

نمونه‌گیری

روشها و کاربردها

پل اس. لهوی

استنلی لمی شو

مترجم

گیتی مختاری امیرمجدی



پژوهشکده آمار

لوی، پل	Levy, Paul S.
نمونه‌گیری: روشها و کاربردها / پل اس. لهوی، استنلی لمی شو؛ مترجم گیتی مختاری امیرمجدی	
-- تهران: مرکز آمار ایران، پژوهشکده آمار، ۱۳۸۱.	
۵۸۵ ص. : جدول، نمودار.	
فهرست‌نویسی بر اساس اطلاعات فیفا.	ISBN 964-365-150-9 ریال : ۳۰۰۰۰
عنوان اصلی: Sampling of Populations: Methods and Applications.	
کتاب‌نامه.	
نمایه	
۱. جمعیت - روشهای آماری. ۲. آمارگیری نمونه‌ای. الف. لمشو، استنلی Lemeshow, Stanley	
ب. مختاری امیرمجدی، گیتی، ۱۳۳۳- ، مترجم. ج. مرکز آمار ایران، پژوهشکده آمار. د. عنوان.	
ان ۹/۴۹/۴۹ HB8۴۹/۴۹/۴۹	۳۰۴/۶۰۱۵۱۹۵۲
۱۳۸۱	
کتابخانه ملی ایران	۴۰۶۴۷-۸۱ م

- مدیریت تولید : گروه پژوهشی طرحهای فنی و روشهای آماری
- ویراستار علمی و ادبی : دکتر علی عمیدی
- حروف‌نگاری، نمونه‌خوانی، صفحه‌بندی، صفحه‌آرایی : محبوبه کاظمی
- ویراستار هنری : فرشید خان‌زاده
- طراحی جلد : نیما دانش‌پرور
- مدیر فنی : علی اصغر حائری مهریزی
- امور فنی و چاپ : مؤسسه انتشارات ستایش

© ۱۳۸۱ پژوهشکده آمار

شماره ۵۲، خیابان شهید فکوری، خیابان باباطاهر، خیابان دکتر فاطمی
تهران ۱۴۱۳۷۱۷۹۱۱، ایران



URL: <http://www.src.ac.ir>

e-mail: src@src.ac.ir

تلفن: ۸۹۵۹۰۲۹ دورنگار: ۸۰۰۷۹۸۹

همه حقوق این اثر برای پژوهشکده آمار محفوظ است. هیچ بخشی از این کتاب را نمی‌توان بدون اجازه کتبی از ناشرش تکثیر یا به هر شکلی و با هر وسیله‌ای ذخیره کرد. استفاده یا تقلید از طرح جلد، ممنوع است.

حروف‌نگاری شده با قلم‌های فارسی لوتوس و تیترو میترا، و قلم لاتین Times New Roman.

چاپ و صحافی شده در ایران.

چاپ یکم

شمارگان: ۶۰۰

پیشنهاد برای نحوه نقل مطلب، جدول یا نمودار از این کتاب، به صورت زیر است:

مختاری امیرمجدی، گیتی (۱۳۸۱). نمونه‌گیری: روشها و کاربردها. لهوی، پل اس.؛ لمی شو، استنلی. ترجمه از انگلیسی به فارسی. تهران: پژوهشکده آمار.

شابک ۹۶۴ - ۳۶۵ - ۱۵۰ - ۹

ISBN 964-365-150-9

بها: سی و پنج هزار ریال

فصل ۷

برآورد نسبتی

در این فصل مفهوم نسبت $\frac{\bar{x}}{\bar{y}}$ دو میانگین نمونه‌ای \bar{x} و \bar{y} را ارائه می‌کنیم. این نسبت به عنوان برآوردی از نسبت $\frac{X}{Y}$ میانگینهای دو متغیر x و y در جامعه به کار می‌رود. ولی مهمتر از آن، کاربرد این نسبت به عنوان ابزاری برای به دست آوردن برآورد X ، مجموع کل جامعه است که می‌توان آن را دقیقتر از برآورد x' به وسیله ضرب ساده مجموع نمونه‌ای x در $\frac{N}{n}$ یعنی عکس کسر نمونه‌گیری تعیین کرد. این روش را برآورد کردن نسبتی و برآوردهای حاصل از آن را برآوردهای نسبتی می‌نامند. برای شروع این مبحث به مثالی از برآورد نسبتی توجه می‌کنیم.

مثال تشریحی: جامعه‌ای را در نظر می‌گیریم که دارای هشت محله است. فرض کنید می‌خواهیم نسبت R کل هزینه‌های دارویی X را به کل هزینه‌های درمانی Y در میان تمام اشخاص موجود در جامعه برآورد کنیم. برای انجام این کار قرار است نمونه تصادفی ساده‌ای متشکل از دو محله گرفته شود و با یکایک خانوارها در هر محله نمونه مصاحبه به عمل آید.

عناصر این طرح نمونه‌گیری، محله‌ها هستند و نسبت $\frac{X}{Y}$ ، کل هزینه‌های دارویی به کل هزینه‌های درمانی، را می‌توان برحسب r برآورد کرد که $r = \frac{x}{y}$ نسبت به دست آمده از نمونه است. فرض کنیم داده‌های مربوط به محله‌ها به صورتی باشند که در جدول ۱.۷ ارائه شده است.

حالا فرض کنید مثلاً محله‌های ۲ و ۵ برای نمونه انتخاب شده‌اند. پس داریم

$$x = 50000 + 150000 = 200000$$

$$y = 200000 + 450000 = 650000$$

و

$$r = \frac{x}{y} = \frac{200000}{650000} = 0.308$$

به این ترتیب، نسبت کل هزینه‌های دارویی به کل هزینه‌های درمانی برابر است با ۰/۳۰۸.

□

جدول ۱.۷ هزینه‌های دارویی و کل هزینه‌های درمانی بین همه ساکنان هشت محله

محله	کل هزینه‌های دارویی (دلار) X	کل هزینه‌های درمانی (دلار) Y
۱	۱۰۰۰۰۰	۳۰۰۰۰۰
۲	۵۰۰۰۰	۲۰۰۰۰۰
۳	۷۵۰۰۰	۳۰۰۰۰۰
۴	۲۰۰۰۰۰	۶۰۰۰۰۰
۵	۱۵۰۰۰۰	۴۵۰۰۰۰
۶	۱۷۵۰۰۰	۵۲۰۰۰۰
۷	۱۷۰۰۰۰	۶۸۰۰۰۰
۸	۱۵۰۰۰۰	۴۵۰۰۰۰
مجموع	۱۰۷۰۰۰۰	۳۵۰۰۰۰۰

۱.۷ برآورد کردن نسبی تحت نمونه‌گیری تصادفی ساده

مثال فوق را می‌توان به شرح زیر تعمیم داد. فرض کنید یک نمونه تصادفی ساده n عنصری از جامعه‌ای N عنصری داریم و می‌خواهیم نسبت R دو مجموع کل X و Y جامعه را برآورد کنیم. واضح است که R از فرمول

$$R = \frac{X}{Y}$$

یا از معادل آن

$$R = \frac{\bar{X}}{\bar{Y}}$$

به دست می‌آید.

نسبت R می‌تواند به وسیله r با استفاده از فرمول زیر

$$r = \frac{x'}{y'}$$

یا از دو فرمول زیر که از نظر جبری معادل‌اند به دست آید

$$r = \frac{\bar{x}}{\bar{y}} \quad \text{و} \quad r = \frac{x}{y}$$

برآورد r ، برآورد نسبتی نامیده می‌شود زیرا هم صورت کسر \bar{x} و هم مخرج کسر \bar{y} در معرض تغییرات نمونه‌گیری قرار دارند.

حالا به مثال ارائه شده در مقدمه این فصل برگردیم و ببینیم که آیا نسبت برآورد شده، یک نسبت نارایب است یا نه.

مثال تشریحی: در مثال بالا $\binom{8}{2} = 28$ نمونه تصادفی ساده ممکن دو عنصری از جامعه متشکل از ۸

عنصر وجود دارند. این نمونه‌های ممکن همراه با مقادیر x ، y ، و r به دست آمده از هر نمونه در جدول ۲.۷ فهرست شده‌اند.

توزیع نمونه‌گیری r ، برآورد نسبت، دارای میانگین $E(r)$ و خطای معیار $SE(r)$ به شرح زیر است (تابلوی ۳.۲ را ببینید)

$$E(r) = \left(\frac{1}{28}\right) \times (0/300 + 0/292 + 0/284 + 0/276) = 0/3054$$

$$SE(r) = \left(\frac{1}{\sqrt{28}}\right) \times \left[(0/300 - 0/3054)^2 + (0/292 - 0/3054)^2 + 0/284 - 0/3054)^2 + (0/276 - 0/3054)^2 \right]^{1/2} = 0/0268$$

مقدار واقعی R از جدول ۱.۷ عبارت است از

$$R = \frac{1070000}{3500000} = 0/3057$$

تفاوت بین $E(r)$ میانگین توزیع نمونه‌گیری r ، برآورد نسبت و نسبت واقعی R جامعه ناشی از خطای گرد کردن نیست. به این ترتیب می‌بینیم که r ، برآورد نسبت الزاماً یک برآورد نارایب از نسبت جامعه‌ای R نیست. ولی در بیشتر موارد اریبی r ، برآورد نسبت کم است و برآوردهایی از این دست موارد استفاده زیادی دارند. □

توجه کنید که برای محاسبه $SE(r)$ در مثال بالا، تولید همه نمونه‌های ممکن و محاسبه خطای معیار توزیع نمونه‌گیری لازم است. البته، در کاربست واقعی، همه توزیع نمونه‌گیری را تولید نمی‌کنیم. ولی برخلاف سایر برآوردهایی که تاکنون در این کتاب بررسی کردیم (مانند میانگینها، مجموعها و نسبتها)، در r ، برآورد نسبت، هم صورت و هم مخرج کسر در معرض تغییرات نمونه‌گیری قرار دارند.

جدول ۲.۷ نمونه‌های ممکن دوعنصری از جامعه‌ای متشکل از هشت عنصر (جدول ۱.۷)

r	y	x	شماره محله‌ها در نمونه
۰/۳۰۰	۵۰۰۰۰۰	۱۵۰۰۰۰	۱, ۲
۰/۲۹۲	۶۰۰۰۰۰	۱۷۵۰۰۰	۱, ۳
۰/۳۳۳	۹۰۰۰۰۰	۳۰۰۰۰۰	۱, ۴
۰/۳۳۳	۷۵۰۰۰۰	۲۵۰۰۰۰	۱, ۵
۰/۳۳۵	۸۲۰۰۰۰	۲۷۵۰۰۰	۱, ۶
۰/۲۷۶	۹۸۰۰۰۰	۲۷۰۰۰۰	۱, ۷
۰/۳۳۳	۷۵۰۰۰۰	۲۵۰۰۰۰	۱, ۸
۰/۲۵۰	۵۰۰۰۰۰	۱۲۵۰۰۰	۲, ۳
۰/۳۱۳	۸۰۰۰۰۰	۲۵۰۰۰۰	۲, ۴
۰/۳۰۸	۶۵۰۰۰۰	۲۰۰۰۰۰	۲, ۵
۰/۳۱۳	۷۲۰۰۰۰	۲۲۵۰۰۰	۲, ۶
۰/۲۵۰	۸۸۰۰۰۰	۲۲۰۰۰۰	۲, ۷
۰/۳۰۸	۶۵۰۰۰۰	۲۰۰۰۰۰	۲, ۸
۰/۳۰۶	۹۰۰۰۰۰	۲۷۵۰۰۰	۳, ۴
۰/۳۰۰	۷۵۰۰۰۰	۲۲۵۰۰۰	۳, ۵
۰/۳۰۵	۸۲۰۰۰۰	۲۵۰۰۰۰	۳, ۶
۰/۲۵۰	۹۸۰۰۰۰	۲۴۵۰۰۰	۳, ۷
۰/۳۰۰	۷۵۰۰۰۰	۲۲۵۰۰۰	۳, ۸
۰/۳۳۳	۱۰۵۰۰۰۰	۳۵۰۰۰۰	۴, ۵
۰/۳۳۵	۱۱۲۰۰۰۰	۳۷۵۰۰۰	۴, ۶
۰/۲۸۹	۱۲۸۰۰۰۰	۳۷۰۰۰۰	۴, ۷
۰/۳۳۳	۱۰۵۰۰۰۰	۳۵۰۰۰۰	۴, ۸
۰/۳۳۵	۹۷۰۰۰۰	۳۲۵۰۰۰	۵, ۶
۰/۲۸۳	۱۱۳۰۰۰۰	۳۲۰۰۰۰	۵, ۷
۰/۳۳۳	۹۰۰۰۰۰	۳۰۰۰۰۰	۵, ۸
۰/۲۸۸	۱۲۰۰۰۰۰	۳۴۵۰۰۰	۶, ۷
۰/۳۳۵	۹۷۰۰۰۰	۳۲۵۰۰۰	۶, ۸
۰/۲۸۳	۱۱۳۰۰۰۰	۳۲۰۰۰۰	۷, ۸

