

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 ۱۲

برآورد کردن واریانس در آمارگیریهای نمونه‌ای پیچیده

در فصلهای قبل، روشهای برآورد واریانسهای مشخصه‌های برآورد شده بیشتر طرحهای نمونه‌گیری گوناگون را که مورد بحث قرار گرفته بودند معرفی کردیم. در بسیاری از این طرحها، واریانس نظری برآوردگری ویژه، تابعی خطی از پارامترهای جامعه است. مثلاً در نمونه‌گیری تصادفی ساده، واریانس یک مجموع برآورد شده، x' ، از فرمول زیر به دست می‌آید

$$Var(x') = \frac{N^2}{n} \sigma_x^2 \left(\frac{N-n}{N-1} \right)$$

واضح است که چون عبارت بالا تابعی خطی از واریانس جامعه‌ای σ_x^2 است با جایگزین کردن برآورد نااریب $\hat{\sigma}_x^2$ از σ_x^2 در عبارت مربوط به $Var(x')$ که در فصل ۲ ارائه شد، می‌توان برآوردی نااریب به دست آورد.

ولی در بسیاری از طرحهای نمونه‌ای که عملاً به کار برده می‌شوند، فرایند برآورد کردن ممکن است مستلزم طبقه‌بندی، چندین مرحله نمونه‌گیری خوشه‌ای، برآورد نسبتی یا رگرسیونی، پس طبقه‌بندی برای مجموعهای معلوم، و سایر شیوه‌ها باشد و به این نتیجه بینجامد که واریانس برآورد حاصل، یک تابع خطی یا حتی تابعی معلوم از پارامترهای جامعه نباشد. این، در مورد بسیاری

از آمارگیریهای نمونه‌ای که به وسیلهٔ ادارهٔ سرشماری امریکا^۱، مرکز ملی آمارهای بهداشتی^۲، دفتر آمار کار^۳، و سایر سازمانها اجرا می‌شود مصداق دارد.

برای برآورد کردن واریانسهای برآوردهای حاصل از آمارگیریهای نمونه‌ای که مستلزم نمونه‌گیریها و شیوه‌های پیچیده برای برآورد کردن هستند ضروری است یکی از دو ردهٔ کلی از روشهایی که مخصوصاً برای این منظور بسط داده شده‌اند، یعنی روشهای خطی سازی و تکرار، به کار برده شود.

فن خطی‌سازی به وسیلهٔ کی‌فیتس [۱]، وودراف [۲] و دیگران برای آمارگیریهای نمونه‌ای و براساس تقریب سری تیلور^۴ بسط داده شده است. در اواخر دههٔ ۱۹۷۰، پژوهشگران مؤسسهٔ مثلث تحقیق^۵ شروع به ابداع مطلبی کردند که سرانجام تبدیل به یک نرم‌افزار راحت براساس خطی‌سازی شد [۳]. این نرم‌افزار، شکل نخستین SUDAAN بود که پس از تکامل بسیار به نرم‌افزاری بدل شد که کاربرد وسیعی دارد و در سراسر این کتاب مورد اشاره و بحث بوده است. قابل استفاده بودن و جاذبهٔ این نرم‌افزار و سایر نرم‌افزارها باعث شده است که خطی‌سازی شاید اکنون پراستفاده‌ترین روش برآورد واریانس برای آمارگیریهای پیچیده باشد.

تکرار نیز یک ردهٔ کلی از روشهایی است که با آن برآوردی از واریانس یک پارامتر جامعه‌ای برآورد شده به دست می‌آید، به این ترتیب که پارامتر جامعه‌ای برآورد شده به صورت مجموع یا میانگینی از چندین آماره بیان می‌شود که هر یک بر زیرمجموعه‌ای از مشاهدات نمونه‌ای متکی است. واریانس پارامتر جامعه‌ای برآورد شده را سپس می‌توان با به دست آوردن برآوردی از واریانس این آماره‌های «جزء نمونه‌ای» برآورد کرد. مثال ساده‌ای از این فن در فصل ۴، با توجه به برآورد واریانس یک مشخصهٔ جامعه‌ای برآورد شده، تحت نمونه‌گیری سیستماتیک مکرر نشان داده شد. گونه‌هایی از این رده از فنون طی چندین سال مورد استفاده قرار گرفته است. هسن و همکاران [۵] در کتاب درسی خود که در اوایل دههٔ ۱۹۵۰ تألیف کردند به این فنون تحت عنوان روشهای گروه تصادفی اشاره کرده‌اند. این روشها در دههٔ ۱۹۶۰ توسط مک‌کارتی [۶] و [۷] دقیقتر شده و در سطح گسترده‌ای در مرکز ملی آمارهای بهداشتی و سایر سازمانهای فدرال مورد استفاده قرار گرفتند. مطالعاتی توسط فرانکل [۸] و بین [۹] انجام گرفته و نشریاتی متعدد به وسیلهٔ لمی شو و همکاران [۱۰] تا [۱۶] به رشتهٔ تحریر درآمده‌اند که سعی در ارزشیابی این روشها داشته‌اند.

¹ U.S. Bureau of the Census

² National Center for Health Statistics

³ Bureau of Labor Statistics

⁴ Taylor series approximation

⁵ Research Triangle Institute

هدف از این فصل، معرفی منطق نهفته در پس الگوریتمهای مورد استفاده در بسته‌های نرم‌افزاری گوناگونی است که در حال حاضر برای برآورد واریانسهای برآوردهای حاصل از آمارگیریهای دارای طرحهای پیچیده در دسترس‌اند. در بحثی که به دنبال خواهد آمد، توجه خود را بر روشهای خطی‌سازی و تکرار متمرکز می‌کنیم که فعلاً پراستفاده‌ترین روشهای برآورد کردن واریانس به شمار می‌روند.

۱.۱۲ خطی‌سازی

در آمارگیریهای نمونه‌ای پیچیده، برآورد کردن واریانسهای مشخصه‌های جامعه‌ای برآورد شده می‌تواند مشکل باشد چه به علت راهی که نمونه با آن راه گرفته شده است و چه به دلیل راهی که برآوردهای مشخصه‌های جامعه‌ای با آن ساخته شده‌اند. مثلاً، یک مجموع برآورد شده براساس چندین مرحله نمونه‌گیری خوشه‌ای، تعدیل نسبی، و پس طبقه‌بندی غالباً ممکن است به راحتی یک برآوردگر مناسب واریانس را به دست ندهد. در چنین وضعیتهایی می‌توان از فن خطی‌سازی برای ساختن تقریبی برای شکل تابعی مشخصه جامعه‌ای برآورد شده استفاده کرد که تابعی خطی از مشاهدات اصلی است و لذا برای ساختن برآوردگر واریانس مناسب است.

بحث روش‌شناسی خطی‌سازی که در پی خواهد آمد اساساً گزیده و فشرده مطالبی است که در پیوست فنی کتابچه راهنمای SUDAAN ارائه شده است که همراه با نسخه رایانه شخصی این محصول نرم‌افزاری است [۳].

برای سهولت کار، فرض می‌کنیم جامعه‌ای داریم که از مشاهدات x_i و y_i درباره i امین واحد شمارش تشکیل شده است و باز فرض می‌کنیم می‌خواهیم، θ ، نوعی مشخصه جامعه‌ای را برآورد کنیم و یک برآوردگر، $\hat{\theta}$ ، هم وجود دارد که تابع $f(x, y)$ از متغیرهای x و y است که هر دو تابعی خطی از مشاهدات نمونه‌ای هستند. ولی تابع $f(x, y)$ تابع خطی مشاهدات نمونه‌ای نیست.

برای مثال، $\hat{\theta}$ ممکن است $\frac{\bar{x}}{\bar{y}}$ ، نسبت میانگینهای نمونه‌ای، \bar{x} و \bar{y} ، باشد. واضح است که \bar{x} و \bar{y} تابعهای خطی مشاهدات نمونه‌ای هستند در حالی که $f(\bar{x}, \bar{y})$ تابع غیرخطی دو میانگین نمونه‌ای است و از این رو است که در مشاهدات نمونه‌ای نیز غیرخطی است.

محاسبه واریانس $\hat{\theta}$ با روش خطی‌سازی، شیوه‌ای دومرحله‌ای است. مرحله اول یا مرحله خطی‌سازی مستلزم تقریب کردن $f(x, y)$ به وسیله سری تیلور مرتبه اول است. این مرحله به تقریبی می‌انجامد که نسبت به آماره‌های نمونه‌ای \bar{x} و \bar{y} ، خطی است و به همین علت در مشاهدات نمونه‌ای نیز خطی خواهد بود. همین که این تقریب سری تیلور مرتبه اول برای $\hat{\theta}$ به دست آمد، مرحله دوم

متضمن استفاده از روشهای مبتنی بر طرح، از قبیل روشهای توصیف شده در فصلهای ۳ تا ۱۱، برای برآورد واریانس آن خواهد بود. در زیر توضیح می‌دهیم که چگونه این کار به صورت مکانیکی انجام می‌پذیرد.

z_i را که مقدار خطی شده $f(\bar{x}, \bar{y})$ برای مشاهده i نامیده می‌شود به صورت زیر تعریف می‌کنیم:

$$z_i = (\partial f_{\bar{x}})x_i + (\partial f_{\bar{y}})y_i$$

که در آن

$$\partial f_{\bar{x}} = \frac{\partial f(\bar{x}, \bar{y})}{\partial \bar{x}}, \quad (\bar{X}, \bar{Y}) \text{ در نقطه}$$

و

$$\partial f_{\bar{y}} = \frac{\partial f(\bar{x}, \bar{y})}{\partial \bar{y}}, \quad (\bar{X}, \bar{Y}) \text{ در نقطه}$$

(که در آنها \bar{X} و \bar{Y} پارامترهای جامعه‌ای نامعلوم‌اند که \bar{x} و \bar{y} برآوردگرهای نارایب آنها هستند).

