

این مراحل به شکل زیر نشان داده می‌شود: *روش‌های آماری*
مرحله اول- (۱۸) را خطی کرده و سهم متغیر را از آن استخراج نمایید. واریانس برآورد نسبت با آن سهم خطی شده تقریب زده می‌شود.

$$(19) \quad \text{var}(X') \approx \text{var}\left(\sum_a^A \gamma_a \left(x'_a - \frac{\tilde{X}_a}{\tilde{Y}_a} y'_a\right)\right)$$

از (۱۷) و (۱۹) می‌توان نوشت:

$$\text{var}(X') \approx \text{var}\left(\sum_a^A \gamma_a \left((x_{aSR} + \sum_c^C \frac{x_{acNS}}{Z_c} Z_c) - \frac{\tilde{X}_a}{\tilde{Y}_a} (y_{aSR} + \sum_c^C \frac{y_{acNS}}{Z_c} Z_c)\right)\right)$$

$$\text{که } \tilde{Y}_a = E(y'_a) \text{ و } \tilde{X}_a = E(x'_a), \gamma_a = \frac{\hat{Y}_a}{\tilde{Y}_a}$$

عبارات مربوط به طبقات SR و طبقات NS را به طور مجزا جمع‌آوری می‌کنیم و برآورد نسبتی را با جمع زدن روی c برای دومین مرتبه خطی می‌کنیم، سپس به دست می‌آید:

$$\begin{aligned} \text{var}(\hat{X}) &\approx \text{var}\left(\sum_a \gamma_a (x_{aSR} - \frac{\tilde{X}_a}{\tilde{Y}_a} y_{aSR})\right) \\ (20) \quad &+ \sum_a \gamma_a \sum_c \frac{Z_c}{\tilde{Z}_c} (x_{acNS} - \frac{\tilde{X}_a}{\tilde{Y}_a} y_{acNS} + \frac{\tilde{X}_a}{\tilde{Y}_a} \frac{\tilde{Y}_{ac}}{\tilde{Z}_c} z_c - \frac{\tilde{X}_{ac}}{\tilde{Z}_c} z_c) \end{aligned}$$

که

$$x_{aSR} = \sum_s \sum_i \sum_j^{n_{rs}} w_{sij} x_{asij}, \quad \tilde{X}_{ac} = E(x_{acNS}), \quad \tilde{Y}_{ac} = E(y_{acNS}), \quad \tilde{Z}_c = E(z_c)$$

$$\begin{aligned} x_{acNS} &= \sum_s \sum_i \sum_j^{n'_{rs}} w_{sij} x_{acsij}, \quad y_{aSR} = \sum_s \sum_i \sum_j^{n_{rs}} w_{sij} y_{asij} \\ z_c &= \sum_s \sum_i \sum_j^{n'_{rs}} w_{sij} z_{csij}, \quad y_{acNS} = \sum_s \sum_i \sum_j^{n'_{rs}} w_{sij} y_{acsij} \end{aligned}$$

L_1 تعداد SR-PSUها، n_{rs} تعداد واحدهای نمونه‌گیری مرحله دوم (SSUs) در s امین SR-PSU و n_{rsi} تعداد واحدهای نمونه‌گیری مرحله سوم (TSUs) در i امین SSU می‌باشد. n'_1 تعداد NS-PSUهای نمونه‌گیری شده است.

مرحله دوم- با استفاده از تعاریف بالا از x_{aSR} ، y_{aSR} ، x_{acNS} ، y_{acNS} و z_c با مجموع‌های سه‌گانه، ما علامت‌های جمع ناشی از نمونه‌گیری را به جلوی علامت‌های برآورد نسبتی برای سن-جنسیت-رنگ و سکونت-نژاد، منتقل می‌کنیم. برای سن-جنسیت-رنگ و سکونت-نژاد، روی a و c جمع می‌بندیم، فقط علائم مجموع مربوط به نمونه‌گیری باقی می‌ماند، (۲۰) را می‌توان به صورت زیر نوشت:

(۲۱) $\text{var}(X') \cong$

$$\text{var}\left(\sum_S \sum_i \sum_j^{n_{sr}, n_{sri}} (B_{+sij} + C_{+sij}) + \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} (B'_{++sij} - C'_{++sij} + D'_{++sij} - D_{++sij})\right)$$

که می تواند به صورت زیر بیان شود:

(۲۲) $\text{var}(X') = \text{var}\left(\sum_S \sum_i \sum_j^{L_1, n_{sr}, n_{sri}} t_{sij} + \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} t'_{sij}\right)$

که t_{sij} و t'_{sij} در آن ها به خوبی تعریف شده اند. شش عبارت در (۲۱) به صورت زیر تعریف می شوند:

(۲۱-۱):

$$\sum_a \gamma_a x_{aSR} = \sum_S \sum_i \sum_j^{L_1, n_{sr}, n_{sri}} \sum_a B_{asij} = \sum_S \sum_i \sum_j^{L_1, n_{sr}, n_{sri}} B_{+sij}, B_{asij} = \gamma_a w_{sij} x_{asSRij};$$

(۲۱-۲):

$$\sum_a \frac{\tilde{X}_a}{\tilde{Y}_a} \gamma_a y_{aSR} = \sum_S \sum_i \sum_j^{L_1, n_{sr}, n_{sri}} \sum_a C_{asij} = \sum_S \sum_i \sum_j^{L_1, n_{sr}, n_{sri}} C_{sij}, C_{asij} = \frac{\tilde{X}_a}{\tilde{Y}_a} \gamma_a w_{sij} y_{asSRij};$$

برای عبارت دوم،

(۲۱-۳):

$$\sum_a \sum_c \gamma_a x_{ac} = \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} \sum_a \sum_c B'_{acsij} = \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} B'_{++sij}, B'_{acsij} = \gamma_a w_{sij} x_{acNSsij};$$