از این رو نمی‌توان عبارت دقیقی برای خطای معیار آن به دست آورد. ولی هنس و همکاران [۱] پیشنهاد کرده‌اند که اگر ضریب تغییرات $V(\bar{y})$ مخرج کسر، \bar{y} ، یا r (با استفاده از حالت $r = \frac{\bar{x}}{\bar{y}}$) کوچکتر از یا برابر با ۰/۰۵ باشد، آنگاه می‌توان خطای معیار r ، یعنی $SE(r)$ را از عبارت زیر تقریب زد.

$$SE(r) \approx \left(\frac{R}{\sqrt{n}} \right) \times (V_x^2 + V_y^2 - 2\rho_{xy} V_x V_y)^{1/2} \times \sqrt{\frac{N-n}{N-1}} \quad (۱.۷)$$

که در آن ρ_{xy} ، ضریب همبستگی بین x و y است و به صورت زیر تعریف می‌شود

$$\rho_{xy} = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})/N}{\sigma_x \sigma_y} \quad (۲.۷)$$

حال بعضی محاسبات را برای مثالی که با آن کار می‌کردیم انجام می‌دهیم.

مثال تشریحی: برای داده‌های جدول ۱.۷ پارامترهای جمعیتی زیر را داریم [رجوع کنید به عبارتهای تابلوی ۱.۲ و معادله‌های (۷.۲) و (۲.۷)]:

$\sigma_x = ۴۹۴۱۸/۵$	$\bar{X} = ۱۳۳۷۵۰$	$R = ۰/۳۰۵۷$
$\sigma_y = ۱۵۲۷۰۴/۸$	$\bar{Y} = ۴۳۷۵۰۰$	$V_x = ۰/۳۶۹۵$
$N = ۸$	$\rho_{xy} = ۰/۹۲۷۲$	$V_y = ۰/۳۴۹۰$

برای داده‌های جدول ۲.۷، ضریب تغییرات مخرج کسر، \bar{y} ، برای r عبارت است از

$$V(\bar{y}) = \left(\frac{1}{\bar{Y}} \right) \times \left(\frac{\sigma_y}{\sqrt{n}} \right) \times \sqrt{\frac{N-n}{N-1}} = \left(\frac{1}{۴۳۷۵۰۰} \right) \times \left(\frac{۱۵۲۷۰۴/۸}{\sqrt{۲}} \right) \times \sqrt{\frac{۸-۲}{۸-1}} = ۰/۲۲۸۵$$

محاسبه با معادله (۱.۷) مقدار

$$SE(r) \approx \left(\frac{۰/۳۰۵۷}{\sqrt{۲}} \right) \times [(۰/۳۶۹۵)^2 + (۰/۳۴۹۰)^2 - 2(۰/۹۲۷۲)(۰/۳۶۹۵)(۰/۳۴۹۰)]^{1/2} \times \sqrt{\frac{۸-۲}{۸-1}} = ۰/۰۲۷۷$$

را در مقایسه با مقدار واقعی آن که ۰/۰۲۶۸ است نتیجه می‌دهد. چون ضریب تغییرات \bar{y} بزرگتر از ۰/۰۵ است، تقریب ارائه شده در معادله (۱.۷) به طور نرمال نباید برای $SE(r)$ مورد استفاده قرار بگیرد. ولی در این مورد به نظر می‌رسد که به صورتی معقول خوب عمل می‌کند.

حالا انتخاب نمونه‌هایی با اندازه $n=7$ را از جامعه مورد بررسی در نظر بگیرید. توزیع دقیق r برای نمونه‌هایی با $n=7$ عنصر در جدول ۳.۷ نشان داده شده است.

خطای معیار دقیق $SE(r)$ برای r که با شمارش کلیه نمونه‌های ممکن به دست آمده است ۰/۰۰۶۳ است. محاسبات نشان می‌دهد که

$$V(\bar{y}) = \left(\frac{1}{437500} \right) \times \left(\frac{152704/8}{\sqrt{7}} \right) \times \sqrt{\frac{8-7}{8-1}} = 0.0499$$

بنابراین، چون $V(\bar{y}) = 0.0499 \leq 0.05$ انتظار داریم تقریب معادله (۱.۷) تقریب خوبی از مقدار واقعی $SE(r)$ فراهم کند. این محاسبه به شرح زیر است:

$$SE(r) \approx \left(\frac{0.3057}{\sqrt{7}} \right) \times \left[(0.3695)^2 + (0.3490)^2 - 2(0.9272)(0.3695)(0.3490) \right]^{1/2} \times \sqrt{\frac{8-7}{8-1}} = 0.061$$

به این ترتیب می‌بینیم که تقریب (۱.۷) در این مثال تطابق بسیار نزدیکی با خطای معیار واقعی r دارد.

جدول ۳.۷ نمونه‌های هفت‌تایی از جامعه جدول ۱.۷

r	y	x	نمونه
۰/۳۰۱۶	۳۰۵۰۰۰۰	۹۲۰۰۰۰	۱, ۲, ۳, ۴, ۵, ۶, ۷
۰/۳۱۹۱	۲۸۲۰۰۰۰	۹۰۰۰۰۰	۱, ۲, ۳, ۴, ۵, ۶, ۸
۰/۳۰۰۳	۲۹۸۰۰۰۰	۸۹۵۰۰۰	۱, ۲, ۳, ۴, ۵, ۷, ۸
۰/۳۰۱۶	۳۰۵۰۰۰۰	۹۲۰۰۰۰	۱, ۲, ۳, ۴, ۶, ۷, ۸
۰/۳۰۰۰	۲۹۰۰۰۰۰	۸۷۰۰۰۰	۱, ۲, ۳, ۵, ۶, ۷, ۸
۰/۳۱۰۹	۳۲۰۰۰۰۰	۹۹۵۰۰۰	۱, ۲, ۴, ۵, ۶, ۷, ۸
۰/۳۰۹۱	۳۳۰۰۰۰۰	۱۰۲۰۰۰۰	۱, ۳, ۴, ۵, ۶, ۷, ۸
۰/۳۰۳۱	۳۲۰۰۰۰۰	۹۷۰۰۰۰	۲, ۳, ۴, ۵, ۶, ۷, ۸

□

خطای معیار برآورد یک نسبت را می‌توان از روی این داده‌ها و با جایگزین کردن $\hat{\sigma}_x^2$ به جای σ_x^2 ، $\hat{\sigma}_y^2$ به جای σ_y^2 ، $\hat{\rho}_{xy}$ به جای ρ_{xy} ، \bar{x} به جای \bar{X} ، \bar{y} به جای \bar{Y} و r به جای R در

رابطه (۱.۷) برآورد کرد. برآورد حاصل، $\hat{SE}(r)$ ، از فرمول زیر به دست می‌آید

$$\hat{SE}(r) = \left(\frac{r}{\sqrt{n}} \right) \times \left(\hat{V}_x^2 + \hat{V}_y^2 - 2\hat{\rho}_{xy} \hat{V}_x \hat{V}_y \right)^{1/2} \times \sqrt{\frac{N-n}{N-1}} \quad (3.7)$$

که در آن،

$$\hat{V}_x^2 = \left(\frac{N-1}{N} \right) \left(\frac{s_x^2}{\bar{x}^2} \right)$$

$$\hat{V}_y^2 = \left(\frac{N-1}{N} \right) \left(\frac{s_y^2}{\bar{y}^2} \right)$$

$$\hat{\rho}_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1) s_x s_y}$$

به طور کلی این تقریب تنها هنگامی مورد استفاده قرار می‌گیرد که نابرابری زیر برقرار باشد:

$$\frac{s_y}{\sqrt{n} \times \bar{y}} \times \sqrt{\frac{N-n}{N}} \leq 0.05$$

برای سهولت کار، فرمولهایی که می‌توانند برای برآورد کردن نسبتی مورد استفاده قرار گیرند در تابلوی ۱.۷ خلاصه شده‌اند.

تابلوی ۱.۷ فرمولهای برآورد کردن نسبتی تحت نمونه‌گیری تصادفی ساده

پارامترهای جامعه

$$R = \frac{X}{Y} = \frac{\bar{X}}{\bar{Y}}$$

$$SE(r) \approx \left(\frac{R}{\sqrt{n}} \right) \times \left(V_x^{\hat{}} + V_y^{\hat{}} - 2\rho_{xy} V_x V_y \right)^{1/2} \times \sqrt{\frac{N-n}{N-1}}$$

برآوردهای نمونه‌ای

$$r = \frac{x'}{y'} = \frac{x}{y} = \frac{\bar{x}}{\bar{y}}$$

$$\hat{SE}(r) = \left(\frac{r}{\sqrt{n}} \right) \times \left(\hat{V}_x^{\hat{}} + \hat{V}_y^{\hat{}} - 2\hat{\rho}_{xy} \hat{V}_x \hat{V}_y \right)^{1/2} \times \sqrt{\frac{N-n}{N-1}}$$

بازه اطمینان $(1-\alpha) \times 100$ درصدی را می‌توان به صورت زیر ساخت

$$r - z_{1-(\alpha/2)} \hat{SE}(r) \leq R \leq r + z_{1-(\alpha/2)} \hat{SE}(r)$$

V_x و V_y در معادله (۷.۲) تعریف شده‌اند. \bar{X} ، X ، σ_x ، Y ، \bar{Y} و σ_y به همان صورت تعریف شده در تابلوی ۱.۲ هستند. ρ_{xy} به صورت تعریف شده در معادله (۲.۷) است. x ، \bar{x} ، x' ، y ، \bar{y} ، y' و $S_x^{\hat{}}$ و $S_y^{\hat{}}$ به صورت تعریف شده در تابلوی ۲.۲ هستند. پس

$$\hat{V}_x^{\hat{}} = \left(\frac{N-1}{N} \right) \left(\frac{s_x^{\hat{}}}{\bar{x}^{\hat{}}} \right)$$

$$\hat{V}_y^{\hat{}} = \left(\frac{N-1}{N} \right) \left(\frac{s_y^{\hat{}}}{\bar{y}^{\hat{}}} \right)$$

$$\hat{\rho}_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1) s_x s_y}$$

مثال تشریحی: فرض کنید نمونه متشکل از محله‌های ۱، ۲، ۳، ۴، ۵، ۶ و ۸ را از جدول ۱.۷ انتخاب کرده‌ایم. محاسبه می‌کنیم که

$$\frac{s_y}{\sqrt{ny}} \times \sqrt{\frac{(N-n)}{N}} = 0.0468 \quad \text{و} \quad \bar{y} = 402857/1, \quad s_y = 141033/6$$

که کمتر از ۰/۰۵ است. به این ترتیب انتظار داریم که تقریب (۳.۷) برآورد خوبی از خطای معیار نسبت r ، برآورد فراهم نماید. برای این مثال داریم

$$\begin{aligned} r &= 0.3191 & \hat{\rho}_{xy} &= 0.9900 \\ s_x &= 54826/6 & s_y &= 141033/6 \\ \bar{x} &= 1288571/4 & \bar{y} &= 402857/1 \\ \hat{V}_x^2 &= 0.1591 & \hat{V}_y^2 &= 0.1072 \end{aligned}$$

پس

$$\hat{SE}(r) = \left(\frac{0.3191}{\sqrt{7}} \right) \times \left[0.1591 + 0.1072 - 2(0.9900) \sqrt{0.1591} \times \sqrt{0.1072} \right]^{1/2} \times \left(\frac{8-7}{8-1} \right)^{1/2} = 0.0040$$

در مقام مقایسه، مقدار واقعی جامعه‌ای $SE(r)$ برابر با ۰/۰۰۶۳ است. این تفاوت هیچ تعجبی ندارد زیرا معلوم شده است که برآورد واریانسها به شدت متغیرند.

