

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 ۵

طبقه‌بندی و نمونه‌گیری تصادفی

طبقه‌بندی شده

در دو نوع نمونه‌گیری که تا اینجا مورد بحث قرار گرفتند یعنی در نمونه‌گیری تصادفی ساده و نمونه‌گیری سیستماتیک، لازم است که نمونه‌ای از جامعه به صورت یک کل گرفته شود و در هیچ یک از آنها لازم نیست که زیرحوزه‌ها یا زیرگروه‌های جامعه قبل از گرفتن نمونه شناسایی شوند. ولی گاهی اوقات می‌توان چارچوب نمونه‌گیری را به گروه‌ها یا طبقه‌هایی افراز و نمونه‌گیری را جداگانه در داخل هر طبقه اجرا کرد. این طرح نمونه‌گیری را نمونه‌گیری طبقه‌بندی شده می‌نامند. اگر برای انتخاب نمونه در داخل هر طبقه از نمونه‌گیری تصادفی ساده استفاده شود طرح نمونه را نمونه‌گیری تصادفی طبقه‌بندی شده می‌نامند.

در این فصل برخی «مبانی» طبقه‌بندی و نمونه‌گیری تصادفی طبقه‌بندی شده را معرفی می‌کنیم. اینها شامل بحث درباره روشهای گرفتن نمونه، تعریف اصطلاحها و نمادگذاریهای مشخصه‌های جامعه و نمونه که در بحث طرحهای نمونه با طبقه‌بندی به کار رفته‌اند، و برآوردهای بنیادی مشخصه‌های جامعه خواهند بود. در فصل بعد با جزئیات بیشتری به روشهای استفاده از طرحهای طبقه‌بندی از راههایی که احتمال تولید قابل اعتمادترین برآوردها بیشتر است می‌پردازیم:

۱.۵ نمونه تصادفی طبقه‌بندی شده چیست؟

همان طور که در بالا اشاره شد، نمونه تصادفی طبقه‌بندی شده یک برنامه نمونه‌گیری است که در آن جامعه به L طبقه دو به دو ناسازگار و جامع طبقه‌بندی می‌شود، و یک نمونه تصادفی ساده متشکل از n_h عنصر از داخل هر طبقه h گرفته می‌شود. نمونه‌گیری در داخل هر طبقه به طور مستقل اجرا می‌شود. در واقع می‌توانیم طرح نمونه‌گیری تصادفی ساده را به صورت L نمونه تصادفی ساده مجزا تصور کنیم.

از نظر عملیاتی، نمونه تصادفی طبقه‌بندی شده به همان طریق نمونه تصادفی ساده گرفته می‌شود ولی نمونه‌گیری در داخل هر طبقه به طور جداگانه و مستقل انجام می‌گیرد. اگر N_1, \dots, N_p, N_L معرف تعداد واحدهای نمونه‌گیری در داخل هر طبقه و n_1, \dots, n_p, n_L معرف تعداد واحدهای نمونه‌گیری انتخاب شده به طور تصادفی در داخل هر طبقه باشند، آنگاه، تعداد کل نمونه‌های تصادفی طبقه‌بندی شده ممکن برابر خواهد بود با:

$$\binom{N_1}{n_1} \times \binom{N_2}{n_2} \times \dots \times \binom{N_L}{n_L}$$

که از $\binom{N}{n}$ ، کل تعداد نمونه‌های تصادفی ساده ممکن، کوچکتر بوده و یا با آن برابر است.