مثلاً اگر $\hat{\theta} = \bar{x}/\bar{y}$ باشد، آن‌گاه

$$\partial f_{\bar{x}} = \frac{1}{\bar{y}}, \quad \partial f_{\bar{y}} = -\frac{\bar{x}}{\bar{y}^2}, \quad z_i = \frac{1}{\bar{y}}x_i - \frac{\bar{x}}{\bar{y}^2}y_i$$

بسط سری تیلور مرتبه اول برای $\hat{\theta}$ به صورت زیر است

$$\begin{aligned} \hat{\theta} &= f(\bar{x}, \bar{y}) \approx (\partial f_{\bar{x}})(\bar{x} - \bar{X}) + (\partial f_{\bar{y}})(\bar{y} - \bar{Y}) \\ &= (\partial f_{\bar{x}})\bar{x} + (\partial f_{\bar{y}})\bar{y} - ((\partial f_{\bar{x}})\bar{X} + (\partial f_{\bar{y}})\bar{Y}) \end{aligned}$$

و

$$Var(\hat{\theta}) \cong Var((\partial f_{\bar{x}})\bar{x} + (\partial f_{\bar{y}})\bar{y})$$

زیرا عبارت $((\partial f_{\bar{x}})\bar{X} + (\partial f_{\bar{y}})\bar{Y})$ یک مقدار ثابت است. ولی

$$(\partial f_{\bar{x}})\bar{x} + (\partial f_{\bar{y}})\bar{y} = (\partial f_{\bar{x}})\sum_{i=1}^n \frac{x_i}{n} + (\partial f_{\bar{y}})\sum_{i=1}^n \frac{y_i}{n} = \sum_{i=1}^n \frac{z_i}{n} = \bar{z}$$

به این ترتیب، مشکل یافتن واریانس تقریبی $\hat{\theta}$ به مسئله تعیین واریانس یک ترکیب خطی از مشاهدات نمونه‌ای تبدیل می‌شود و می‌تواند با تعیین طرح نمونه‌ای و استفاده از فرمول مناسب به دست آید. مثلاً اگر طرح، نمونه‌گیری تصادفی ساده باشد، در آن صورت، واریانس $\hat{\theta}$ به صورت زیر برآورد می‌شود

$$\widehat{Var}(\hat{\theta}) = \frac{\sum_{i=1}^n (\tilde{z}_i - \bar{\tilde{z}})^2}{n(n-1)} \left(\frac{N-n}{N} \right) \quad (1.12)$$

که در آن

$$\tilde{z}_i = (\partial f_{\bar{x}})x_i + (\partial f_{\bar{y}})y_i$$

$$\partial f_{\bar{x}} = \frac{\partial f(\bar{x}, \bar{y})}{\partial \bar{x}}, \quad (\bar{x}, \bar{y}) \text{ محاسبه شده در نقطه}$$

و

$$\partial f_{\bar{y}} = \frac{\partial f(\bar{x}, \bar{y})}{\partial \bar{y}}, \quad (\bar{x}, \bar{y}) \text{ محاسبه شده در نقطه}$$

(زیرا \bar{X} و \bar{Y} پارامترهای مجهول اند).

مثال تشریحی: یک مثال بسیار کوچک ساده را انتخاب می‌کنیم تا نشان دهیم خطی سازی دقیقاً چگونه عمل می‌کند. فرض کنید هزینه‌های بیمه کمکهای پزشکی را برای یک بیمار حسابرسی می‌کنیم و یک نمونه تصادفی ساده متشکل از ۱۰ درخواست از مجموع ۶۵ درخواست حق بیمه را که از جانب بیمار تنظیم شده است انتخاب کرده‌ایم. داده‌های حاصل از این حسابرسی در جدول ۱.۱۲ نشان داده شده‌اند. می‌خواهیم نسبت مجموع دلارهای اضافی را که از طرف این بیمار پرداخت شده است برآورد کنیم.

جدول ۱.۱۲ پرداختها و اضافه پرداختهای بیمه کمکهای پزشکی مربوط به

۱۰ درخواست حق بیمه که برای بیمار صادر شده است

مشاهده خطی شده	اضافه پرداخت (x)	پرداخت (y)	درخواست
\tilde{z}_i	به دلار	به دلار	حق بیمه
۰/۴۶۷۵	۲۱۰	۲۱۰	۱
-۰/۱۰۹۴	۰	۷۸	۲
-۰/۰۳۴۶	۱۲۳	۳۴۳	۳
۰/۱۵۱۸	۱۵۷	۲۹۸	۴
-۰/۴۸۹۴	۰	۳۴۹	۵
۰/۴۶۷۵	۲۱۰	۲۱۰	۶
-۰/۷۵۱۶	۰	۵۳۶	۷
۰/۰۸۴۶	۱۳۵	۲۸۹	۸
-۰/۱۳۷۴	۰	۹۸	۹
۰/۳۵۰۸	۲۳۰	۳۴۵	۱۰

از روی این داده‌ها داریم $\bar{x} = 106/5$ و $\bar{y} = 275/6$. چون از برآورد نسبتی استفاده می‌کنیم، متغیر خطی شده

$$\tilde{z}_i = \frac{1}{275/6} x_i - \frac{106/5}{(275/6)^2} y_i$$

به صورتی که در بالا توصیف شد به دست خواهد آمد. مثلاً

$$\tilde{z}_i = \frac{1}{275/6} \times 210 - \frac{106/5}{(275/6)^2} \times 210 = 0.4675$$

\tilde{z}_i ها در ستون چهارم جدول ۱.۱۲ نشان داده شده‌اند. با $N = 65$ و $n = 10$ ، خطای معیار برآورد نسبتی، r ، را از برابری (۱.۱۲) محاسبه می‌کنیم که برابر است با 0.1158 .

□

از SUDAAN و STATA، هر دو، برای برآورد کردن واریانسهای آماره‌های برآورد شده از خطی سازی استفاده می‌کنند. برای استفاده از هر یک از این نرم‌افزارها لازم است متغیرهای $N (= 65)$ و $w (= \frac{65}{10} = 6.5)$ را به هر یک از سابقه‌ها در مجموعه داده‌های جدول ۱.۱۲ اضافه کنیم. به متغیر \tilde{z}_i نیاز نخواهیم داشت، زیرا توسط برنامه محاسبه خواهد شد. مجموعه داده‌های حاصل برای استفاده به وسیله SUDAAN یا STATA در زیر نشان داده می‌شود:

claim	payment	ovpymnt	N	w
1	210	210	65	6.5
2	78	0	65	6.5
3	343	123	65	6.5
4	298	157	65	6.5
5	349	0	65	6.5
6	210	210	65	6.5
7	536	0	65	6.5
8	289	135	65	6.5
9	98	0	65	6.5
10	345	230	65	6.5

فرمانهای زیر برای برآورد کردن نسبت x/y و خطای معیار آن توسط STATA به کار می‌روند:

```
use"a:\exmp12_2.dta", clear
. svyset fpc N
. svyset pweight w
. svyratio ovpaymnt/payment
```

این فرمانها، خروجی زیر را تولید می‌کنند:

Survey ratio estimation					
pweight:	W	Number of obs	=	10	
Strata:	< one >	Number of strata	=	1	
PSU:	< observations >	Number of PSUs	=	10	
FPC:	N	Population size	=	65	
	Ratio	Estimate	Std. Err.	[95% Conf. Interval]	Deff
	ovpaymnt/payment	.3864296	.1158187	.1244294 .6484298	1
<p>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.</p>					

فرمانهای زیر برای برآورد کردن با SUDAAN به کار می‌روند:

```

1 PROC RATIO DATA = EXMP12_2 FILETYPE = SAS DESIGN = WOR;
2 NEST_ONE_;
3 WEIGHT W;
4 TOTCNT N;
5 NUMER OVPAYMNT;
6 DENOM PAYMENT;
7 SETENV COLWIDTH = 15;
8 SETENV DECWIDTH = 3;
    
```

این فرمانها، خروجی زیر را تولید می‌کنند:

Number of observations read	:	10	Weighted count :	65
Number of observations skipped	:	0		
(WEIGHT variable nonpositive)				
Denominator degrees of freedom	:	9		
by: Variable, One.				
Variable			One	
			1	
OVPAYMNT/PAYMENT	Sample Size			10.000
	Weighted Size			65.000
	Weighted X-sum			17914.000
	Weighted Y-sum			6922.500
	Ratio Est.			0.386
	SE Ratio			0.116

۲.۱۲ روشهای تکرار

روشهای تکرار (که گاهی روشهای باز نمونه‌گیری خوانده می‌شوند) واریانس نمونه‌گیری یک آماره را با محاسبه آن آماره برای زیرمجموعه‌ای از نمونه و بررسی تغییرپذیری آن روی مجموعه‌ها برآورد می‌کنند. دو رهیافت کلی برای تکرار که طی سه دهه گذشته بسط داده شده‌اند، رهیافت جک‌نایف^۱ و رهیافت تکرار مکرر متعادل^۲ (که به نام رهیافت نیمه نمونه‌ای متعادل^۳ نیز خوانده می‌شود) هستند. در این بحث از فنون تکرار، تمرکز ما منحصرأ بر روشهای تکرار مکرر متعادل (BRR) است. زیرا از لحاظ تاریخی بیشتر از روشهای جک‌نایف در کارهای آمارگیری نمونه‌ای مورد استفاده بوده‌اند.