□

r ، برآورد نسبت، به طور کلی یک برآورد اریب از R است. ولی هنگامی که اندازه‌های نمونه به صورتی معقول بزرگ باشند اریبی عموماً کوچک است و بازه‌های اطمینان تقریبی برای نسبت R مجهول جامعه را می‌توان با استفاده از برآورد خطای معیار r ، یعنی $\hat{SE}(r)$ ساخت.

برآوردهای نسبی در آمارگیریهای نمونه‌ای اهمیت دارند، بخصوص هنگامی که واحدهای شمارش همان واحدهای اولیه نیستند. برای مثال، اگر یک نمونه تصادفی ساده به منظور برآورد میانگین تعداد روزهای غیبت از کار به علت بیماری سخت به ازای هر نفر گرفته شده باشد، یک برآورد نسبی را به صورتی که در بالا توصیف شد می‌توان به کار برد که صورت کسر آن تعداد روزهای غیبت از کار به ازای هر خانوار و مخرج کسر آن تعداد اشخاص در خانوار خواهد بود. هم صورت و هم مخرج کسر در معرض تغییرپذیری نمونه‌گیری خواهند بود.

در بخش ۵.۳ در برآورد کردن میانگینهای جامعه‌ای برای زیرگروههای جامعه تحت نمونه‌گیری تصادفی ساده بحث کردیم. برآورد مناسب برای این وضعیت، حالتی خاص از نسبت برآورد شده است که مخرج کسر آن تعداد اعضای زیرگروه در نمونه است. اگر برنامه نمونه‌گیری، انتخاب نمونه تصادفی ساده‌ای از عناصر باشد، این صورت برآورد نسبی، برآوردی نارایب از میانگین زیرگروه مناسب است.

برآورد کردن نسبتها و خطاهای معیار آنها را می‌توان با نرم‌افزار SUDAAN و با استفاده از مدول *PROC RATIO* در این نرم‌افزار اجرا کرد. آنچه در پی می‌آید مجموعه فرمانهایی را نشان می‌دهد که برای مثال تشریحی قبل مورد استفاده قرار می‌گیرد.

```
1 PROC RATIO DATA=TAB7PT1 FILETYPE=SAS DESIGN=STRWOR;
2 TOTCNT TOTCNT;
3 WEIGHT WT1;
4 NEST_ONE_;
5 NUMBER PHARMEXP;
6 DENOM TOTMEDEX;
7 SETENV COLWIDTH = 1;
8 SETENV DECWIDTH = 4;
```

فرمان اول نشان می‌دهد که شیوه *PROC RATIO* قرار است مورد استفاده قرار گیرد و داده‌های نمونه در پرونده داده‌های SAS به نام *TAB7PT1.SSD* قرار دارند. شکل ظاهری این پرونده به صورت زیر است:

AREA	PHARMEXP	TOTMEDEX	TOTCNT	WT1
1	100000	300000	8	1.1428571
2	50000	200000	8	1.1428571
3	75000	300000	8	1.1428571
4	200000	600000	8	1.1428571
5	150000	450000	8	1.1428571
6	175000	520000	8	1.1428571
8	150000	450000	8	1.1428571

پرونده داده‌ها در بالا دارای ۷ سابقه (رکورد)، یعنی یک سابقه برای هر محله نمونه است. هر سابقه

دارای پنج متغیر به شرح زیر است:

area شماره شناسایی محله.

pharmexp کل هزینه‌های دارویی مربوط به هر محله.

totmedex کل هزینه‌های درمانی مربوط به هر محله.

totcnt کل تعداد N محله‌ها در جامعه.

wt1 وزن نمونه‌گیری، N/n .

فرمان *DESIGN=STRWOR* همراه با فرمان چهارم *NEST=_ONE_* نشان می‌دهند که طرح نمونه‌گیری، نمونه‌گیری تصادفی ساده بدون جایگذاری است. فرمان دوم نشان می‌دهد که N ، اندازه جامعه، به صورت متغیر *TOTCNT* در هر سابقه قرار دارد و در این مثال برابر با ۸ است. فرمان سوم

مشخص می‌سازد که وزن نمونه‌گیری $\frac{1}{N} = 1/1429$ در هر سابقه نمونه به صورت متغیر *WT1* ظاهر

می‌شود. بالاخره فرمانهای پنجم و ششم، متغیرهای صورت و مخرج کسر را برای محاسبه نسبت

مشخص می‌کنند و فرمان هفتم پهنای ستون و تعداد رقمهای دهدهی را در خروجی حاصل تعیین می‌کند.

خروجی حاصل از فرمانهای بالا در پایین نشان داده شده است. توجه کنید که برآورد خطای معیار، دقیقاً برابر با خطای معیار حاصل از فرمول نشان داده شده در تابلوی ۱.۷ است. جای تعجب نیست، زیرا SUDAAN دقیقاً از همان فرمول استفاده می‌کند.

Variable		One 1
PHARMEXP/TOTMEDEX	Sample Size	7.0000
	Weighted Size	8.0000
	Weighted X-Sum	3222978.0000
	Weighted Y-Sum	1028610.0000
	Ratio Est.	0.3191
	SE Ratio	0.0040

برآورد کردن با استفاده از STATA می‌تواند روی پرونده اطلاعاتی STATA به نام *tab7pt1.dta* با فرمانهای زیر انجام شود:

```
. use a:tab7pt1
. svyset fpc totcnt
. svyset pweight wt1
. svyratio pharmexp/totmedex
```

پرونده *tab7pt1.dta* دارای همان ساختاری است که پرونده *tab7pt1.ssd* در تحلیل SUDAAN مورد استفاده قرار داد.

این فرمانها خروجی زیر را تولید می‌کنند.

Survey ratio estimation

Pweight :	wt1	Number of obs =	7
Strata :	<one>	Number of strata =	1
PSU :	<observations>	Number of PSUs =	7
FPC :	totcnt	Population size =	8

Ratio	Estimate	Std.Err.	[95% Conf. Interval]	Deff
pharmexp/totmedex	.3191489	.0040067	.309345 .3289529	1

تصحیح جامعه متناهی (FPC)، با این فرض است که نمونه‌گیری تصادفی ساده بدون جایگذاری از واحدهای نمونه‌گیری اولیه (PSU) در داخل هر طبقه بدون نمونه‌گیری فرعی در داخل واحدهای نمونه‌گیری اولیه انجام می‌گیرد.

۲.۷ برآورد کردن نسبتها در زیرحوزه‌ها، تحت نمونه‌گیری تصادفی ساده

در بخش ۳.۶ به طور نسبتاً غیررسمی از برآورد کردن میانگینها برای زیرحوزه‌ها تحت نمونه‌گیری تصادفی ساده بحث کردیم. متوجه شدیم که میانگین یک زیرحوزه، در واقع یک برآورد نسبتی است که صورت کسر آن برابر است با مجموع نمونه‌ای متغیر مورد برآورد در همه عناصر نمونه متعلق به زیرحوزه و مخرج کسر آن برابر است با تعداد عناصر نمونه متعلق به زیرحوزه. در این بخش، برآورد نسبتها برای زیرحوزه‌ها را به وضعیتهایی تعمیم می‌دهیم که در آنها مخرج کسر نسبت الزاماً تعداد عناصر زیرحوزه در نمونه نیست.

فرض کنید یک نمونه تصادفی ساده n عنصری از جامعه‌ای N عنصری می‌گیریم و می‌خواهیم نسبت $R_{(k)}$ را برآورد کنیم که به صورت زیر تعریف می‌شود:

$$R_{(k)} = \frac{X_{(k)}}{Y_{(k)}}$$

که در آن $X_{(k)}$ و $Y_{(k)}$ ، به ترتیب مجموع متغیرهای X و Y روی همه عناصر موجود در یک زیرحوزه تعریف شده k هستند. با استفاده از تبدیلی از متغیرهای صورت و مخرج کسر، می‌توانیم از شیوه‌های برآورد که در بخش ۱.۷ شرح دادیم و از فرمولهایی که در تابلوی ۱.۷ نشان دادیم استفاده کنیم. این شیوه در زیر نشان داده می‌شود.

قرار دهید

$$\delta_{ik} = \begin{cases} 1 & \text{اگر عنصر } i \text{ در زیرحوزه } k \text{ باشد،} \\ 0 & \text{در غیر صورت بالا،} \end{cases}$$

و

$$r_{(k)} = \frac{x_{(k)}}{y_{(k)}}$$

$r_{(k)}$ ، برآورد خطای معیار نسبت زیرحوزه را، می‌توان با فرمولی که در تابلوی ۱.۷ نشان دادیم برآورد کرد. این مطلب با مثال زیر نشان داده می‌شود.

مثال تشریحی: فرض کنید در ناحیه‌ای که دارای ۱۰۱ بیمارستان است، نمونه‌ای متشکل از ۵۶ بیمارستان می‌گیریم و برای هر بیمارستان نمونه، پرونده همه نوزادانی را که در سال تقویمی گذشته در آن بیمارستان متولد شده‌اند مرور می‌کنیم تا تعیین کنیم که آیا وضعیت مادر از لحاظ پادگن سطحی هپاتیت **B** (HBsAg) در پرونده پزشکی بیماری نوزاد ثبت شده است یا نه. وجود این اطلاع در پرونده نوزاد از این نظر مهم است که اگر پادگن مادر مثبت باشد باید نوزاد گلوبولین مصون‌سازی دریافت کند تا از بروز هپاتیت **B** پیشگیری شود.

بیمارستانها براساس سطح خدمات بیماریهای زنان و زایمان طبقه‌بندی شده‌اند (۱= خدمات «اولیه»، ۲= «میانی»، ۳= «درجه ۳») و هدف، برآورد کردن نسبت نوزادانی است که وضعیت پادگن سطحی هپاتیت B مادران آنها در پرونده پزشکی نوزاد ثبت شده است. داده‌های مربوط به این مثال تشریحی در پرونده *BHRATIO.DAT* هستند. برای اهداف مثال، متغیرهای مربوط به ۱۰ سابقه اول آن پرونده در جدول زیر نشان داده شده‌اند.

بیمارستان نمونه	تعداد پرونده‌هایی که پادگن مادر در آنها گزارش شده است	تعداد متولد شده‌ها	سطح زنان و زایمان بیمارستان	نشانگرمتغیر برای زیرحوزه δ_i	متغیر صورت کسر تبدیل یافته (x_i)	متغیر مخرج کسر تبدیل یافته (y_i)
۱	۲۸۹۸	۲۸۹۸	۲	۱	۲۸۹۸	۲۸۹۸
۲	۱۰۹۵	۱۳۰۴	۲	۱	۱۰۹۵	۱۳۰۴
۳	۱۸۶۰	۲۰۲۲	۲	۱	۱۸۶۰	۲۰۲۲
۴	۱۲۲۷	۱۳۹۵	۲	۱	۱۲۲۷	۱۳۹۵
۵	۶۱۸	۷۷۳	۲	۱	۶۱۸	۷۷۳
۶	۱۴۹۹	۱۶۳۰	۲	۱	۱۴۹۹	۱۶۳۰
۷	۱۳۲۱	۱۴۳۶	۲	۱	۱۳۲۱	۱۴۳۶
۸	۰	۲۵۲۵	۲	۱	۰	۲۵۲۵
۹	۳۵۹۳	۴۶۷۴	۳	۰	۰	۰
۱۰	۱۷۳۲	۲۲۷۹	۲	۱	۱۷۳۲	۲۲۷۹

در بین ۱۰ سابقه اول نمونه که در جدول نشان داده شده‌اند فقط یک سابقه (مربوط به بیمارستان نمونه ۹) تحت تأثیر تبدیل قرار گرفته است، زیرا این سابقه در بین ۱۰ بیمارستان اول از تنها بیمارستانی گرفته شده است که در زیرحوزه موردنظر نیست. در کل نمونه متشکل از ۵۶ بیمارستان، ۱۳ بیمارستان از نظر خدمات زنان و زایمان در سطح «۲» قرار ندارند و تحت تأثیر تبدیل قرار می‌گیرند، ۴۳ بیمارستان دیگر در سطح «۲» قرار دارند. آماره‌های خلاصه برای نمونه متشکل از ۵۶ بیمارستان به شرح زیرند:

$$N = 101 \quad n = 56$$

$$x_{(k)} = \sum_{i=1}^{56} x_i \delta_{ik} = 43910 \quad y_{(k)} = \sum_{i=1}^{56} y_i \delta_{ik} = 58082$$

$$r_{(k)} = \frac{x_{(k)}}{y_{(k)}} = 0.7560$$

$$s_{x(k)}^2 = 547393/0.79 \quad s_{y(k)}^2 = 639484/913$$

$$\bar{x}_{(k)} = \frac{x_{(k)}}{56} = 784/107 \quad \bar{y}_{(k)} = \frac{y_{(k)}}{56} = 1037/179$$

$$\hat{V}_{x(k)} = \sqrt{\frac{100}{101} \frac{S_{x(k)}}{\bar{x}_{(k)}}} = 0.9389$$

$$\hat{V}_{y(k)} = 0.7672 \hat{\rho}_{x(k), y(k)} = 0.825$$

در این صورت، از فرمول تابلوی ۱.۷ خواهیم داشت

$$\hat{SE}(r_{(k)}) = 0.0360$$

به این ترتیب برآورد می‌شود که ۷۵/۶٪ همه نوزادان در بیمارستانهای سطح «۲» از نظر خدمات زنان و زایمان، وضعیت پادگن سطحی هیاتیت B مادرشان در پرونده پزشکی آنها ثبت شده و خطای معیار این برآورد ۰/۰۳۶۰ است.