(۲۱-۴):

$$\sum_a \sum_c \gamma_a \frac{\tilde{X}_a}{\tilde{Y}_{ac}} y_{ac} = \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} \sum_a \sum_c C'_{acsij} = \sum_S \sum_i \sum_j^{n'_{sr}, n_{sri}} C'_{++sij}, C'_{acsij} = \gamma_a \frac{\tilde{X}_a}{\tilde{Y}_{ac}} w_{sij} y_{acNSsij};$$

(۲۱-۵):

$$\sum_a^A \sum_c^C \gamma_a \frac{\tilde{X}_a}{\tilde{Y}_a} \frac{\tilde{Y}_{ac}}{\tilde{Z}_c} Z_c = \sum_S^{n'_s} \sum_i^{n_{SY}} \sum_j^{n_{STj}} \sum_a^A \sum_c^C D'_{acsij} = \sum_S^{n'_s} \sum_i^{n_{SY}} \sum_j^{n_{STj}} D'_{++sij},$$

$$D'_{csij} = \gamma_a \frac{\tilde{X}_a}{\tilde{Y}_a} \frac{\tilde{Y}_{ac}}{\tilde{Z}_c} w_{sij} z_{acsij};$$

(۲۱-۶):

$$\sum_a^A \sum_c^C \gamma_a \frac{\tilde{Y}_{ac}}{\tilde{Z}_c} z_c = \sum_S^{n'_s} \sum_i^{n_{SY}} \sum_j^{n_{STj}} \sum_a^A \sum_c^C D_{acsij} = \sum_S^{n'_s} \sum_i^{n_{SY}} \sum_j^{n_{STj}} D_{++sij}$$

$$D_{acsij} = \gamma_a \frac{\tilde{Y}_{ac}}{\tilde{Z}_c} w_{sij} z_{csij};$$

مرحله سوم- واریانس X' - عبارت‌های اول و دوم به ترتیب حاصل از SR-PSUها و برای NS-PSUها هستند. می‌توانیم واریانس‌های آن‌ها را به‌طور جداگانه محاسبه کرده و با یکدیگر جمع کنیم.

از آن‌جا که برای طبقه‌های غیر خود- نماینده، نمونه‌ها برای دو مرحله اول با طرح PPS با جایگذاری و برای مرحله آخر با احتمال‌های برابر بدون جایگذاری انتخاب شده است، وقتی که داده‌های طبقات خود- نماینده، همان‌گونه که در (۷) نشان داده شد، فقط از دو مرحله حاصل شده باشند، می‌توانیم (۶) را برای برآورد واریانس طبقات غیر خود- نماینده به کار ببریم. واریانس X' عبارت است از:

$$\text{var}\left(\sum_S^{L_1} \sum_i^{n_{SY}} \sum_j^{n_{STj}} t_{sij}\right) + \text{var}\left(\sum_S^{n'_s} \sum_i^{n_{SY}} \sum_j^{n_{STj}} t'_{sij}\right) = \sigma_{sSR}^2 + \sigma_{sNS}^2$$

عبارت اول مجموع L_1 واحد نمونه‌گیری مقدماتی (PSUs) خود- نماینده است، در حالی که عبارت دوم مجموع n'_s واحد نمونه‌گیری مقدماتی (PSUs) غیر خود- نماینده است.

ممکن است فردی بخواهد واریانس نسبت X'/Y' را داشته باشد، که X' همان‌طور که نشان داده شد خطی شده است، Y' برآورد نسبتی دیگری است که با

استفاده از (۱۷) حاصل می‌شود و $\text{var}(X'/Y')$ نیز به طریقی مشابه به دست می‌آید. ممکن است که به‌طور مکرر از روش دلتا برای $\text{var}(X')$ استفاده کنیم. با استفاده از این مثال، کاربرد نخست می‌تواند برای رسته‌های سن-جنسیت-نژاد به کار رود، کاربرد دوم برای سلول‌های سکونت-رنگ و در نهایت کاربرد سوم برای طرح سه مرحله‌ای مورد استفاده قرار گیرند. در اینجا فرض‌های دیگری، بجز آن‌هایی که از تقریب سری‌های تیلور حاصل می‌شود، فرض دیگری وجود ندارد.

روش‌هایی که در این مقاله ارائه شد، از این جهت که می‌توانیم روش‌های نمونه‌گیری و برآورد را روی واریانس نمونه منعکس کنیم، ممکن است یکی از بهترین روش‌های برآورد واریانس باشد. عملکرد این روش با استفاده از یک روش تجربی قابل تشخیص است.

این روش ممکن است که برای متغیرهای گسسته نیز به‌شکل متغیرهای پیوسته به کار گرفته شود.

مرجع‌ها

- [1] Durbin, J. (1953). Some results in sampling theory when the units are selected with unequal probabilities. *Journal of the Royal Statistical Society, B*, 15, 262.
- [2] Hidiroglou, M. A. and Rao, J.N.K. (1983). Chi-Square Tests for the Analysis of Three Way Contingency Tables from the Canada Health Survey. *Statistics Canada, Ottawa, Canada*.
- [3] Shah, B.V. (1981). SESUDAAN: Standard Errors Program for Computing of Standardized Rates from Sample Survey Data. Unpublished document, RTI.
- [4] Kendall M. G. and Stuart A. S. (1968). *The Advanced theory of Statistics*, Vol. 3. Hafner Publishing Company, New York.
- [5] Woodruff, Ralph S. (1971). Simple Method for Approximating Variance of a Complicated Estimate. *Journal of the American Statistical Association*, Volume 66, June, pp411-414.