۱.۲.۱۲ روش تکرار مکرر متعادل

روش تکرار مکرر متعادل بیشترین کاربرد را در رده‌ای از طرحهای نمونه‌ای دارد که در سطح گسترده‌ای مورد استفاده قرار می‌گیرند و در آنها واحدهای نمونه‌گیری اولیه (PSU ها) در L طبقه گروه‌بندی می‌شوند، از هر طبقه، نمونه‌هایی متشکل از دو واحد نمونه‌گیری اولیه انتخاب می‌شود، و از هر واحد نمونه‌گیری اولیه، برآوردهایی مستقل از مشخصه‌های طبقه به دست می‌آید.

حال، روش تکرار مکرر متعادل را با نشان دادن چگونگی برآورد کردن واریانس یک برآورد نسبتی با استفاده از این شیوه به صورتی دقیقتر نشان می‌دهیم. فرض کنید $x'_{h1}, x'_{h2}, y'_{h1}, y'_{h2}$ برآوردهایی از هر یک از دو واحد نمونه‌گیری اولیه برای پارامترهای طبقه‌ای X_h و Y_h در هر یک از L طبقه هستند. (برای این بحث فرض می‌کنیم که طبقه‌ها برای تعیین برآوردهای جامعه‌ای دارای وزن برابرند). پس، نتیجه می‌گیریم که برآورد نسبتی r برای نسبت جامعه‌ای R بر مبنای همه $2L$ برآورد از فرمول زیر به دست می‌آید

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

که در آن

$$x' = \sum_{h=1}^L \sum_{i=1}^2 (x'_{hi} / 2)$$

و

$$y' = \sum_{h=1}^L \sum_{i=1}^2 (y'_{hi} / 2)$$

¹ Jackknifing Approach

² Balanced Repeated Replication Approach (BRR)

³ Balanced Half-sample Approach

برآورد نیمه نمونه‌ای نسبت R را به صورت زیر تعریف می‌کنیم:

$$r_{(k)} = \frac{\sum_{h=1}^L [\delta_{kh} x'_{h1} + (1 - \delta_{kh}) x'_{h2}]}{\sum_{h=1}^L [\delta_{kh} y'_{h1} + (1 - \delta_{kh}) y'_{h2}]}$$

که در آن $(\delta_{k1}, \dots, \delta_{kL})$ یک بردار L -بعدی است که عناصر آن برابر با یک یا صفرند. به عبارت دیگر، هر $r_{(k)}$ برآوردی از R است که با گرفتن یکی از دو برآورد از هر طبقه تشکیل شده است. واضح است که 2^L بردار ممکن (δ_k) به صورتی که در بالا نشان داده شده است وجود دارند و از این رو، در کل 2^L برآورد نیمه نمونه‌ای ممکن $r_{(k)}$ موجودند. برآوردی مناسب از واریانس r که از روی 2^L نیمه نمونه ممکن محاسبه شده باشد چنین خواهد بود

$$\hat{Var}(r) = \left(\frac{1}{2^L} \right) \sum_{k=1}^{2^L} (r_{(k)} - \bar{r})^2$$

معمولاً 2^L به قدری بزرگ است که محاسبه همه $r_{(k)}$ های ممکن عملی نیست. زیرمجموعه‌ای از K برآورد نیمه نمونه‌ای را (که در آن $K \leq 2^L$) می‌توان با به دست آوردن δ_{ij} از روی ماتریس متعامد از نوعی که پلاکت و برمن توصیف کرده‌اند [۱۷] ساخت. جزئیات این شیوه توسط مک‌کارتی [۶] شرح و بسط داده شده است که نشان می‌دهد این امکان وجود دارد که زیرمجموعه کوچکی از نیمه نمونه‌ها با روش بسیار دقیقی انتخاب شود و با وجود این برآوردهای رضایتبخشی از واریانس به دست آید. در واقع، او نشان داده است که استفاده از زیرمجموعه‌ای کوچک برای برآوردهای خطی، از قبیل میانگینها و مجموعها، دقیقاً همان برآوردی از واریانس را نتیجه می‌دهد که از همه 2^L نیمه نمونه به دست خواهد آمد. مجموعه برآوردهای نیمه نمونه‌ای حاصل را متعامد می‌نامند زیرا نشان داده شده است که سهم آنها در واریانس بین طبقه‌ای خنثی می‌شود. فرض کنید مجموعه متعادلی شامل K برآورد نیمه نمونه‌ای از این نوع را در نظر می‌گیریم و برآورد نهایی \bar{r} خود را به صورت زیر تعریف می‌کنیم:

$$\bar{r} = \frac{1}{K} \sum_{k=1}^K r_{(k)} \quad (۲.۱۲)$$

برآوردگر واریانس این برآورد نسبتی \bar{r} از فرمول زیر به دست می‌آید

$$\hat{Var}(\bar{r}) = \left(\frac{1}{K} \right) \sum_{k=1}^K (r_{(k)} - \bar{r})^2 \quad (۳.۱۲)$$

منطق و خواص آماری این برآوردگر واریانسی نیمه نمونه‌ای متعادل توسط مک‌کارتی [۶] و بین [۹] و نیز لمی‌شو و له‌وی [۱۰] به تفصیل شرح داده شده‌اند. ما نیز برای این برآوردگر واریانسی، منطقی رهگشا در پیوست فنی این فصل ارائه می‌کنیم و حالا استفاده از این روش را با مثال ساده‌ای نمایش می‌دهیم.

مثال تشریحی: شهری به سه ناحیه خدمات فوریت‌های پزشکی^۱ تقسیم شده که هر ناحیه شامل پنج ایستگاه آمبولانس است. می‌خواهیم نسبت کسانی را که در کل شهر دچار ایست قلبی شده و پس از دیده شدن توسط پیراپزشکان یکی از ۱۵ ایستگاه آمبولانس توصیف شده در بالا زنده وارد بیمارستان شده‌اند برآورد کنیم. یک آمارگیری نمونه‌ای اجرا شده است که طی آن دو ایستگاه آمبولانس به طور تصادفی از هر یک از این سه ناحیه خدمات اورژانسی نمونه‌گیری شده‌اند. سوابع موجود در هر یک از این شش ایستگاه آمبولانس نمونه‌گیری شده مورد بررسی قرار گرفته و برآوردهایی از تعداد ایستهای قلبی و تعداد و نسبت بیمارانی که دچار ایست قلبی شده و زنده به بیمارستان رسیده‌اند برای ناحیه خدمات اورژانسی تهیه شده است. نتایج این آمارگیری نمونه‌ای در زیر آمده‌اند:

ناحیه خدمات فوریت‌های پزشکی (ESA)	ایستگاه آمبولانس	برآورد تعداد ایستهای قلبی برای ESA	برآورد تعداد بیماران دچار ایست قلبی که زنده به بیمارستان رسیده‌اند برای ESA
		y'_{hi}	x'_{hi}
۱	۱	۱۲۰	۲۵
	۲	۷۸	۲۴
۲	۱	۱۸۵	۳۰
	۲	۲۲۸	۴۹
۳	۱	۶۷۰	۸۰
	۲	۵۳۰	۷۰

باید تأکید شود که x'_{hi} و y'_{hi} که در بالا نشان داده شده‌اند برآوردهایی از تعداد ایستهای قلبی و تعداد بیمارانی که زنده به بیمارستان رسیده‌اند براساس داده‌های حاصل از ایستگاه ویژه آمبولانس هستند.

^۱ Emergency Service Area (ESA)

واضح است که طرح این آمارگیری نمونه‌ای از طبقه‌بندی بر حسب ناحیه خدمات فوریت‌های پزشکی و برآورد کردن برای هر ناحیه خدمات اورژانسی از روی دو ایستگاه آمبولانس که به طور تصادفی انتخاب شده‌اند استفاده می‌کند.