□

۳.۷ برآوردهای نسبتی پس طبقه‌بندی شده تحت نمونه‌گیری تصادفی ساده

فن پس طبقه‌بندی که عموماً در روش‌شناسی نمونه‌گیری به کار می‌رود یکی از فنون برآورد کردن است که از اطلاعات معلوم جامعه در مورد زیرگروهها یا زیرحوزهها استفاده می‌کند تا برآوردهایی تولید کند که از بهبود دقت برخوردار باشند. در فصل ۶، بخش ۶.۶، برآوردهای پس طبقه‌بندی شده میانگینهای جامعه را تحت نمونه‌گیری تصادفی ساده شرح دادیم. مطالب مورد بحث در آن قسمت برای سایر برآوردهای خطی از قبیل مجموعها و نسبتها نیز به همان اندازه قابل اجراست. در این بخش، استفاده از پس طبقه‌بندی را در برآورد کردن نسبتی تحت نمونه‌گیری تصادفی ساده مورد بحث قرار خواهیم داد.

فرض کنیم یک نمونه تصادفی ساده n تایی از جامعه‌ای متشکل از N عنصر داریم و می‌خواهیم

$$\text{نسبت } R = \frac{X}{Y} \text{ را برآورد کنیم.}$$

برآورد نسبتی، r را که در قسمت ۱.۷ تصریح شد می‌توان به صورت زیر نوشت.

$$r = \frac{x'}{y'} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i y_i}$$

که در آن

$$w_i = \frac{N}{n}$$

اگر k زیرگروه یا زیرحوزه دو به دو ناسازگار و فراگیر وجود داشته باشند آنگاه فرمول زیر برای x' درست است:

$$x' = \sum_{k=1}^K \sum_{i=1}^n w_i \delta_{ik} x_i = \frac{N}{n} \sum_{k=1}^K x_{(k)}$$

که در آن

$$x_{(k)} = \sum_{i=1}^n \delta_{ik} x_i$$

و

$$\delta_{ik} = \begin{cases} 1 & \text{اگر عنصر } i \text{ در زیرحوزه } k \text{ باشد} \\ 0 & \text{در غیر صورت بالا} \end{cases}$$

نهاد $x_{(k)}$ به سادگی مجموع نمونه‌ای متغیر X برای همه عناصر نمونه است که در زیرحوزه k قرار دارند. به همین ترتیب، نهاد

$$n_{(k)} = \sum_{i=1}^n \delta_{ik}$$

تعداد عناصر نمونه است که در زیرحوزه k قرار دارند. نسبت $\frac{x_{(k)}}{n_{(k)}}$ سطح X را به ازای هر عنصر در حوزه k از روی نمونه برآورد می‌کند و اگر کل تعداد عناصر، $N_{(k)}$ ، در حوزه k ی جامعه معلوم باشد، در آن صورت نسبت $(N_{(k)}/n_{(k)})x_{(k)}$ برآورد «مشخص»تری از کل X در حوزه k خواهد بود. بالاخره، با استفاده از این «تصحیح» برای هر حوزه و نیز برای متغیر Y ، برآورد نسبتی پس طبقه‌بندی شده r_{pstr} را به دست می‌آوریم که در زیر نشان داده شده است:

$$r_{pstr} = \frac{\sum_{k=1}^K \frac{N_k}{n_{(k)}} x_{(k)}}{\sum_{k=1}^K \frac{N_k}{n_{(k)}} y_{(k)}}$$

که در آن

$$x_{(k)} = \sum_{i=1}^n \delta_{ik} x_i$$

$$y_{(k)} = \sum_{i=1}^n \delta_{ik} y_i$$

و

$$n_{(k)} = \sum_{i=1}^n \delta_{ik}$$

چون نهادهای $x_{(k)}$ ، $y_{(k)}$ و $n_{(k)}$ همه در معرض تغییرپذیری نمونه‌گیری قرار دارند نتیجه می‌گیریم که برآورد نسبتی پس طبقه‌بندی شده r_{pstr} در واقع نسبت مجموع برآوردهای نسبتی K شامل $x_{(k)}$ و $n_{(k)}$ تقسیم بر مجموع برآوردهای نسبتی K شامل $y_{(k)}$ و $n_{(k)}$ است. خطای معیار r_{pstr} صورت ساده‌ای ندارد و در اینجا ارائه نخواهد شد ولی می‌توان آن را با استفاده از خطی‌سازی سری تیلور (فصل ۱۲ را ببینید) برآورد کرد که در نرم‌افزار SUDAAN موجود است.

مثال تشریحی: نمونه متشکل از ۵۶ بیمارستان را در نظر می‌گیریم که در بخش ۲.۷ مورد بحث قرار گرفت (داده‌های حاصل از این نمونه در پرونده ASCII، BHRATIO.DAT قرار دارند). استفاده سراسر از فرمولهای برآورد کردن نسبتی تابلوی ۱.۷، برآوردی معادل $۰/۷۵۶$ با برآورد انحراف معیار $۰/۰۳۶$ را برای نسبت نوزادانی که وضعیت پادگن سطحی هپاتیت B مادرانشان در پرونده پزشکی آنها ثبت شده است نتیجه می‌دهد.

به علاوه، می‌دانیم که ۲۹ ($۵۱/۸\%$) بیمارستان از ۵۶ بیمارستان نمونه دارای سیاستهای مکتوبی هستند مبنی بر اینکه تا زمانی که وضعیت مادر از لحاظ پادگن سطحی هپاتیت B در پرونده نوزاد درج نشود نوزاد نباید از بیمارستان مرخص شود. همچنین می‌دانیم که ۶۶ بیمارستان ($۶۵/۴\%$) از ۱۰۱ بیمارستان موجود در جامعه متناظر که این نمونه از آن گرفته شده است این سیاست مدون را دارند. چون منطقی است که بیمارستانهای پیرو این سیاست دارای نسبت بالاتری از نوزادانی باشند که وضعیت پادگن سطحی هپاتیت B مادرانشان در پرونده‌های آنان درج می‌شود، و چون توزیع جامعه بیمارستانهایی که این سیاست را دارند و ندارند معلوم است، شرایطی فراهم شده است که می‌توان به صورتی معقول از برآورد نسبتی پس طبقه‌بندی شده استفاده کرد. زیرحوزه ۱ را زیرحوزه دارای سیاست مدون در نظر می‌گیریم که مستلزم داشتن اطلاعاتی درباره وضعیت پادگن سطحی هپاتیت B مادر در پرونده پزشکی نوزاد قبل از ترخیص از بیمارستان است و زیرحوزه ۲ را زیرحوزه بیمارستانهایی می‌گیریم که چنین سیاستی را دنبال نمی‌کنند. به این ترتیب آماره‌های خلاصه زیر را به دست می‌آوریم.

تعداد بیمارستانهای نمونه متعلق به زیرحوزه ۱	$n_{(1)} = 29$
تعداد بیمارستانهای نمونه متعلق به زیرحوزه ۲	$n_{(2)} = 27$
کل تعداد نوزادان در زیرحوزه ۱ بیمارستانهای نمونه که وضعیت پادگن سطحی هپاتیت B مادرانشان در پرونده پزشکی ثبت شده است	$x_{(1)} = 42749$
کل تعداد نوزادان در زیرحوزه ۲ بیمارستانهای نمونه که وضعیت پادگن سطحی هپاتیت B مادرانشان در پرونده پزشکی ثبت نشده است	$x_{(2)} = 22560$

کل تعداد نوزادان در نمونه زیرحوزه ۱ بیمارستانها	$y_{(1)} = 47850$
کل تعداد نوزادان در نمونه زیرحوزه ۲ بیمارستانها	$y_{(2)} = 38929$
تعداد جامعه معلوم بیمارستانهای زیرحوزه ۱	$N_1 = 66$
تعداد جامعه معلوم بیمارستانهای زیرحوزه ۲	$N_2 = 35$

از روی این آماره‌های خلاصه، برآورد نسبی پس طبقه‌بندی شده را به شرح زیر می‌سازیم:

$$r_{pstr} = \frac{\frac{66}{29} 42749 + \frac{35}{27} 22560}{\frac{66}{29} 47850 + \frac{35}{27} 38929} = 0.794$$

در زیر برآورد نسبی پس طبقه‌بندی شده با برآورد نسبی معمولی در مورد این مثال مقایسه شده است:

خطای معیار برآورد شده	برآورد نسبی	برآورد معمولی
$SE(r) = 0.036$	$r = 0.756$	معمولی
$SE(r_{pstr}) = 0.022$	$r_{pstr} = 0.794$	پس طبقه‌بندی شده

در این مثال، برآورد نسبی پس طبقه‌بندی شده آشکارا با برآورد نسبی معمولی تفاوت دارد و خطای معیار آن به مراتب کمتر است. دلیل این امر آن است که نسبتهای زیرحوزه‌ها تفاوت قابل ملاحظه‌ای دارند و نسبت بیمارستانهای زیرحوزه ۲ در نمونه بیشتر از آن است که در جامعه وجود دارد.

□

۴.۷ برآورد نسبی مجموعه‌ها تحت نمونه‌گیری تصادفی ساده

برآوردهای مجموعه‌ها یا انبوه‌ها که تا اینجا مورد بحث قرار گرفتند از ضرب ساده مجموع نمونه در نسبت $\frac{N}{n}$ به دست می‌آید. درباره این واقعیت بحث کردیم که برآوردهای x' از این نوع، برآوردهایی نارایب‌اند و فرمولی برای خطای معیار این قبیل برآوردها و روشی برای ساختن بازه‌های اطمینان تقریبی مجموع نامعلوم X تحت نمونه‌گیری تصادفی ساده شرح دادیم. ولی گاهی اوقات در موقعیتی قرار می‌گیریم که می‌توانیم از اطلاعات اضافی موجود برای مقاصد ساختن برآوردگر دیگری از مجموع کل جامعه استفاده کنیم که بر پایه برآورد نسبی متکی است، به نحوی که میانگین توان دوم خطای برآوردگر به دست آمده از برآوردگر توری ساده x' کمتر است. این ایده را با نگاهی به یک مثال بررسی می‌کنیم.

مثال تشریحی: فرض کنیم روستایی دارای شش ناحیه سرشماری است که جمعیت سال ۱۹۹۰ آن در جدول ۴.۷ ارائه شده است. میزان ثبت‌نام مدارس این روستا (که پیش از اجرای آمارگیری نامعلوم تلقی می‌شود) نیز در جدول ۴.۷ ارائه شده است.

فرض کنیم که می‌خواهیم کل ثبت‌نام شدگان فعلی مدارس را با انتخاب یک نمونه تصادفی ساده از دو ناحیه سرشماری و معلوم کردن ثبت‌نام شدگان در هر ناحیه نمونه با آمارگیری از هر یک از مدارس موجود در ناحیه، شمارش افراد، و ضرب مجموعهای نمونه‌ای در $\frac{N}{n}$ به طوری که در

فصل ۳ شرح داده شده است برآورد کنیم. $\binom{6}{2} = 15$ نمونه ممکن همراه با مقادیر x' ، برآورد مجموع، در جدول ۵.۷ فهرست شده‌اند.