تکرارهای $۲^۳ = ۸$ نیمه نمونه زیر را می‌توان برای سه طبقه تهیه کرد:

$۱ - \delta_{k3}$	δ_{k3}	$۱ - \delta_{k2}$	δ_{k2}	$۱ - \delta_{k1}$	δ_{k1}	تکرار (k)
۰	۱	۰	۱	۰	۱	۱
۱	۰	۰	۱	۰	۱	۲
۰	۱	۱	۰	۰	۱	۳
۱	۰	۱	۰	۰	۱	۴
۰	۱	۰	۱	۱	۰	۵
۱	۰	۰	۱	۱	۰	۶
۰	۱	۱	۰	۱	۰	۷
۱	۰	۱	۰	۱	۰	۸

مثلاً تکرار ۱، مقدار $r_{(k)}$ را به صورت زیر ارائه می‌دهد:

$$r_{(1)} = \frac{\sum_{h=1}^3 [\delta_{ih} x'_{h1} + (1 - \delta_{ih}) x'_{h2}]}{\sum_{h=1}^3 [\delta_{ih} y'_{h1} + (1 - \delta_{ih}) y'_{h2}]} = \frac{۱۳۵}{۹۷۵} = ۰/۱۳۸۵$$

هشت تکرار، مقادیر زیر از $r_{(k)}$ را به دست می‌دهند:

$r_{(k)}$	$y'_{(k)}$	$x'_{(k)}$	k
۰/۱۳۸۵	۹۷۵	۱۳۵	۱
۰/۱۴۹۷	۸۳۵	۱۲۵	۲
۰/۱۵۱۳	۱۰۱۸	۱۵۴	۳
۰/۱۶۴۰	۸۷۸	۱۴۴	۴
۰/۱۴۳۶	۹۳۳	۱۳۴	۵
۰/۱۵۶۴	۷۹۳	۱۲۴	۶
۰/۱۵۶۸	۹۷۶	۱۵۳	۷
۰/۱۷۱۱	۸۳۶	۱۴۳	۸

از داده‌هایی که در بالا نشان داده شده‌اند پی می‌بریم که برآورد نسبت \bar{r} و واریانس برآورد شده آن، $\hat{Var}(\bar{r})$ ، براساس هشت تکرار به صورت زیر به دست می‌آید

$$\bar{r} = \frac{1}{8} \sum_{k=1}^8 r_{(k)} = 0.1539$$

و

$$\hat{Var}(\bar{r}) = 0.000098$$

پلاکت و برمن [۱۷]، ماتریسهای $m \times m$ را با ابعاد $m = \{4, 8, 12, \dots, 100\}$ به استثنای $m = 92$ ارائه داده‌اند. بعد ماتریس مورد استفاده در یک مسئله خاص به تعداد طبقه‌ها بستگی دارد. مک‌کارتی [۶] از ماتریسهایی استفاده کرده است که در آن m مضربی از ۴ و $L \leq m \leq L + 3$ ، که در آن در همه موارد $m \ll 2^L$ ، $L > 2$ است. همین که ماتریس $m \times m$ انتخاب شد همه سطرها بجز تنها اولین ستونهای L مورد استفاده قرار می‌گیرند.

یک مجموعه متعادل از چهار تکرار را می‌توان از مجموعه هشت تکرار ممکن با اختیار کردن تکرارهای ۱، ۴، ۶ و ۷ به شرح زیر به دست آورد:

$1 - \delta_{k2}$	δ_{k3}	$1 - \delta_{k2}$	δ_{k2}	$1 - \delta_{k1}$	δ_{k1}	تکرار k
۰	۱	۰	۱	۰	۱	۱
۱	۰	۱	۰	۰	۱	۴
۱	۰	۰	۱	۱	۰	۶
۰	۱	۱	۰	۱	۰	۷

توجه کنید که در این مجموعه، هر برآورد طبقه‌ای x'_{hk} یا y'_{hk} به دفعات برابر ظاهر می‌شود (دو بار در این مثال کوچک)، و با هر برآورد دیگر به تعداد دفعات یکسان (در این مثال یک بار) پدیدار می‌شود. حال برآورد \bar{r} و واریانس آن $\hat{Var}(\bar{r})$ را برای این مجموعه چهار تکرار بررسی می‌کنیم:

$$\bar{r} = \frac{(r_{(1)} + r_{(4)} + r_{(6)} + r_{(7)})}{4} = 0.1539$$

و

$$\hat{Var}(\bar{r}) = 0.000088$$

□

چون هر برآورد طبقه‌ای نیمه نمونه‌ای با تعداد دفعاتی یکسان در مجموعه K نیمه نمونه متعادل ظاهر می‌شود، برآورد \bar{y} با برآورد حاصل از مجموعه کامل 2^L نیمه نمونه برابر خواهد بود. همچنین، مک‌کارتی [۶] نشان داده است که پایداری برآوردگر واریانسی حاصل از مجموعه‌ای از تکرارهای نیمه نمونه‌ای متعادل تقریباً همان پایداری برآوردگر واریانسی است که از مجموعه کامل 2^L تکرار نیمه نمونه‌ای ساخته می‌شود.

باید توجه داشت که غالباً آمارگیریهایی طراحی می‌شوند که در آنها $2L$ طبقه وجود دارند و از هر طبقه یک واحد نمونه‌گیری اولیه (PSU) انتخاب می‌شود. برای به کار بردن فن تکرار مکرر متعادل (BRR) واحدهای نمونه‌گیری اولیه، جفت جفت طبقه‌بندی می‌شوند تا L «شبه طبقه» تشکیل شود و فرایند برآورد کردن به همان صورتی که قبلاً شرح داده شد اجرا می‌شود. معمولاً در جفت کردن واحدها دقت می‌شود تا واحدهای نمونه‌گیری اولیه از نظر اطلاعات جمعیت‌شناختی معلوم، حتی‌الامکان شبیه باشند.

استفاده از SUDAAN برای به دست آوردن برآوردهای تکرار مکرر متعادل (BRR): آخرین نسخه SUDAAN در زمان نگارش متن حاضر (Release 7.5) می‌تواند برآوردهایی از خطاهای معیار را با استفاده از روش تکرار مکرر متعادل (BRR) و نیز برآوردهایی از راه خطی‌سازی به دست آورد. در مثال بالا که در آن از مجموعه‌ای متعادل با چهار تکرار از هشت تکرار ممکن برای برآورد کردن نسبت بیماران دچار ایست قلبی که سالم به بیمارستان رسیده بودند استفاده شده بود، مجموعه داده‌های مورد نیاز برای تحلیل توسط SUDAAN، در زیر نشان داده شده است:

ESA	AMBSTAT	CARDARRS	ALIVE	WT	REPWT1	REPWT4	REPWT6	REPWT7
1	1	120	25	2.5	5	5	0	0
1	2	78	24	2.5	0	0	5	5
2	1	185	30	2.5	5	0	5	0
2	2	228	49	2.5	0	5	0	5
3	1	670	80	2.5	5	0	0	5
3	2	530	70	2.5	0	5	5	0

توجه کنید که در این پرونده شش سابقه، یعنی یک سابقه برای هر نقطه نمونه، و ۹ متغیر وجود دارد. دو متغیر اول، یعنی ESA و $AMBSTAT$ نشانه ناحیه خدمات اورژانسی و ایستگاه آمبولانس است. دو متغیر بعدی، $CARDARRS$ و $ALIVE$ به ترتیب نشانه تعداد ایستهای قلبی که در هر ایستگاه آمبولانس نمونه رسیدگی شده و تعداد این قبیل بیماران است که زنده به بیمارستان رسیده‌اند. متغیر WT وزن نمونه، N/n ، است که به هر ایستگاه آمبولانس اختصاص یافته است. برای هر سابقه نمونه در این مثال، $N=15$ ، $n=6$ و $WT=2/5$ است. متغیر $REPWT1$ نشانه وزنی است که در نیمه نمونه

متعادل ۱ به هر مشاهده نمونه‌ای مورد استفاده در نیمه نمونه داده شده است. چون نیمه نمونه متعادل ۱ فقط شامل سه نقطه نمونه‌ای است، برای سه نقطه نمونه‌ای که در به دست آوردن آن برآورد نیمه نمونه‌ای به کار رفته است $REPWT1 = \frac{15}{3} = 5$ و برای سه نقطه نمونه‌ای که در نیمه نمونه قرار نداشته‌اند برابر با صفر خواهد بود. سه متغیر دیگر، یعنی $REPWT4$ ، $REPWT6$ و $REPWT7$ نیز به همین قیاس تعریف می‌شوند.

فرمانهای زیر برای به دست آوردن برآوردهای تکرار مکرر متعادل (BRR) با استفاده از SUDAAN به کار می‌روند:

```
PROC RATIO DATA = AMBLNCE2 FILETYPE = SAS DESIGN = BRR;
WEIGHT WT;
REPWGT REPWT1 REPWT4 REPWT6 REPWT7;
NUMER ALIVE;
DENOM CARDARRS;
SETENV COLWIDTH = 14;
SETENV DECWIDTH = 5;
```

فرمان اول، نام و نوع پرونده داده‌ها را معرفی می‌کند و حاکی از آن است که برآورد نسبتی قرار است به روش تکرار مکرر متعادل اجرا شود. فرمان دوم، متغیر شامل وزنهای نمونه‌گیری را نشان می‌دهد و فرمان سوم، متغیرهای شامل وزنهای نمونه‌گیری را برای آن نیمه نمونه بخصوص در هر نیمه نمونه متعادل نشان می‌دهد. باقیمانده فرمانها بیانگر متغیرهایی هستند که شامل صورت و مخرج کسر نسبت و قالب خروجی است. خروجی حاصل از این فرمانها در زیر نشان داده شده است:

Variance Estimation Method: BRR by: Variable, One.		
Variable		One --
ALIVE/CARDARRS	Sample Size	6.00000
	Weighted Size	15.00000
	Weighted X-sum	4527.50000
	Weighted Y-sum	695.00000
	Ratio Est.	0.15351
	SE Ratio	0.00943

برآورد نسبت و خطای معیار آن که از SUDAAN به دست می‌آیند با آنچه قبلاً نشان داده شد تطبیق می‌کنند (جز در مورد تفاوتی در سومین رقم دهدهی که ناشی از خطای گرد کردن است).

حالا به مقایسه برآوردهای بالا که با استفاده از SUDAAN به دست آمده است با برآوردهای حاصل از همین داده‌ها با استفاده از روش خطی‌سازی می‌پردازیم. طرح، یک نمونه تصادفی طبقه‌بندی شده است که نواحی خدمات اورژانسی، طبقه‌ها را تشکیل می‌دهند و نمونه‌ای متشکل از ۲ ایستگاه آمبولانس از جامعه متشکل از ۵ ایستگاه آمبولانس در هر ناحیه خدمات اورژانسی گرفته شده است. پرونده داده‌ها برای به دست آوردن برآوردهای خطی‌سازی در زیر نشان داده شده است.

ESA	AMBSTAT	CARDARRS	ALIVE	WT	NAMBSTAT
1	1	120	25	2.5	5
1	2	78	24	2.5	5
2	1	185	30	2.5	5
2	2	228	49	2.5	5
3	1	670	80	2.5	5
3	2	530	70	2.5	5

توجه کنید که متغیر *NAMBSTAT* به هر سابقه اضافه شده است. این متغیر از آن رو لازم است که نمونه‌گیری در داخل نواحی خدمات اورژانسی بدون جایگذاری است و نشانه آن است که اندازه جامعه در هر طبقه برابر با پنج است. فرمانهای مورد استفاده برای تهیه برآوردها ذیلاً نشان داده می‌شوند:

```
PROC RATIO DATA = AMBLNCE2 FILETYPE = SAS DESIGN = STRWOR;
  NEST ESA;
  TOTCNT NAMBSTAT;
  WEIGHT WT;
  NUMER ALIVE;
  DENOM CARDARRS;
  SETENV COLWIDTH = 14;
  SETENV DECWIDTH = 5;
```

حکم مربوط به طرح *design = strwor* نشان می‌دهد که طرح، یک نمونه تصادفی طبقه‌بندی شده یک‌مرحله‌ای است که واحدهای شمارش بدون جایگذاری انتخاب شده‌اند. خروجی حاصل از این مجموعه فرمانها در پایین نشان داده شده است:

Variance Estimation Method: Taylor Series (STRWOR)		
by: Variable, One.		
Variable		One
		--
ALIVE/CARDARRS	Sample Size	6.00000
	Weighted Size	15.00000
	Weighted X-sum	4527.50000
	Weighted Y-sum	695.00000
	Ratio Est.	0.15351
	SE Ratio	0.00760

توجه کنید که خطای معیار نسبت برآورد شده که از خطی‌سازی به دست می‌آید کوچکتر از خطای معیاری است که از روش تکرار مکرر متعادل به دست می‌آید و این دور از انتظار نیست، زیرا روش خطی‌سازی در نظر می‌گیرد که نمونه‌گیری بدون جایگذاری است (و بنابراین از تصحیح جامعه متناهی استفاده می‌کند)، در حالی که در روش تکرار مکرر متعادل چنین چیزی لحاظ نمی‌شود. اگر از روش خطی‌سازی استفاده کرده ولی طرح نمونه‌گیری طبقه‌بندی شده با جایگذاری را به کار برده بودیم (حکم طرح $design = strwr$ می‌شد) برآورد خطای معیار نسبت برآورد شده برابر با $0/00981$ به دست می‌آمد که به خطای معیار برابر با $0/00943$ که از روش تکرار مکرر متعادل به دست می‌آمد بسیار نزدیک است.

۲.۲.۱۲ برآورد کردن به روش جک‌نایف

جک‌نایف روش دیگری است که برای برآورد کردن خطاهای معیار برآوردهای حاصل از آمارگیریهای نمونه‌ای پیچیده به کار می‌رود و مانند روش تکرار مکرر متعادل، مستلزم محاسبه برآورد موردنظر (مانند مجموع، میانگین، نسبت) برای چند نیمه نمونه و سپس محاسبه واریانس این برآوردها روی مجموعه نیمه نمونه‌هاست. برآورد کردن به روش جک‌نایف تاریخچه‌ای طولانی داشته است که به نگارش مقاله‌ای توسط کنویی در دهه ۱۹۵۰ برمی‌گردد [۲۰] و صورتهای بسیار متعددی از آن در دست است. در اینجا روش‌شناسی ویژه جک‌نایف را برای برآورد نسبتی که در حال حاضر در SUDAAN مورد استفاده قرار می‌گیرد شرح خواهیم داد و سپس روش‌شناسی برآورد کردن با SUDAAN را مجدداً با استفاده از نمونه ایستگاههای آمبولانس نشان خواهیم داد.

فرض می‌کنیم که L طبقه و n_h نمونه از واحدهای نمونه‌گیری اولیه در داخل طبقه h در اختیار داریم ($h = 1, \dots, L$). می‌خواهیم خطای معیار برآورد نسبتی، r ، را برآورد کنیم که در آن

$$r = \frac{\sum_{h=1}^L \sum_{i=1}^{n_h} x'_{hi}}{\sum_{h=1}^L \sum_{i=1}^{n_h} y'_{hi}}$$

و x'_{hi} و y'_{hi} برآوردهایی بر مبنای داده‌های حاصل از آن واحد نمونه‌گیری اولیه خاص‌اند.

برای هر واحد نمونه‌گیری اولیه در مجموعه داده‌ها، برآورد $r_{(hi)}$ را برای نسبت جامعه‌ای مبتنی بر همه مشاهدات به استثنای آنهایی که در واحد نمونه‌گیری اولیه hi قرار دارند به دست می‌آوریم. $r_{(hi)}$ برای یک برآورد نسبتی از فرمول زیر به دست می‌آید

$$r_{(hi)} = \frac{\sum_{h' \neq h} \sum_{i=1}^{n_h} x'_{h'i} + \sum_{i' \neq i} \frac{n_h}{n_h - 1} x'_{hi'}}{\sum_{h' \neq h} \sum_{i=1}^{n_h} y'_{h'i} + \sum_{i' \neq i} \frac{n_h}{n_h - 1} y'_{hi'}}$$

توجه کنید که برآورد، $r_{(hi)}$ ، از برآوردهای طبقه‌ای حاصل از همه واحدهای نمونه‌گیری اولیه بجز واحد نمونه‌گیری اولیه hi تشکیل شده است و برآورد طبقه‌ای حاصل از آن واحد نمونه‌گیری اولیه بر برآوردهای حاصل از سایر واحدهای نمونه‌گیری اولیه در داخل آن طبقه متکی است که به طریقی مناسب با عامل $(n_h / (n_h - 1))$ تعدیل شده است تا نبود داده‌های مربوط به آن واحد نمونه‌گیری اولیه را بازتاب دهد.

همین که $r_{(hi)}$ ها محاسبه شدند، برآورد خطای معیار r از فرمول زیر به دست می‌آید

$$\hat{SE}(r) = \sqrt{\sum_{h=1}^L \sum_{i=1}^{n_h} \frac{n_h - 1}{n_h} (r_{(hi)} - r)^2} \quad (۴.۱۲)$$

مثال تشریحی: با در نظر گرفتن مجدد مثال تشریحی قبل، داده‌های زیر را در اختیار داریم.

ESA(h)	AMBSTAT(i)	x'_{hi}	y'_{hi}	$r_{(hi)}$	r	$(r_{(hi)} - r)^2 / 2$
1	1	692.5	4422.5	.156586	.153506	4.74×10^{-6}
1	2	697.5	4632.5	.150567	.153506	4.32×10^{-6}
2	1	742.5	4635.0	.160194	.153506	2.24×10^{-5}
2	2	647.5	4420.0	.146493	.153506	2.46×10^{-5}
3	1	670.0	4177.5	.160383	.153506	2.36×10^{-5}
3	2	720.0	4877.5	.147617	.153506	1.73×10^{-5}
Total						9.7×10^{-5}

x'_{hi} ، y'_{hi} ، و $r_{(hi)}$ همان طور که در مورد x'_{11} ، y'_{11} و $r_{(11)}$ نشان داده شد به دست می‌آیند:

$$x'_{11} = 2/5 \times (2 \times 24 + 30 + 49 + 80 + 70) = 692/5$$

$$y'_{11} = 2/5 \times (2 \times 78 + 185 + 228 + 670 + 530) = 4422/5$$

$$r_{(11)} = \frac{692/5}{4422/5} = 0.156586$$

همان‌طور که قبلاً نشان داده شد، $r = 0.153506$ ، و از برابری (۴.۱۲) داریم

$$\hat{SE}(r) = \sqrt{9/7 \times 10^{-5}} = 0.009849$$

□

این برآورد کردن را می‌توان با استفاده از SUDAAN با فرمانهای زیر به دست آورد:

```
PROC RATIO DATA = AMBLNCE2 FILETYPE = SAS DESIGN = JACKKNIFE;
  NEST ESA;
  WEIGHT WT;
  NUMER ALIVE;
  DENOM CARDARRS;
  SETENV COLWIDTH = 14;
  SETENV DECWIDTH = 7;
```

که خروجی زیر را تولید می‌کند:

Variance Estimation Method: Jackknife
by: Variable, One.