جدول ۴.۷ جمعیت (سرشماری ۱۹۹۰) و ثبت‌نام شدگان فعلی مدارس برحسب ناحیه سرشماری

ناحیه سرشماری	جمعیت سال ۱۹۹۰	ثبت‌نام شدگان فعلی*
۱	۶۶۵۷	۲۲۶۹
۲	۴۰۵۷	۱۳۲۴
۳	۳۶۴۲	۹۵۲
۴	۵۳۲۰	۱۵۵۸
۵	۴۴۸۰	۱۳۵۲
۶	۵۸۸۰	۱۷۹۶
مجموع	۳۰۰۳۶	۹۲۵۱

* تا زمان آمارگیری از مدارس نامعلوم است.

جدول ۵.۷ نمونه‌های ممکن متشکل از دو مدرسه از روی داده‌های جدول ۴.۷

شماره مدرسه‌ها	برآورد ثبت‌نام شدگان	شماره مدرسه‌ها	برآورد ثبت‌نام شدگان
در نمونه	در مدارس	در نمونه	در مدارس
۱, ۲	۱۰۷۷۹	۲, ۶	۹۳۶۰
۱, ۳	۹۶۶۳	۳, ۴	۷۵۳۰
۱, ۴	۱۱۴۸۱	۳, ۵	۶۹۱۲
۱, ۵	۱۰۸۶۳	۳, ۶	۸۲۴۴
۱, ۶	۱۲۱۹۵	۴, ۵	۸۷۳۰
۲, ۳	۶۸۲۸	۴, ۶	۱۰۰۶۲
۲, ۴	۸۶۴۶	۵, ۶	۹۴۴۴
۲, ۵	۸۰۲۸		

از نتایجی که درباره نمونه‌گیری تصادفی ساده به دست آوردیم می‌دانیم که برآورد مجموع x' یک برآورد نااریب از کل جمعیت است. خطای معیار x' (که با شمارش همه نمونه‌های فهرست شده زیر محاسبه می‌شود) بنابر تابلوی ۳.۲ عبارت است از

$$SE(x') = 1568/4577$$

به جای استفاده از برآورد x' که در مثال بالا انجام دادیم می‌توانیم کل ثبت‌نام شدگان X را با استفاده از داده‌های مربوط به جمعیت سال ۱۹۹۰ که برای هر ناحیه سرشماری موجود است برآورد کنیم. هر نمونه‌ای از نواحی سرشماری به ما امکان می‌دهد که نه تنها ثبت‌نام مدارس را از روی نمونه برآورد کنیم بلکه کل جمعیت را نیز از نمونه برآورد کنیم. فرض کنید از نمادهای زیر استفاده شود:

$$\begin{aligned} Y & \text{ تعداد کل افراد در جامعه} \\ Y_i & \text{ تعداد افراد در ناحیه سرشماری } i \\ y_i & \text{ تعداد افراد در ناحیه نمونه } i \\ y = \sum_{i=1}^n y_i & \text{ کل تعداد افراد در نواحی سرشماری نمونه} \end{aligned}$$

پس برآورد کل جمعیت، y' ، از فرمول زیر به دست می‌آید

$$y' = \left(\frac{N}{n} \right) \times y$$

که در آن،

$$N = \text{کل تعداد نواحی در جامعه}$$

$$n = \text{کل تعداد نواحی در نمونه}$$

هر نمونه، برآورد x' را از X و y' را از Y به دست می‌دهد. ولی چون مقدار واقعی Y از روی داده‌های سرشماری برای ما معلوم است، پس معقول خواهد بود که در صورت وجود همبستگی شدید بین میزان ثبت‌نام مدارس و اندازه جمعیت، فرض را بر این بگیریم که برآوردگر x' از X با خود X به همان نسبتی تفاوت خواهد داشت که برآوردگر y' با Y تفاوت دارد. یعنی $\frac{X}{x'} = \frac{Y}{y'}$. این، انگیزه

استفاده از برآوردگر x'' برای X می‌شود که از فرمول

$$x'' = \left(\frac{x'}{y'} \right) \times Y \quad (4.7)$$

یا از معادل جبری آن

$$x'' = r \times Y \quad (5.7)$$

به دست می‌آید، که در آن r ، یک برآورد نسبتی است.

اگر این برآوردگر به صورت زیر نوشته شود شناخت بیشتری از آن آشکار خواهد شد:

$$x'' = \left(\frac{Y}{y'} \right) \times x'$$

اگر برآورد اندازه جمعیت، y' ، کمتر از مجموع واقعی معلوم Y باشد، آنگاه y' مقدار Y را کم برآورد می‌کند. در این صورت انتظار خواهیم داشت که برآورد x' که از همین نمونه به دست آمده است نیز تا حد مشابهی مجموع واقعی نامعلوم X را کم برآورد کند، زیرا x و y همبستگی دارند. ولی توجه کنید که نسبت $\frac{Y}{y'}$ در این مورد بیشتر از یک است و x'' از برآورد تورمی ساده x' بزرگتر خواهد بود. هرگاه y' بزرگتر از Y باشد جریان معکوس خواهد شد. به این ترتیب، می‌بینیم که حتی در مورد نمونه‌هایی که به طور خاص بدند (یعنی نمونه‌هایی که مقادیری از x' را نتیجه می‌دهند که از لحاظ قدر مطلق تفاوت فاحشی با X دارند) نسبت $\frac{Y}{y'}$ آن را اصلاح می‌کند و برآوردگر x'' به اندازه x' با X تفاوت نخواهد داشت.

اینک همه این ایده‌ها را با چند مثال بررسی می‌کنیم.



مثال تشریحی: برای نمونه‌های $n=2$ ناحیه سرشماری که از جامعه‌ای که در جدول ۴.۷ نشان داده شده گرفته شده است مقادیر نمونه‌ای x' و x'' را بررسی می‌کنیم. این مقادیر در جدول ۶.۷ فهرست شده‌اند.

از توزیع نمونه‌گیری که در جدول ۶.۷ نشان داده‌ایم می‌بینیم که میانگین $E(x'')$ ، خطای معیار

$SE(x'')$ و میانگین توان دوم خطا $MSE(x'')$ عبارت‌اند از (رجوع کنید به تابلوی ۳.۲)

$$E(x'') = \frac{10073 + 9394 + 000 + 9127}{15} = 9211.07$$

$$SE(x'') = \left[\frac{(10073 - 9211.07)^2 + (9394 - 9211.07)^2 + 000 + (9127 - 9211.07)^2}{15} \right]^{1/2} = 466.38$$

$$MSE(x'') = Var(x'') + B^2(x'') = 217512/33 + (9211.07 - 9251)^2 = 219106/73$$

جدول ۶.۷ مقادیر x' و x'' برای نمونه‌های جدول ۵.۷

x''	x'	نمونه
۱۰۰۷۳	۱۰۷۷۹	۱, ۲
۹۳۹۴	۹۶۶۳	۱, ۳
۹۵۹۷	۱۱۴۸۱	۱, ۴
۹۷۶۶	۱۰۸۶۳	۱, ۵
۹۷۳۹	۱۲۱۹۵	۱, ۶
۸۸۷۹	۶۸۲۸	۲, ۳
۹۲۳۱	۸۶۴۶	۲, ۴
۹۴۱۵	۸۰۲۸	۲, ۵
۹۴۳۱	۹۳۶۰	۲, ۶
۸۴۱۲	۷۵۳۰	۳, ۴
۸۵۲۰	۶۹۱۲	۳, ۵
۸۶۶۸	۸۲۴۴	۳, ۶
۸۹۱۹	۸۷۳۰	۴, ۵
۸۹۹۵	۱۰۰۶۲	۴, ۶
۹۱۲۷	۹۴۴۴	۵, ۶

چون برآورد تورمی ساده x' برآوردی نارایب از X است میانگین توان دوم خطای MSE آن برابر با واریانس آن است:

$$MSE(x') = Var(x') = (1568/4577)^2 = 2460059/56$$

به این ترتیب، با این که x'' برآوردی رایب برای X است میانگین توان دوم خطای آن در این مثال خیلی کمتر از x' است.

□

مثال تشریحی: از داده‌های جدول ۱.۵ نیز می‌توان برای تشریح برآورد کردن یک مجموع به وسیله یک نسبت استفاده کرد. اگر بخواهیم کل تعداد X کامیون - مایلهای طی شده را از یک نمونه تصادفی ساده از $n=3$ قسمت برآورد کنیم، می‌توانیم نسبت r کامیون - مایلهای طی شده را به تصادفاتی که در آنها کامیون مداخله داشته است از روی نمونه برآورد کنیم. اگر کل تعداد Y تصادفی که کامیون در آن مداخله داشته است برای کل جاده معلوم باشد می‌توان برآورد x'' را به دست آورد. برای مثال اگر قسمتهای ۱، ۳ و ۴ به عنوان نمونه انتخاب شوند محاسبات زیر را خواهیم داشت (رجوع کنید به تابلوی ۱.۳ و تابلوی ۱.۷):

برآورد تورمی ساده کامیون - مایل طی شده $x' = \left(\frac{\wedge}{\text{ن}}$

برآورد تورمی ساده تعداد تصادفات با مداخله کامیون $y' = \left(\frac{\wedge}{\text{ن}}$

برآورد نسبت کامیون - مایل طی شده به تصادفات با مداخله کامیون $r = \frac{x'}{y'} = ۸۷۸/۹۲۳۰ =$

کل تعداد تصادفات با مداخله کامیون، Y ، معلوم است، دقیقاً $Y = ۵۰$

پس

$$x'' = \left(\frac{x'}{y'}\right) \times Y = ۸۷۸/۹۲۳۰ \times ۵۰ = ۴۳۹۴۶/۱۵ \text{ مایل کامیون - مایل}$$

این رقم به مجموع واقعی ($X = ۳۴۰۵۴$) نزدیکتر است تا برآورد تورمی ساده x' . □

برآورد x'' به صورت rY است که در آن r ، نسبت برآورد شده و Y ، مجموع معلوم جامعه است. از این رو نتیجه می‌گیریم که خطای معیار $SE(x'')$ برابر است با $Y \times SE(r)$ که می‌توان آن را از فرمول زیر تقریب کرد، به شرطی که $V(\bar{y})$ ضریب تغییرات \bar{y} کمتر از ۰/۰۵ باشد:

$$SE(x'') = \left(\frac{Yr}{\sqrt{n}}\right) \times \left(V_x^2 + V_y^2 - 2\rho_{xy}V_xV_y\right)^{1/2} \times \sqrt{\frac{N-n}{N-1}} \quad (۶.۷)$$

به همین ترتیب، برآورد $\hat{SE}(x'')$ از $SE(x'')$ را می‌توان از روی داده‌ها با جایگزین کردن $\hat{SE}(r)$ [رابطه (۳.۷)] به جای $SE(r)$ در فرمول (۶.۷) به شرح زیر به دست آورد:

$$\hat{SE}(x'') = \left(\frac{Yr}{\sqrt{n}}\right) \times \left(\hat{V}_x^2 + \hat{V}_y^2 - 2\hat{\rho}_{xy}\hat{V}_x\hat{V}_y\right)^{1/2} \times \sqrt{\frac{N-n}{N-1}} \quad (۷.۷)$$

از فرمول (۷.۷) می‌توان در به دست آوردن بازه‌های اطمینان تقریبی برای یک مجموع نامعلوم X استفاده کرد، به شرطی که $\hat{V}(\bar{y})$ ضریب برآورد شده تغییرات \bar{y} کمتر از ۰/۰۵ باشد. به مثالی دیگر نگاه کنیم.