Variable		One
		--
ALIVE/CARDARRS	Sample Size	6.0000000
	Weighted Size	15.0000000
	Weighted X-sum	4527.5000000
	Weighted Y-sum	695.0000000
	Ratio Est.	0.1535064
	SE Ratio	0.0098492

توجه کنید که برآورد خطای معیار نسبت، با آنچه که از استفاده مستقیم از برابری (۴.۱۲) به دست آمده است مطابقت دارد و مقدار ۰/۰۰۹۸۴۲ آن بسیار شبیه به مقداری است که قبلاً از روش تکرار مکرر متعادل به دست آمده بود (برآورد خطای معیار r برابر است با ۰/۰۰۹۴۳). هر دوی این مقادیر به مقدار برآورد شده ۰/۰۰۹۴۳ که با استفاده از روش خطی‌سازی با فرض گرفتن طرح نمونه‌گیری با جایگذاری به دست می‌آید شبیه‌اند ولی از مقدار به دست آمده از روش خطی‌سازی تحت فرض نمونه‌گیری بدون جایگذاری (که برآورد خطای معیار r برابر است با ۰/۰۰۷۶۰) بیشترند.

فرض پایه‌ای برای استفاده از یکی از دو شیوه تکرار مکرر متعادل یا جک‌نایف برای برآورد کردن واریانس در نمونه‌گیری یک‌مرحله‌ای، آن است که یا نمونه‌گیری با جایگذاری است یا کسر نمونه‌گیری کوچک (یعنی کمتر از ۱۰ درصد) است. برای نمونه‌گیری چندمرحله‌ای باید نمونه‌گیری در مرحله اول نیز با جایگذاری باشد (یا کسر نمونه‌گیری کوچک باشد). نمونه‌گیری بدون جایگذاری در مراحل بعدی مجاز است.

۳.۲.۱۲ برآورد کردن تغییرپذیری مصاحبه‌گر با استفاده از طرحهای نمونه‌گیری تکرار شده (نمونه‌های نافذ بر هم)

از روشهای تکرار می‌توان برای برآورد آن بخشی از تغییرپذیری استفاده کرد که ناشی از تفاوت‌هایی در فرایند اندازه‌گیری در میان یکایک افرادی است که مسئولیت جمع‌آوری داده‌ها را به عهده دارند. گاهی

برای اندازه‌گیری این تغییرپذیری «مصاحبه‌گر»، از روشی موسوم به نمونه‌های نافذ بر هم استفاده می‌شود. این روش را شاید بتوان به بهترین وجه دربارهٔ وضعیتی توضیح داد که در آن نمونه تصادفی ساده‌ای متشکل از n عنصر از جامعه‌ای با N عنصر انتخاب می‌شود. اگر قرار باشد داده‌ها توسط M مصاحبه‌گر جمع‌آوری شوند در آن صورت، به جای گرفتن یک نمونه تصادفی ساده از n عنصر، M نمونه تصادفی ساده مستقل انتخاب می‌کنیم که هر یک شامل $\bar{n} = n/M$ عنصر باشد و هر مجموعه متشکل از \bar{n} عنصر طی فرایندی تصادفی به یک مصاحبه‌گر خاص اختصاص داده می‌شود به طوری که احتمال تخصیص هر یک از M نمونه به هر مصاحبه‌گر یکسان است.

تحت این طرح، هر مصاحبه‌گر برآورد جداگانه‌ای، \bar{x}_i ، از سطح میانگین \bar{X} برای مشخصه X در جامعه تهیه خواهد کرد و واریانس، σ_x^2 ، مربوط به توزیع X می‌تواند به وسیله s_{bx}^2 برآورد شود با این فرض که در فرایند اندازه‌گیری هیچ تفاوتی میان مصاحبه‌گران وجود نداشته است:

$$s_{bx}^2 = \frac{\bar{n} \sum_{i=1}^M (\bar{x}_i - \bar{x})^2}{M-1}$$

که در آن میانگین \bar{x} میانگین \bar{x}_i است.

برآورد دیگری از σ_x^2 که چه اثرهای مصاحبه‌گر وجود داشته یا نداشته باشند همچنان معتبر است از فرمول زیر به دست می‌آید

$$s_{wx}^2 = \frac{\sum_{i=1}^M s_{ix}^2}{M}$$

که در آن، برای $i = 1, \dots, M$

$$s_{ix}^2 = \frac{\sum_{j=1}^{\bar{n}} (x_{ij} - \bar{x}_i)^2}{\bar{n}-1}$$

تحت فرض صفر مبنی بر اینکه هیچ اثر مصاحبه‌گر وجود ندارد، نسبت s_{bx}^2/s_{wx}^2 با $(M-1)$ و $(n-M)$ درجه آزادی از توزیع F پیروی خواهد کرد و این آمارهٔ آزمون می‌تواند برای آزمون آن فرض به کار رود.

تحت این طرح نمونه‌ای تکرار شده، هر مصاحبه‌گر را می‌توان به صورت خوشه‌ای متشکل از \bar{n} مشاهده نیز در نظر گرفت. با این تفسیر، همبستگی میان رده‌ای را می‌توان به عنوان نشانگری از نسبت واریانس کل ناشی از تغییرپذیری در میان مصاحبه‌گران از لحاظ فرایند اندازه‌گیری نیز به کار برد. δ_x ، برآوردگر ضریب همبستگی میان رده‌ای که توسط کالتون [۱۹] و دیگران پیشنهاد شده است به صورت زیر است:

$$\delta_x = \frac{\left(\frac{s_{bx}^2}{s_{wx}^2} - 1 \right)}{\left(\frac{s_{bx}^2}{s_{wx}^2} - 1 + \bar{n} \right)} \quad (5.12)$$

مثال تشریحی: فرض کنید نمونه‌ای متشکل از ۲۰ کاربر یک بسته نرم‌افزاری حسابداری از روی فهرست کسانی که در شرکت تولید کننده این نرم‌افزار ثبت‌نام کرده‌اند گرفته شده است. این ۲۰ کاربر به طور تصادفی به ۴ مصاحبه‌گر، یعنی هر مصاحبه‌گر ۵ کاربر، اختصاص داده شده بودند که از روی فهرست اختصاصی خود با هر کاربر مصاحبه تلفنی انجام می‌دادند. متغیر اصلی موردنظر، میزان رضایت از نرم‌افزار مزبور در یک مقیاس پنج درجه‌ای بود که «۱» نشانه نارضی و «۵» نشانه بسیار راضی بود. داده‌های مربوط به هر یک از این چهار مصاحبه در پایین نشان داده شده‌اند:

مصاحبه‌گر				
۴	۳	۲	۱	
۴	۳	۲	۱	
۳	۳	۲	۱	
۵	۴	۲	۲	
۵	۴	۳	۳	
۵	۵	۴	۲	
۴/۴۰	۳/۸۰	۲/۶۰	۱/۸۰	\bar{x}_i
۰/۸۰	۰/۷۰	۰/۸۰	۰/۷۰	s_{ix}^2

با $M = 4$ و $\bar{n} = 5$ داریم

$$\bar{x} = 3/15 \quad s_{bx}^2 = 6/85 \quad s_{wx}^2 = 0/75$$

نسبت F چنین است

$$\frac{s_{bx}^2}{s_{wx}^2} = \frac{6/85}{0/75} = 9/133$$

که با ۳ و ۱۶ درجه آزادی، بسیار معنی‌دار است. ضریب همبستگی میان رده‌ای، δ_x ، به صورت زیر به دست می‌آید

$$\delta_x = \frac{9/133 - 1}{9/133 - 1 + 5} = 0/6192$$

به این ترتیب می‌توان برآورد کرد که بیش از ۶۱ درصد از کل واریانس، ناشی از اثرهای مصاحبه‌گر است.



۳.۱۲ خلاصه

در این فصل از مطالب مربوط به برآورد کردن واریانسها یا خطاهای معیار آماره‌های حاصل از آمارگیریهای نمونه‌ای پیچیده بحث کردیم. به خصوص درباره برآوردهای واریانس حاصل از سه روش که در سطح گسترده‌ای مورد استفاده قرار می‌گیرند، یعنی خطی‌سازی، تکرار مکرر متعادل، و روش جک‌نایف به تفصیل بحث کردیم و مثالهایی آوردیم. مطالعات تجربی چنین نتیجه داده‌اند که هیچ روشی همواره بهتر از دیگر روشها عمل نکرده است و هر روش می‌تواند برآوردهایی نسبتاً معقول برای واریانس تهیه کند، به شرط آنکه اندازه‌های نمونه‌ای موردنظر به قدر کافی بزرگ باشند.

تمرین

۱.۱۲ \bar{x} را برآوردی از سطح میانگین متغیر x و \bar{y} را برآوردی از سطح میانگین متغیر y در نظر بگیرید که هر دو از نمونه‌گیری تصادفی ساده به دست آمده‌اند. فرض کنید $\bar{z} = \bar{x}\bar{y}$. از خطی‌سازی در یافتن عبارتی برای برآورد واریانس توزیع \bar{z} استفاده کنید.

۲.۱۲ در یک ایالت بزرگ ۳۴ بخش وجود دارند که در نواحی کلانشهری قرار ندارند. به منظور برآورد کردن کل تعداد بیمارانی که با سندروم اکتسابی نارسایی سیستم ایمنی (ایدز) طی سال تقویمی ۱۹۹۰ در بیمارستانهای این ۳۴ بخش پذیرش شده‌اند، بخشهای مزبور در شش طبقه گروه‌بندی می‌شوند و یک نمونه تصادفی ساده متشکل از دو بخش از هر طبقه انتخاب می‌شود. در هر یک از بخشهای نمونه، یک بیمارستان به صورت نمونه تصادفی ساده گرفته می‌شود و همه سوابق پزشکی بیمارستانی بیمارانی که با تشخیص بیماری ایدز ترخیص شده‌اند در هر بیمارستان نمونه بررسی، خلاصه‌برداری، و شمارش می‌شوند. در جدول زیر، بیمارستانها برحسب طبقه، بخش، و تعداد تخت فهرست شده‌اند.

طبقه‌بندی برای آمارگیری از پذیرش بیماران دچار ایدز

تخت	بیمارستان	بخش	طبقه
۷۲	۱	۱	۱
۸۷	۱	۲	۱
۱۰۴	۲	۲	۱
۳۴	۳	۲	۱
۱۷	۱	۳	۱
۱۴۰	۲	۳	۱
۱۰۴	۱	۴	۱

طبقه‌بندی برای آمارگیری از پذیرش بیماران دچار ایدز (ادامه)

طبقه	بخش	بیمارستان	تخت
۲	۱	۱	۴۴
۲	۲	۱	۹۹
۲	۳	۱	۸۶
۲	۳	۲	۹۱
۲	۳	۳	۵۳
۲	۴	۱	۴۸
۳	۱	۱	۱۷۱
۳	۱	۲	۸۵
۳	۲	۱	۱۰۸
۳	۳	۱	۹۹
۳	۴	۱	۱۳۱
۳	۴	۲	۱۸۲
۴	۱	۱	۴۸
۴	۱	۲	۱۸
۴	۲	۱	۵۰
۴	۳	۱	۴۲
۴	۴	۱	۳۸
۵	۱	۱	۴۲
۵	۲	۱	۵۴
۵	۳	۱	۴۵
۵	۴	۱	۳۴
۶	۱	۱	۳۹
۶	۱	۲	۵۹
۶	۲	۱	۷۶
۶	۳	۱	۶۸
۶	۳	۲	۶۸
۶	۴	۱	۶۵

الف. نشان دهید که چگونه تعداد پذیرش بیماران دچار ایدز را برای ناحیه ۳۴ بخشی با در نظر گرفتن تعداد تختها برآورد می‌کنید.

ب. بیمارستانهایی که عملاً در نمونه قرار گرفته‌اند همراه با تعداد پذیرش بیماران دچار ایدز که در هر بیمارستان نمونه شمارش شده‌اند در جدول زیر نشان داده شده‌اند. از روی

این داده‌ها، کل تعداد پذیرش بیماران مبتلا به ایدز را در سراسر ناحیه ۳۴ بخشی برآورد کنید.

بیمارستانها در نمونه

طبقة	بخش	بیمارستان	تخت	کل مبتلایان به ایدز
۱	۱	۱	۷۲	۲۰
۱	۲	۱	۸۷	۴۹
۲	۲	۱	۹۹	۳۸
۲	۴	۱	۴۸	۲۳
۳	۳	۱	۹۹	۳۸
۳	۴	۱	۱۳۱	۷۸
۴	۳	۱	۴۲	۷
۴	۴	۱	۳۸	۲۸
۵	۱	۱	۴۲	۲۶
۵	۴	۱	۳۴	۹
۶	۱	۱	۳۹	۱۸
۶	۲	۱	۷۶	۲۰

پ. از مجموعه‌ای از نیمه نمونه‌های متعادل برای برآورد کردن خطای معیار برآورد کل تعداد پذیرش مبتلایان به ایدز برای ناحیه ۳۴ بخشی استفاده کنید. از ماتریس پلاکت - برمن که در پایین ارائه می‌شود می‌توان برای ساختن برآوردهای نیمه نمونه‌ای استفاده کرد (مک‌کارتی [۶، ص. ۱۷]):

تکرار	۱	۲	۳	۴	۵	۶
۱	+	-	-	+	-	+
۲	+	+	-	-	+	-
۳	+	+	+	-	-	+
۴	-	+	+	+	-	-
۵	+	-	+	+	+	-
۶	-	+	-	+	+	+
۷	-	-	+	-	+	+
۸	-	-	-	-	-	-

ت. از روش جک‌نایف برای به دست آوردن خطای معیار برآورد کل تعداد پذیرش مبتلایان به ایدز در ناحیه ۳۴ بخشی استفاده کنید. این نتیجه را با برآورد حاصل از روش تکرار مکرر متعادل در قسمت (پ) مقایسه کنید.

۳.۱۲ یک آمارگیری نمونه‌ای از فروشگاههای عرضه ویدیو به عمل آمد که هدف از آن برآورد کردن نسبت این قبیل فروشگاهها بود که براساس برنامه‌های پیش پرداخت حق عضویت، انگیزه‌هایی برای مشتریان فراهم می‌کنند. یک نمونه تصادفی ساده متشکل از ۴۸ فروشگاه از این نوع گرفته شد و به هر یک از ۶ مصاحبه‌گر ۸ فروشگاه به طور تصادفی، برای مصاحبه تلفنی، اختصاص داده شد. داده‌های زیر به دست آمدند:

تعداد فروشگاههای دارای برنامه عضویت از پیش پرداخت شده	مصاحبه‌گر
۷	۱
۱	۲
۰	۳
۸	۴
۷	۵
۰	۶

آیا نشانه‌ای از اثر مصاحبه‌گر دیده می‌شود؟ اگر پاسخ «بلی» است، بزرگی این اثر چقدر است؟

۴.۱۲ از روش جک‌نایف در برآورد کردن خطای معیار نسبت اضافه پرداختی برای داده‌های جدول ۱.۱۲ استفاده کنید. نتیجه حاصل چگونه با خطای معیار حاصل از روش خطی‌سازی مقایسه می‌شود؟

پیوست فنی

X'_h و $X'_{h\gamma}$ را برآوردهای نارایب مستقلی از دو واحد نمونه‌گیری اولیه برای پارامتر X_h (مانند میانگین، مجموع) در طبقه h ($h=1, \dots, L$) در نظر می‌گیریم و فرض می‌کنیم که $X'_{(1)}, \dots, X'_{(K)}$ مجموعه‌ای متعادل از برآوردهای نیمه نمونه‌ای باشد. فرض می‌کنیم که L به K قابل قسمت است و برآورد \hat{X} را به صورت زیر تعریف می‌کنیم:

$$\begin{aligned}\hat{X} &= \frac{\sum_{k=1}^K X'_{(k)}}{K} \\ &= \left(\frac{1}{K}\right) \sum_{h=1}^L \frac{K}{\gamma} (X'_{h1} + X'_{h\gamma}) \\ &= \sum_{h=1}^L \frac{(X'_{h1} + X'_{h\gamma})}{\gamma}\end{aligned}$$

واریانس $\hat{Var}(\hat{X})$ را به صورت زیر تعریف می‌کنیم:

$$\hat{Var}(\hat{X}) = \frac{\sum_{k=1}^K (X_{(k)} - \hat{X})^2}{K}$$

نشان خواهیم داد که $\hat{Var}(\hat{X})$ برآوردی نارایب از $Var(\hat{X})$ است.

به پنج نتیجه زیر توجه کنید

۱.

$$\begin{aligned}E(X'_{h1} - X'_{h\gamma})^2 &= Var(X'_{h1} - X'_{h\gamma}) + E^2(X'_{h1} - X'_{h\gamma}) \\ &= Var(X'_{h1} - X'_{h\gamma}) \\ &= Var(X'_{h1}) + Var(X'_{h\gamma})\end{aligned}$$

$$E\left[\sum_{h=1}^L (X'_{h1} - X'_{h\gamma})^2\right] = \sum_{h=1}^L [Var(X'_{h1}) + Var(X'_{h\gamma})] \quad ۲.$$

(این نتیجه پیامد ۱ است.)