مثال تشریحی: اگر قسمتهای ۱، ۳ و ۴ از داده‌های جدول ۱.۵ نمونه‌گیری شده باشند داده‌های نمونه جدول ۷.۷ را در اختیار داریم. پس آماره‌های خلاصه زیر را خواهیم داشت:

$n = ۳$	$\bar{x} = ۷۶۱۷/۳۳$	$s_x = ۱۱۹۶/۸۰$	$\hat{V}_x^2 = ۰/۰۲۱۶$
$N = ۸$	$\bar{y} = ۸/۶۷$	$s_y = ۰/۵۷۷۴$	$\hat{V}_y^2 = ۰/۰۰۳۸۸۳$
$\hat{\rho}_{xy} = ۰/۹۳۳۷$	$Y = ۵۰$	$r = ۸۷۸/۹۲۳۱$	$\bar{x} = ۴۳۹۴۶/۱۵$

برآورد ضریب تغییرات، $\hat{V}(\bar{y})$ ، عبارت است از^۱

$$\hat{V}(\bar{y}) = \left(\frac{s_y}{\bar{y}\sqrt{n}} \right) \times \sqrt{\frac{N-n}{N}} = \left[\frac{0.05774}{1.67\sqrt{3}} \right] \times \sqrt{\frac{8-3}{8}} = 0.0304$$

چون $\hat{V}(\bar{y}) < 0.05$ می‌توان از عبارت (۷.۷) برای برآورد خطای معیار x'' به صورت زیر استفاده کرد:

$$\begin{aligned} \hat{SE}(x'') &= \left(\frac{50 \times 171/5848}{\sqrt{3}} \right) \times [0.0245 + 0.003883 - 2 \times 0.09337 \\ &\quad \times \sqrt{0.0216} \times \sqrt{0.003883}]^{1/2} \times \sqrt{\frac{8-3}{8-1}} = 1963/04 \end{aligned}$$

جدول ۷.۷ داده‌های نمونه از جدول ۱.۵

قسمتهای نمونه	تعداد کامیون - مایل (x_i)	تعداد تصادفات با مداخله کامیون (y_i)
۱	۶۳۲۷	۸
۳	۸۶۹۱	۹
۴	۷۸۳۴	۹

در این صورت بازه اطمینان ۹۵٪ را به شرح زیر برای کل تعداد برآورد شده کامیون - مایل طی شده خواهیم داشت:

$$\begin{aligned} x'' - 1/96 \times \hat{SE}(x'') \leq X \leq x'' + 1/96 \times \hat{SE}(x'') \\ 43946/15 - 1/96 \times 1963/04 \leq X \leq 43946/15 + 1/96 \times 1963/04 \\ 40098 \leq X \leq 47793/71 \end{aligned}$$

توجه کنید که در مورد این نمونه بخصوص، بازه اطمینان ۹۵٪، مجموع واقعی جامعه را نمی‌پوشاند.

□

^۱ فرمول مربوط به $V(\bar{y})$ در معادله (۷.۳) ارائه شده است. با جایگزین کردن (\hat{V}_y) در فرمول (۷.۳) به صورتی که در تابلوی ۱.۷ ارائه شده نتیجه زیر را به دست می‌آوریم

$$\hat{V}(\bar{y}) = \left(\frac{\hat{V}_y}{\sqrt{n}} \right) \times \sqrt{\frac{N-n}{N-1}} = \left[\sqrt{\frac{N-1}{N}} \times \left(\frac{s_y}{\bar{y}} \right)^2 \right]^{1/2} \times \frac{1}{\sqrt{n}} \times \sqrt{\frac{N-n}{N-1}} = \left(\frac{s_y}{\bar{y}\sqrt{n}} \right) \times \sqrt{\frac{N-n}{N}}$$

۵.۷ مقایسه برآورد نسبتی با برآورد تورمی ساده

فرض کنیم تقریب خطای معیار برآورد نسبتی مجموع، تحت نمونه‌گیری تصادفی ساده، که از رابطه (۶.۷) به دست می‌آید معتبر است و اریبی برآورد نسبتی را می‌توان نادیده گرفت. در این صورت می‌توانیم ارزیابی کنیم که آیا برآورد نسبتی x'' احتمال دارد به برآوردی از مجموع منجر شود که نسبت به برآورد تورمی ساده بهبودی بیشتری داشته باشد یا نه. این ارزیابی با بررسی نسبت واریانس x'' به واریانس x' به شرح زیر میسر است:

$$\frac{Var(x'')}{Var(x')} = \frac{V_x^r + V_y^r - 2\rho_{xy}V_xV_y}{V_x^r}$$

چون

$$\begin{aligned} Var(x'') &= Y^r \times Var(r) \\ &= Y^r \times \left(\frac{R^r}{n}\right) \times (V_x^r + V_y^r - 2\rho_{xy}V_xV_y) \times \left(\frac{N-n}{N-1}\right) \end{aligned}$$

و

$$\begin{aligned} Var(x') &= \left(\frac{N^r}{n}\right) \times \sigma_x^r \times \left(\frac{N-n}{N-1}\right) \\ &= \left(\frac{N^r}{n}\right) \times \left(\bar{X}^r V_x^r\right) \times \left(\frac{N-n}{N-1}\right) = \left(\frac{X^r}{n}\right) \times V_x^r \times \left(\frac{N-n}{N-1}\right) \end{aligned}$$

نتیجه می‌گیریم که واریانس x'' وقتی کمتر از واریانس x' خواهد بود که

$$V_x^r + V_y^r - 2\rho_{xy}V_xV_y < V_x^r$$

یا به طور هم‌ارز هرگاه

$$\frac{V_y}{V_x} < 2\rho_{xy} \quad (۸.۷)$$

به این ترتیب برای مقادیر ثابت ضرایب تغییرات X و Y در جامعه، هر چه همبستگی بین X و Y بیشتر باشد احتمال این که برآورد مجموع، x'' ، واریانسی کمتر نسبت به برآورد مجموع، x' ، یعنی برآورد تورمی ساده، داشته باشد بیشتر است.

۶.۷ تقریب خطای معیار برآورد نسبتی مجموع

هرگاه رابطه بین X و Y را بتوان به صورتی نسبتاً دقیق با یک خط راست که از مبدأ عبور کند نشان داد وضعیت خاصی پیش می‌آید که جالب توجه است. به عبارت دیگر، فرض کنید رابطه بین X و Y را بتوان با مدل رگرسیونی زیر بیان کرد

$$X_i = \beta Y_i + \varepsilon_i \quad (9.7)$$

که در آن β ، شیب خط رگرسیونی X و Y است و $\sum_{i=1}^N \varepsilon_i = 0$.

پس عبارت مربوط به ρ_{xy} به صورت زیر درمی آید

$$\begin{aligned} \rho_{xy} &= \frac{\sum_{i=1}^N (\beta Y_i + \varepsilon_i - \beta \bar{Y})(Y_i - \bar{Y}) / N}{\sigma_x \sigma_y} \\ &= \frac{\beta \sum_{i=1}^N (Y_i - \bar{Y})^2 + \sum_{i=1}^N \varepsilon_i (Y_i - \bar{Y})}{N \sigma_x \sigma_y} \end{aligned}$$

اگر عبارت $\sum_{i=1}^N \varepsilon_i (Y_i - \bar{Y})$ نزدیک به صفر باشد (که اگر مدل بالا به خوبی با داده‌ها برازش داشته باشد همین طور هم خواهد شد) آنگاه تقریباً برابر خواهد بود با عبارت

$$\rho_{xy} \approx \frac{\beta \sum_{i=1}^N (Y_i - \bar{Y})^2}{N \sigma_x \sigma_y}$$

چون از عبارت (۹.۷) نتیجه می‌گیریم که $\beta = \frac{\bar{X}}{\bar{Y}}$ و چون

$$\sum_{i=1}^N (Y_i - \bar{Y})^2 = N \sigma_y^2$$

داریم،

$$\rho_{xy} \approx \frac{\bar{X} \sigma_y}{\bar{Y} \sigma_x} = \frac{V_y}{V_x} \quad (10.7)$$

این نتیجه، از رابطه (۶.۷) تقریب خطای معیار x'' ، برآورد مجموع حاصل از نسبت را به ما می‌دهد:

$$SE(x'') \approx \left(\frac{YR}{\sqrt{n}} \right) \times (V_x^2 - V_y^2)^{1/2} \sqrt{\frac{N-n}{N-1}} \quad (11.7)$$

که می‌تواند از روی داده‌ها با فرمول زیر برآورد شود

$$\hat{SE}(x'') \approx \left(\frac{Yr}{\sqrt{n}} \right) \times (\hat{V}_x^2 - \hat{V}_y^2)^{1/2} \sqrt{\frac{N-n}{N-1}} \quad (12.7)$$

عبارتهای بالا بخصوص از این جهت سودمندند که شامل ضریب همبستگی به صورت صریح آن نیستند.

در مثال قبلی که در آن سه قسمت به عنوان نمونه انتخاب شدند، با فرض این که رگرسیون، خطی

است و از مبدأ می‌گذرد، برآورد خطای معیار یعنی، $\hat{SE}(x'')$ که از تقریب (۱۲.۷) به دست می‌آید عبارت است از

$$\hat{SE}(x'') = \left[\frac{50 \times 1878 / 5848}{\sqrt{3}} \right] (0.0245 - 0.004435) \sqrt{\frac{8-3}{8-1}} = 3.36/33$$

که از خطای معیار حاصل از رابطه (۷.۷) بیشتر است.

۷.۷ تعیین اندازه نمونه

برای این که $(1-\alpha) \times 100\%$ مطمئن باشیم که r ، برآورد نسبت یا برآورد نسبتی مجموع، x'' ، در داخل محدوده $100 \times \varepsilon\%$ از مقدار واقعی R (یا X) قرار دارد، تقریب اندازه نمونه مورد نیاز از فرمول زیر به دست می‌آید

$$n = \frac{z_{1-(\alpha/2)}^2 \times N \times (V_x^2 + V_y^2 - 2\rho_{xy} V_x V_y)}{z_{1-(\alpha/2)}^2 \times (V_x^2 + V_y^2 - 2\rho_{xy} V_x V_y) + (N-1)\varepsilon^2} \quad (13.7)$$

اگر پارامترهای عبارت (۱۳.۷) را از روی داده‌های اولیه بتوان حدس زد یا برآورد کرد، می‌توان برآورد اندازه نمونه را به دست آورد.

۸.۷ برآورد رگرسیونی مجموعهها

برآورد x'' از مجموع X بر مبنای یک نسبت r ، موردی خاص از برآورد رگرسیونی مجموع است. برآورد رگرسیونی x''' دارای صورتی است که از فرمول زیر به دست می‌آید

$$x''' = x' + b(Y - y') \quad (14.7)$$

که در آن b ، برآورد شیب و x' و y' ، برآوردهای تورمی ساده‌اند. برآورد شیب به صورت زیر است

$$b = \hat{\rho}'_{xy} \frac{\hat{SE}(x')}{\hat{SE}(y')} \quad (15.7)$$

که در آن $\hat{\rho}'_{xy}$ ، برآورد ضریب همبستگی بین x' و y' ، برآورد مجموعههاست. این فرمول برای نمونه‌گیری تصادفی ساده به صورت زیر تبدیل می‌شود.

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (16.7)$$

برآوردی از واریانس برآورد رگرسیونی x''' که برای نمونه‌های با اندازه‌های بزرگ، تحت نمونه‌گیری تصادفی ساده، معتبر است از فرمول زیر به دست می‌آید

$$\hat{Var}(x''') = \hat{Var}(x') \times (1 - \hat{\rho}_{xy}^2) \quad (17.7)$$

این برآورد تحت شرایطی خاص واریانسی کمتر از برآورد نسبتی خواهد داشت. مثال زیر، کاربرد این برآورد را نشان خواهد داد.