۳.

$$\begin{aligned}E(\hat{X}) &= \left(\frac{1}{K}\right) \sum_{k=1}^K \sum_{h=1}^L E[\delta_{kh} X'_{h1} + (1 - \delta_{kh}) X'_{h\gamma}] \\ &= \left(\frac{1}{K}\right) \sum_{k=1}^K \sum_{h=1}^L X_h \\ &= \left(\frac{1}{K}\right) K \times X \\ &= X\end{aligned}$$

.۴

$$\begin{aligned} \text{Var}(\hat{X}) &= \left(\frac{1}{K^2}\right) \text{Var}\left(\sum_{k=1}^K \sum_{h=1}^L [\delta_{kh} X'_{h1} + (1-\delta_{kh}) X'_{h2}]\right) \\ &= \left(\frac{1}{K^2}\right) \sum_{k=1}^K \sum_{h=1}^L \text{Var}[\delta_{kh} X'_{h1} + (1-\delta_{kh}) X'_{h2}] \end{aligned}$$

(این عبارت آخر به این علت درست است که در یک مجموعه متعادل از نیمه نمونه‌ها، مجموع جملات حاصلضرب متقاطع صفر می‌شود.)

$$= \left(\frac{1}{K^2}\right) \sum_{h=1}^L \text{Var}\left[\left(\frac{K}{2}\right) X'_{h1} + \left(\frac{K}{2}\right) X'_{h2}\right]$$

(زیرا در یک مجموعه از نیمه نمونه‌های متعادل، هر یک از دو برآورد طبقه‌ای به تعداد دفعات $\left[\frac{K}{2}\right]$ ظاهر می‌شود.)

$$= \left(\frac{1}{4}\right) \sum_{h=1}^L (\text{Var}(X'_{h1}) + \text{Var}(X'_{h2}))$$

.۵

$$\begin{aligned} E\left[\sum_{k=1}^K (X'_{(k)} - \hat{X})^2\right] &= \sum_{k=1}^K E(X'_{(k)} - \hat{X})^2 \\ &= \sum_{k=1}^K \left[\text{Var}(X'_{(k)} - \hat{X}) + [E(X'_{(k)} - \hat{X})]^2 \right] \\ &= \sum_{k=1}^K [\text{Var}(X'_{(k)} - \hat{X}) + 0] \end{aligned}$$

(از روی نتیجه ۳ و این واقعیت که $E(X'_{(k)}) = X$)

$$\begin{aligned} &= \sum_{k=1}^K \text{Var}\left[\sum_{h=1}^L \left(\delta_{kh} X'_{h1} + (1-\delta_{kh}) X'_{h2} - \frac{X'_{h1} + X'_{h2}}{2}\right)\right] \\ &= \sum_{k=1}^K \text{Var}\left[\sum_{h=1}^L \left[\left(\delta_{kh} - \frac{1}{2}\right) X'_{h1} + \left(\delta_{kh} - \frac{1}{2}\right) X'_{h2}\right]\right] \\ &= \sum_{k=1}^K \sum_{h=1}^L \left(\delta_{kh} - \frac{1}{2}\right)^2 \text{Var}(X'_{h1} - X'_{h2}) \\ &= \sum_{k=1}^K \frac{1}{4} \sum_{h=1}^L \text{Var}(X'_{h1} - X'_{h2}) \\ &= \frac{K}{4} \sum_{h=1}^L [\text{Var}(X'_{h1}) + \text{Var}(X'_{h2})] \\ &= K \times \text{Var}(\hat{X}) \end{aligned}$$

بنابراین،

$$\begin{aligned} E[Var(\hat{X})] &= E\left[\left(\frac{1}{K}\right)\sum_{k=1}^K (X'_{(k)} - \hat{X})^2\right] \\ &= \left(\frac{1}{K}\right) \times K \times Var(\hat{X}) \\ &= Var(\hat{X}) \end{aligned}$$

برای این مورد خاص نشان دادیم که $Var(\hat{X})$ برآوردگری نااریب از $Var(\hat{X})$ است. با این که برآوردهای ناخطی از قبیل برآوردهای نسبتی و حاصلضربی نااریب نیستند، اریبی آنها غالباً ناچیز است و استدلالی که در بالا عرضه شد غالباً برای این برآوردهای ناخطی به صورت تقریبی مصداق پیدا می‌کند که موجب استفاده از برآوردگر واریانسی نیمه نمونه‌ای متعادل برای این برآوردهای ناخطی می‌شود.

کتابشناسی

The following publications develop linearization methodology.

1. Keyfitz, N., Estimates of sampling variance where two units are selected from each stratum. *Journal of the American Statistical Association*, 52: 503, 1957.
2. Woodruff, R. S., A simple method for approximating the variance of a complicated estimate. *Journal of the American Statistical Association*, 66: 411, 1971.

The following publications are relevant to software for implementation of the linearization method.

3. Research Triangle Institute, *SUDAAN: Professional Software for SURvey DATA ANalysis*, Research Triangle Institute, Research Triangle Park, N.C., 1989.
4. SAS Institute, Inc., *SAS Users Guide: Statistics*, SAS Institute Inc., Cary, N.C., 1982.

The following publications are relevant to the discussion of replication methods in this chapter.

5. Hansen, M. H., Hurwitz, W. N., and Madow, W. G., *Sample Survey Methods and Theory*, Vol. 1, Wiley, New York, 1953.
6. McCarthy, P. J., *Replication: An Approach to the Analysis of Data from Complex Surveys*, Vital and Health Statistics Series 2, No. 14, DHEW PUB No. (HSM) 73 – 1260, Health Services and Mental Health Administration, U.S. Government Printing Office, Washington, D.C., 1966.
7. McCarthy, P. J., Pseudo-replication: Half samples. *Review of the International Statistical Institute*, 37: 239, 1969.
8. Frankel, M. R., *Inference from Sample Surveys: An Empirical Investigation*, Institute for Social Research, Ann Arbor, MI., 1971.
9. Bean, J. A., *Distribution and Properties of Variance Estimators for Complex Multistage Probability Samples: An Empirical Distribution*, Vital and Health

- Statistics Series 2, No. 65, DHEW Pub. No. (HRA) 75-1339, Health Resources Administration, U.S. Government Printing Office, Washington, D.C., 1975.
10. Lemeshow, S., and Levy, P. S., Estimating the variances of ratio estimates in complex sample surveys with two primary sampling units per stratum-A comparison of balanced replication and jackknife techniques. *Journal of Statistical Computation and Simulation*, 8: 191, 1979.
 11. Lemeshow, S., Half-sample techniques. In *Encyclopedia of Statistical Sciences*, Vol. 3, Wiley, New York, 1983.
 12. Lemeshow, S., and Stoddard, A. M., A comparison of alternative estimation strategies for estimating the slope of a linear regression in sample surveys. *Communications in Statistics B: Simulation and Computation*, 13: 153, 1984.
 13. Lemeshow, S., and Epp, R., Properties of the balanced half-sample and jackknife variance estimation techniques in the linear case. *Communications in Statistics: Theory and Methods*, A6(13): 1259, 1977.
 14. Stanek, E., and Lemeshow, S., The behavior of balanced half-sample variance estimates for linear and combined ratio estimates when strata are paired to form pseudo-strata. *Estadística*, 32: 71, 1978.
 15. Lemeshow, S., The use of unique statistical weights for estimating variances with the balanced half-sample technique. *Journal of Statistical Planning and Inference*, 3: 315, 1979.
 16. Lemeshow, S., Hosmer, D. W., and Hislop, D., The effect of non-normality on estimating the variance of the combined ratio estimate. *Communications in Statistics: Simulation and Computation*, B9(4): 371, 1980.
 17. Plackett, R. L., and Burman, J. P., The design of optimal multifactorial experiments. *Biometrika*, 33: 305, 1946.

The following references contain discussions of the replicated sample designs or interpenetrating samples.

18. Deming, W. E., *Sample Design in Business Research*, Wiley, New York, 1960.
19. Kalton, G., *Introduction to Survey Sampling*, Sage Publications, Newbury Park, New York, 1983.

The following is considered the original article on jackknife estimation.

20. Quenouille, M. H., Note on bias in estimation. *Biometrika*, 43, 353-360, 1956.

The following recent review articles are relevant to topics covered in this chapter.

21. Frankel, M. R., Resampling procedures in estimation of sampling error. In *Encyclopedia of Biostatistics*, Armitage, P. A., and Colton, T. D., Eds., Wiley, Chichester, U.K., 1998.
22. Carlson, B. L., Software for statistical analysis of sample survey data. In *Encyclopedia of Biostatistics*, Armitage, P. A., and Colton, T. D., Eds., Wiley, Chichester, U.K., 1998.
23. Shah, B. V., Linearization methods of variance estimation. In *Encyclopedia of Biostatistics*, Armitage, P. A., and Colton, T. D., Eds., Wiley, Chichester, U.K., 1998.
24. Ghosh, S., Interpenetrating samples. In *Encyclopedia of Biostatistics*, Armitage, P. A., and Colton, T. D., Eds., Wiley, Chichester, U.K., 1998.

The following articles show consequences of not using appropriate methods of analysis when analyzing sample survey data.

25. Brogan, D., Pitfalls of using standard statistical software packages for sample

survey data. In *Encyclopedia of Biostatistics*, Armitage, P. A., and Colton, T. D., Eds., Wiley, Chichester, U.K., 1998.

26. Lemeshow, S., Letenneur, L., Lafont, S., Orgogozo, J. M., and Commenges, D., An Illustration of Analysis Taking into Account Complex Survey Considerations: The Association Between Wine Consumption and Dementia in the PAQUID study. *The American Journal of Epidemiology*, 148 No 3., 298-306, 1998.

In addition to the above literature, Dr Alan Zaslafsky of Harvard University maintains a home page on the Worldwide Web which features information on software for analysis of data from complex sample surveys. The address of this site is: <http://www.fas.harvard.edu/~stats/survey-soft/survey-soft.html/>