مثال تشریحی: شهری دارای چهار محله تعریف شده جغرافیایی است. در سال ۱۹۸۵ یک آمارگیری از تأمین‌کنندگان مراقبتهای بهداشتی، از جمله بیمارستانها و تسهیلات مراقبت دراز مدت اجرا شد تا برآوردی از تعداد افرادی به دست آید که در طی سال مذکور به هیأتیت وابسته به استفاده از مواد مخدر تزریقی دچار شده بودند. مایل بودند برآورد مشابهی برای سال ۱۹۹۰ در سراسر شهر به دست آورند، ولی بودجه‌ای فقط برای اجرای بررسی مزبور در نمونه‌ای متشکل از دو محله از چهار محله شهر موجود بود. فرض کنیم داده‌های واقعی (نامعلوم) جامعه به شرح زیر است:

تعداد افراد دچار هیأتیت وابسته به تزریق مواد مخدر

محلّه	۱۹۹۰ (X_i)	۱۹۸۵ (Y_i)
۱	۴۰	۳۰
۲	۱۰۳	۶۰
۳	۱۱۵	۷۰
۴	۲۳	۲۱
مجموع	۲۸۱	۱۸۱

اگر محله‌های ۲ و ۳ انتخاب شده باشند، برآورد رگرسیونی x''' به شرح زیر به دست خواهد آمد:

$$x' = 2 \times (103 + 115) = 436$$

$$y' = 2 \times (60 + 70) = 260$$

$$\bar{x} = \frac{103 + 115}{2} = 109$$

$$\bar{y} = \frac{60 + 70}{2} = 65$$

$$b = \frac{(103 - 109)(60 - 65) + (115 - 109)(70 - 65)}{(60 - 65)^2 + (70 - 65)^2} = 1/20$$

$$Y = 181$$

و

$$x''' = x' + b(Y - y') = 436 + 1/20(181 - 260) = 341/2$$

برآورد نسبتی x'' به صورت زیر به دست می‌آید

$$x'' = \left(\frac{x'}{y'} \right) \times Y = \left(\frac{436}{260} \right) \times 181 = 303/5$$

توزیعهای برآورد رگرسیونی x''' ، برآورد نسبتی x'' ، و برآورد تورمی ساده x' برای شش نمونه ممکن در پایین نشان داده شده است:

نمونه	برآورد تورمی ساده x'	برآورد نسبتی x''	برآورد رگرسیونی x'''
۱, ۲	۲۸۶	۲۸۷/۵۹	۲۸۸/۱۰
۱, ۳	۳۱۰	۲۸۰/۵۵	۲۷۴/۳۸
۱, ۴	۱۲۶	۲۲۳/۵۹	۲۷۵/۲۲
۲, ۳	۴۳۶	۳۰۳/۵۲	۳۴۱/۲۰
۲, ۴	۲۵۲	۲۸۱/۵۶	۲۹۰/۹۷
۳, ۴	۲۷۶	۲۷۴/۴۸	۲۷۴/۱۲

از توزیعهای بالا، واریانسها، اریبیها، و میانگین توان دوم خطاها (MSEs) برای برآوردهای سه گانه عبارتند از:

برآورد	واریانس	2 (اریبی)	میانگین توان دوم خطا
x'	۸۲۹۷	۰	۸۲/۹۷
x''	۶۱۹/۸۳	۳۳/۴۷	۶۴۸/۳۰
x'''	۵۵۶/۳۳	۹۳/۴۳	۶۴۹/۷۶

می توان دید که در این مثال هم برآورد رگرسیونی x''' و هم برآورد نسبتی x'' دارای میانگین توان دوم خطایی هستند که به مراتب از میانگین توان دوم خطای برآورد تورمی ساده x' کمتر است. در این مورد، برآورد مجموع رگرسیونی x''' علی رغم دارا بودن بیشترین اریبی، دارای میانگین توان دوم خطایی است که با میانگین توان دوم خطای مجموع حاصل از برآورد نسبتی تقریباً یکسان است. □

برآورد مجموع نسبتی x'' و برآورد مجموع رگرسیونی x''' هر دو از رابطه بین متغیر مورد نظر x و یک متغیر کمکی y در به دست آوردن برآوردی دقیقتر برای مجموع واقعی X استفاده می کنند. برآورد مجموع رگرسیونی در واقع صورت تعمیم یافته برآورد مجموع نسبتی است و وقتی ضریب رگرسیونی b برابر با $\frac{x'}{y'}$ باشد به همان صورت برآورد مجموع نسبتی تبدیل می شود.

۹.۷ برآورد نسبی در نمونه‌گیری تصادفی طبقه‌بندی شده

وقتی طرح نمونه‌گیری، نمونه‌گیری تصادفی طبقه‌بندی شده است، برای برآورد نسبتها معمولاً از دو روش زیر استفاده می‌شود، یعنی از برآورد نسبی ترکیبی r_{strc} و برآورد نسبی تفکیکی r_{strs} . این دو روش را در زیر شرح می‌دهیم.

برآورد نسبی ترکیبی، r_{strc}

$$r_{strc} = \frac{\bar{x}_{str}}{\bar{y}_{str}} \quad (18.7)$$

که در آن \bar{x}_{str} و \bar{y}_{str} میانگینهای برآورد شده‌ای هستند که برای نمونه‌گیری تصادفی طبقه‌بندی شده که در فصل ۶ شرح داده شد مناسب‌اند. واریانس این برآورد از رابطه زیر تقریب زده می‌شود.

$$Var(r_{strc}) \approx \left(\frac{1}{N^2 \bar{y}^2} \right) \sum_{h=1}^L \frac{N_h^2 (N_h - n_h)}{n_h (N_h - 1)} \sigma_{hz}^2 \quad (19.7)$$

که در آن،

$$\sigma_{hz}^2 = \sigma_{hx}^2 + R^2 \sigma_{hy}^2 - 2R\rho_{hxy} \sigma_{hx} \sigma_{hy}$$

و σ_{hx}^2 و σ_{hy}^2 در داخل محدوده واریانسهای طبقه‌ای X و Y قرار دارند و ρ_{hxy} ضریب همبستگی بین X و Y است.

برآورد نسبی تفکیکی r_{strs}

$$r_{strs} = \frac{\sum_{h=1}^L x_h''}{Y} \quad (20.7)$$

که در آن

$$x_h'' = Y_h \times \frac{\bar{x}_h}{\bar{y}_h}$$

\bar{x}_h و \bar{y}_h ؛ $h = 1, \dots, L$ برآورد میانگینهای ویژه طبقه هستند

Y_h ؛ $h = 1, \dots, L$ مجموعهای معلوم ویژه طبقه برای Y هستند

و

$$Y = \sum_{h=1}^L Y_h$$

واریانس تقریبی آن عبارت است از

$$Var(r_{strs}) \approx \left(\frac{1}{Y^2} \right) \sum_{h=1}^L Var(x_h'') \quad (21.7)$$

که در آن، $Var(x_h'')$ برای $h=1, \dots, L$ ، واریانس برآورد مجموع نسبتی در داخل طبقه h به شرحی است که در بخش ۲.۷ مورد بحث قرار گرفت.

برآورد نسبتی تفکیکی، r_{strs} ، برای این که بتواند مورد استفاده قرار گیرد نیازمند شناخت مجموعهای Y_h طبقه است. شرایطی که طی آن یکی از این روشها بر دیگری ترجیح دارد پیچیده است و با جزئیات بیشتر، هنس و همکاران [۱] و کوکران [۲] از آن بحث کرده‌اند. به طور کلی، هنگامی که نسبتهای درون طبقه‌ای، R_h ، بیش از حد متنوع نباشند فایده‌چندانی بر استفاده از برآورد نسبتی تفکیکی مترتب نخواهد بود. برآوردهایی از واریانسهای برآورد نسبتی ترکیبی و برآورد نسبتی تفکیکی هر دو را می‌توان با جایگزین کردن آماره‌های نمونه‌ای مناسب برای پارامترهای نشان داده شده در رابطه (۱۹.۷) و رابطه (۲۱.۷) ساخت.

مثال تشریحی: یک دانشگاه بزرگ به کارکنان خود حق انتخاب می‌دهد تا به عنوان بخشی از برنامه مزایای کارکنان به یکی از دو سازمان حفظ بهداشت (HMO) بپیوندند. برای تعیین هزینه‌ای که دانشگاه باید به ازای هر برخورد با بیمار متقبل شود، یک نمونه تصادفی طبقه‌بندی شده متشکل از ۲۰۰ نفر در داخل هر یک از این دو سازمان HMO گرفته شد و داده‌های زیر به دست آمد:

Y_h	$\hat{\rho}_{hxy}$	$\hat{\sigma}_{hy}$	\bar{y}_h	$\hat{\sigma}_{hx}$	\bar{x}_h	n_h	N_h	HMO
۶۹۰	۰/۶۹	۰/۹	۳/۲	۵۶۴ دلار	۱۳۴۶ دلار	۲۰۰	۱۲۹۵	۱
۴۷۲	۰/۶۵	۰/۶	۲/۵	۳۴۵ دلار	۷۸۵ دلار	۲۰۰	۲۳۱۰	۲

در این فهرست

$$N_h = \text{تعداد کارکنان دارای اشتراک } HMO_h (h=1,2)$$

$$n_h = \text{تعداد کارکنان نمونه‌گیری شده در } HMO_h (h=1,2)$$

$$\bar{x}_h = \text{برآورد میانگین هزینه سالانه به ازای هر نفر در } HMO_h (h=1,2)$$

$$\hat{\sigma}_{hx} = \text{برآورد انحراف معیار توزیع میانگین هزینه سالانه به ازای هر نفر در } HMO_h (h=1,2)$$

$$\bar{y}_h = \text{برآورد تعداد ویزیت‌های سالانه به ازای هر نفر در } HMO_h (h=1,2)$$

$$\hat{\sigma}_{hy} = \text{برآورد انحراف معیار توزیع تعداد ویزیتها در سال به ازای هر نفر در } HMO_h (h=1,2)$$

$$\hat{\rho}_{hxy} = \text{برآورد همبستگی بین هزینه سالانه به ازای هر نفر و ویزیت‌های سالانه به ازای هر نفر در } HMO_h (h=1,2)$$

$$Y_h = \text{کل ویزیت‌های سالانه در بین همه افراد در } HMO_h (h=1,2)$$

برآورد نسبتی تفکیکی، r_{strs} ، به شرح زیر محاسبه می‌شود:

$$r_{strs} = \left(\frac{1}{1162} \right) \left(\frac{1346}{3/2} \times 690 + \frac{785}{2/5} \times 472 \right) = 377/31$$

از رابطه‌های (۱۲.۷) و (۲۱.۷) خطای معیار آن به صورت زیر برآورد می‌شود

$$\hat{SE}(r_{strs}) = 5/72$$

برآورد نسبتی ترکیبی، r_{strc} ، از فرمول زیر به دست می‌آید

$$r_{strc} = \frac{\bar{x}_{str}}{\bar{y}_{str}} = \frac{986/5243}{2/751456} = 358/55$$

و از رابطه (۱۹.۷) خطای معیار آن به صورت زیر برآورد می‌شود

$$\hat{SE}(r_{strc}) = 14/92$$

در این مثال بخصوص، چون نسبتهای تکی طبقه متنوع است ($r_1 = 420/62$ ، $r_2 = 314/00$) برآورد نسبتی تفکیکی بر برآورد نسبتی ترکیبی ترجیح داده می‌شود و این امر در خطاهای معیار برآورد شده این دو برآورد بازتاب یافته است.

۱۰.۷ خلاصه

در این فصل، مفاهیم اساسی برآورد نسبتی را شرح و بسط دادیم که فنی است که غالباً برای برآورد نسبتها و نیز مجموعه‌های جامعه‌ای مورد استفاده قرار می‌گیرد. برآورد کردن نسبتی را می‌توان با هر برنامه نمونه‌گیری به کار برد و هرگاه برای برآورد مجموعه‌های جامعه‌ای به کار رود غالباً برآوردهایی تولید می‌کند که میانگین توان دوم خطای آن کمتر از آن است که از برآوردهای تورمی ساده تولید می‌شود. در این فصل برآورد نسبتی تحت نمونه‌گیری تصادفی ساده مورد تأکید خاصی قرار گرفت. توزیعهای نمونه‌گیری برآوردهای نسبتی را برای این نوع نمونه‌گیری مورد بحث قرار دادیم و روشهایی برای برآورد کردن خطاهای معیار، برای تعیین اندازه‌های نمونه، و برای ساختن بازه‌های اطمینان ارائه دادیم. بالاخره، برآورد مجموع رگرسیونی را مورد بحث قرار دادیم که از این جهت شبیه برآورد مجموع نسبتی است که از اطلاعات کمکی برای تولید برآورد دقیقتری از مجموع جامعه استفاده می‌کند.

تمرین

۱.۷ یک آمارگیری نمونه‌ای در دست برنامه‌ریزی است که می‌خواهند با آن، متوسط نسبت هزینه‌های درمانی به درآمد خانواده در یک شهر بزرگ با ۲۳۴۷۸۵ خانواده را برآورد کنند. براساس داده‌های همراه با تمرین ۲.۷، اگر برآورد این پارامتر با ۹۵٪ اطمینان در داخل محدوده ۵٪ مقدار واقعی آن موردنظر باشد چند خانواده باید برای نمونه انتخاب شوند؟

۲.۷ جدول همراه با این تمرین که براساس یک نمونه تصادفی ساده متشکل از ۳۳ خانواده از جامعه‌ای شامل ۶۰۰ خانواده گرفته شده است اندازه خانواده، درآمد هفتگی خالص خانواده، و مخارج هفتگی هزینه‌های درمانی شامل دارو را ارائه می‌دهد. این جامعه شامل ۲۷۰۰ نفر است.

الف. متوسط هزینه درمان هفتگی هر خانواده را با بازه اطمینان ۹۵٪ برآورد کنید.

ب. متوسط هزینه درمان هفتگی هر نفر را با بازه اطمینان ۹۵٪ برآورد کنید. روش برآورد کردن را توجیه کنید.

پ. متوسط نسبت درآمد خانواده را که برای مخارج درمانی هزینه شده است با بازه اطمینان ۹۵٪ برآورد کنید.

ت. کل هزینه‌های درمانی هفتگی پرداخت شده توسط جامعه را با بازه اطمینان ۹۵٪ برآورد کنید.

ث. متوسط نسبت درآمد خانواده را که برای مخارج درمانی هزینه شده است با بازه اطمینان ۹۵٪ برای خانواده‌هایی که درآمد خالص هفتگی آنان کمتر از ۴۰۰ دلار است برآورد کنید. همین کار را برای خانواده‌هایی که درآمد خالص هفتگی آنان بیشتر از ۴۰۰ دلار است انجام دهید.

شماره خانواده	اندازه خانواده	درآمد هفتگی خالص (دلار)	هزینه‌های درمانی در هفته (دلار)
۱	۲	۳۷۲	۲۸/۶۰
۲	۳	۳۷۲	۴۱/۶۰
۳	۳	۵۲۲	۴۵/۴۰
۴	۵	۳۹۰	۶۱/۰۰
۵	۴	۳۴۸	۸۲/۴۰
۶	۷	۵۵۲	۵۶/۴۰

۴۸/۴۰	۵۲۸	۲	۷
۶۰/۱۰۰	۵۷۴	۴	۸
۴۸/۴۰	۴۹۸	۲	۹
۸۸/۸۰	۳۷۲	۵	۱۰
۲۶/۸۰	۳۷۸	۳	۱۱
۳۹/۶۰	۳۷۲	۶	۱۲
۵۸/۸۰	۳۶۰	۴	۱۳
۵۴/۲۰	۴۵۰	۴	۱۴
۴۴/۲۰	۵۴۰	۲	۱۵
۷۵/۴۰	۴۵۰	۵	۱۶
۴۵/۲۰	۴۱۴	۳	۱۷
۷۲/۱۰۰	۴۹۸	۴	۱۸
۲۱/۲۰	۵۱۰	۲	۱۹
۵۵/۴۰	۴۳۸	۴	۲۰
۵۱/۹۰	۳۹۶	۲	۲۱
۴۶/۶۰	۳۴۸	۵	۲۲
۷۹/۶۰	۴۶۲	۳	۲۳
۳۳/۶۰	۴۱۴	۴	۲۴
۷۵/۶۰	۳۹۰	۷	۲۵
۶۹/۶۰	۴۶۲	۳	۲۶
۵۷/۴۰	۴۱۴	۳	۲۷
۱۲۶/۱۰۰	۵۷۰	۶	۲۸
۳۹/۱۰	۴۶۲	۲	۲۹
۴۳/۲۰	۴۱۴	۲	۳۰
۳۶/۴۰	۴۱۴	۶	۳۱
۴۰/۲۰	۴۰۲	۴	۳۲
۴۱/۴۰	۳۷۸	۲	۳۳

۳.۷ برای شهر مذکور در تمرین ۲.۷، می‌خواهند برآورد کل مبلغی را که برای مخارج درمانی هزینه شده است نیز با اطمینان قطعی در محدوده ۱۰ درصد مقدار واقعی آن به دست آورند. اگر قرار باشد این ویژگی تأمین شود، چند خانواده باید برای نمونه انتخاب شوند؟

۴.۷ برای داده‌های جدول همراه با تمرین ۲.۷، برآورد رگرسیونی از کل هزینه‌های درمانی هفتگی پرداخت شده توسط جامعه را به دست آورید. براساس این برآورد رگرسیونی، برای مقدار واقعی این پارامتر بازه‌های اطمینان ۹۵٪ ارائه دهید.

۵.۷ مطلوب است برآورد مجموع X ساعتی که یک پرستار متخصص برای مراقبت مستقیم از بیماران در یک سازمان بزرگ حفظ بهداشت طی سال گذشته صرف کرده است. این کار باید با گرفتن یک نمونه تصادفی ساده از بیماران و تعیین تعداد ساعتهایی که پرستار متخصص برای هر ویزیت در طی سال صرف کرده است انجام شود. می‌دانیم که تعداد اعضای سازمان حفظ بهداشت ۳۵۲۴ نفر و تعداد ویزیت‌های انجام شده در آن سال ۸۹۵۰ مورد است. نمونه‌های مقدماتی کوچکی از ۱۰ بیمار، داده‌های زیر را نتیجه داده است:

بیمار	ویزیت	ساعتهای صرف شده توسط پرستار متخصص
۱	۰	۰
۲	۵	۳
۳	۱	۰
۴	۲	۶
۵	۳	۳
۶	۷	۰
۷	۱	۰
۸	۱	۲
۹	۰	۰
۱۰	۴	۳

براساس این داده‌های مقدماتی، برآورد تورمی ساده را به عنوان روش برآورد کردن برای این آمارگیری نمونه‌ای توصیه می‌کنید یا برآورد نسبتی را؟ دلایل خود را برای این توصیه با مدرک ثابت کنید.

۶.۷ براساس داده‌های تمرین ۵.۷، اگر قرار شود برآورد تورمی ساده به کار رود چند بیمار باید نمونه‌گیری شوند؟

۷.۷ براساس داده‌های تمرین ۵.۷، اگر قرار شود برآورد نسبتی به کار رود چند بیمار باید نمونه‌گیری شوند؟

۸.۷ یک آمارگیری نمونه‌ای قرار است اجرا شود که به وسیله آن می‌خواهند نسبت اشخاص بالاتر از ۷۰ ساله‌ای را که دارای نشانه‌هایی از اختلالات شناختی هستند برآورد کنند. این کار قرار است با گرفتن نمونه‌ای از خانوارها و آزمون ساده کارکرد شناختی همه اعضای بالاتر از ۷۰ سال خانوار انجام شود. یک بررسی مقدماتی از ۲۵ خانوار، به طور متوسط ۱/۲ نفر ۷۰ سال به بالا به ازای هر خانوار با انحراف معیاری برابر ۰/۸ و به طور متوسط ۰/۲۴ نفر ۷۰ سال به بالا به ازای هر خانوار با نشانه‌های اختلالات شناختی با انحراف معیاری برابر ۰/۷۶ را نتیجه داده است. می‌دانیم که این جامعه دارای ۳۰۵۸ خانوار است و ۲۹۴۹ نفر در آن هستند که بیشتر از ۷۰ سال سن دارند. اگر بخواهیم برآورد نسبت اشخاص ۷۰ سال به بالا که علائمی از اختلالات شناختی از خود نشان داده‌اند با ۹۵٪ اطمینان در محدوده ۲۰٪ مقدار واقعی باشد چند خانوار باید نمونه‌گیری شوند؟

۹.۷ نشان دهید که در نمونه‌گیری تصادفی ساده، همبستگی ρ'_{xy} بین برآورد مجموعها، x' و y' ، از مشخصه‌های x و y برابر است با ρ'_{xy} ، همبستگی بین x و y .

۱۰.۷ تمرین زیر در مقاله‌ای که در نشریه گزارشهای بهداشت عمومی [۳]، به چاپ رسیده پیشنهاد شده است. در یک مرکز انتقال خون واقع در یک شهر بزرگ برای تعیین میزان شیوع تغییرات از لحاظ HIV (ایدز) در سرم خون کسانی که برای اولین بار خون اهدا می‌کنند یک آمارگیری اجرا شد. در طول یک ماه خاص ۱۸۰ نفر برای اولین بار در این مرکز خون دادند. نتیجه آزمایش خون ۱۷۵ نفر از این تعداد منفی بود. نمونه‌ای متشکل از ۶۰ نفر از افرادی که نتیجه اولین خون‌دهی آنان منفی بود انتخاب و طی شش ماه بعد به آنان وقت داده شد تا دوباره برای خون دادن مراجعه کنند. از این اشخاص داده‌های زیر به دست آمد.

وضعیت از لحاظ HIV		تعداد ماهها از اولین اهدای خون	
منفی	مثبت	تعداد اشخاص	
۱۶	۲	۱۸	۳
۸	۲	۱۰	۴
۷	۱	۸	۵
۷	۰	۷	۶
۴	۲	۶	۷
۵	۰	۵	۸
۳	۱	۴	۹
۱	۱	۲	۱۰

از روی این داده‌ها نرخ شیوع سالانه تغییرات سرمی از لحاظ HIV را در میان اهداکنندگان خون که برای اولین بار به مرکز مراجعه می‌کنند برآورد کنید. برای این نرخ شیوع بازه اطمینان ۹۵ درصد ارائه دهید.

۱۱.۷ یک کارخانه بزرگ ۱۰۰۰ کارگر دارد. یک نمونه تصادفی ساده متشکل از ۲۵ کارگر گرفته شده است تا نسبت روزهای غیبت از کار به کل روزهای اشتغال در طول سال تقویمی گذشته به دست آید:

کل روزهای غیبت از کار در میان ۲۵ فرد نمونه‌گیری شده: ۲۵۰

کل روزهای اشتغال در میان ۲۵ فرد نمونه‌گیری شده: ۴۳۷۵

الف. برآورد نسبت روزهای غیبت از کار به کل روزهای اشتغال در بین کارگران این کارخانه چقدر است؟

ب. برای برآورد خطای معیار این برآورد به چه اطلاعات دیگری نیاز خواهید داشت؟

پ. از روی اطلاعات بالا کل تعداد روزهای غیبت از کار را در بین ۱۰۰۰ کارگر این گروه برآورد کنید.

ت. چنین رخ داده است که این ۱۰۰۰ کارگر در طی سال تقویمی گذشته در مجموع ۱۴۰۰۰۰ روز در این کارخانه مشغول به کار بوده‌اند. براساس این اطلاع، درباره برآورد کل روزهای غیبت از کار که در بخش پ محاسبه شد چه استنباط می‌کنید؟ برای پاسخ خود دلیل بیاورید.

۱۲.۷ در مورد مثال تشریحی قسمت ۳.۷ فرض کنید تعداد معلوم بیمارستانها در زیرحوزه ۱ برابر با ۵۰ و در زیرحوزه ۲ برابر با ۵۱ بیمارستان باشد.

الف. برآورد نسبتی پس طبقه‌بندی شده نسبت نوزادانی را تعیین کنید که با درج وضعیت پادگن سطحی هپاتیت B مادرانشان در پرونده آنها از بیمارستان مرخص شده‌اند.

ب. دلایلی ارائه دهید که چرا این برآورد از برآورد پس طبقه‌بندی شده که در مثال تشریحی به دست آمد به برآورد معمولی این نرخ نزدیکتر است؟

کتابشناسی

The following texts contain more detailed discussions of ratio and regression estimation.

1. Hansen, M. H., Hurwitz, W. N., and Madow, W. G. *Sample Survey Methods and Theory*, Vols. 1 and 2, Wiley, New York, 1953.
2. Cochran, W. G., *Sampling Techniques*, 2nd ed, Wiley, New York, 1962.

The following article gives an example of the use of a ratio estimate in an HIV seroconversion survey.

3. Petersen, L. R., Dodd, R., and Dondero, T. J., Jr., Methodological approaches to surveillance of HIV infection among blood donors, *Public Health Reports* 105: 153, 1990.

The following three review articles present an overview of ratio and regression estimators along with a long list of references including "classic" articles.

4. Rao, P. S. R. S., Ratio and regression estimators. In *Handbook of Statistics*, Krishnaian, P. R. and Rao, C. R., Eds., Vol. 6, *Sampling*, Chap. 18, pp 449-468, Elsevier, Amsterdam and New York, 1988.
5. Rao, J. N. K., Ratio estimators. In *Encyclopedia of Statistical Sciences*, Kotz, S., and Johnson, N. L., Eds., Vol. 7, pp. 639-646, New York, Wiley, 1986.
6. Cumberland, W. G., Ratio and regression estimators. In *Encyclopedia of Biostatistics*, Armitage, P. and T., Colton, Ed., Wiley, Chichester, U.K., 1998.