برای مثال، اگر سه طبقه داشته باشیم و $N_1 = 3$ ، $N_2 = 5$ و $N_3 = 6$ باشد تعداد کل نمونه‌های ممکن متشکل از $n_1 = 1$ عنصر از طبقه ۱، $n_2 = 2$ عنصر از طبقه ۲ و $n_3 = 4$ عنصر از طبقه ۳ چنین است:

$$\binom{3}{1} \times \binom{5}{2} \times \binom{6}{4} = 450$$

تعداد کل نمونه‌های تصادفی ساده متشکل از ۷ عنصر از ۱۴ عنصر جامعه عبارت است از

$$\binom{14}{7} = 3432$$

احتمال انتخاب شدن یک عنصر در نمونه بستگی دارد به طبقه خاصی که آن عنصر در آن گروه‌بندی شده است و می‌توان نشان داد که برابر است با $\frac{n_h}{N_h}$ (در صورتی که آن عنصر در طبقه h باشد). در

مثالی که هم اکنون شرح دادیم، احتمال انتخاب شدن یک عنصر برابر است با $\frac{1}{3}$ برای عناصر طبقه ۱،

$\frac{2}{5}$ برای عناصر طبقه ۲ و $\frac{4}{6}$ برای عناصر طبقه ۳.

۲.۵ چگونگی گرفتن نمونه تصادفی طبقه‌بندی شده

همان طور که در بالا شرح داده شد، نمونه تصادفی طبقه‌بندی شده صرفاً با گرفتن نمونه‌های تصادفی ساده متشکل از n_h عنصر ($h=1, \dots, L$) به طور مستقل در داخل هر طبقه تهیه می‌شود. انتخاب عملی نمونه، با استفاده از هر یک از روشهایی که قبلاً برای گرفتن نمونه تصادفی ساده شرح داده شد اجرا می‌شود.

مثلاً، اگر کسی بخواهد از گزارشهای ذخیره شده در پرونده اطلاعاتی رایانه نمونه‌ای بگیرد می‌تواند از یک الگوریتم برای گرفتن نمونه تصادفی طبقه‌بندی شده که در کتاب راهنمای SAS [۷] توصیف شده است استفاده کند. همچنین مدول تکمیلی SAMPLE, SYSTAT نیز می‌تواند برای گرفتن نمونه تصادفی طبقه‌بندی شده مورد استفاده قرار گیرد.

۳.۵ چرا نمونه‌گیری طبقه‌بندی شده؟

نمونه‌گیری طبقه‌بندی شده به این دلیل در انواع خاصی از آمارگیریها به کار می‌رود که آسانی درک نمونه‌گیری تصادفی ساده را با افزایش قابل توجهی در قابلیت اعتماد بالقوه ترکیب می‌کند. هر گاه بخواهیم برآوردهای جداگانه‌ای از پارامترهای جامعه برای هر یک از زیرحوزه‌ها در داخل کل جامعه به دست آوریم و علاوه بر آن بخواهیم مطمئن باشیم که نمونه ما نماینده جامعه است استفاده از این فن راحت است. مثلاً، فرض کنید می‌خواهیم تعداد کل تختها را در بیمارستانهای یک ایالت مشخص برآورد کنیم. می‌دانیم که اکثر بیمارستانها، اندازه‌های کوچک یا متوسط دارند و تنها چندتایی بیمارستانهای بزرگ موجودند. می‌دانیم که این بیمارستانهای بزرگ بخش قابل توجهی از کل تعداد تختها را دارند.

حالا فرض کنید تصمیم بر این است که یک نمونه تصادفی ساده از بیمارستانهای این ایالت انتخاب کنیم، تعداد تختها را در هر یک از نمونه‌هایی که به این ترتیب انتخاب شده‌اند تعیین، و با استفاده از روشهای فصل ۳، کل تعداد تختها را در تمام بیمارستانهای سراسر ایالت برآورد کنیم. مشکل این شیوه آن است که شانس زیادی وجود دارد که نمونه ما تعداد بسیار زیاد یا بسیار کمی از بیمارستانهای بسیار بزرگ را در خود جای دهد. در نتیجه ممکن است نمونه به صورتی مناسب معرف جامعه نباشد.

راه حل ما برای این مشکل آن است که واحدهای نمونه‌گیری (بیمارستانها) را قبل از نمونه‌گیری بر اساس اندازه (یعنی کوچک، متوسط، بزرگ) در سه گروه طبقه‌بندی کنیم و سپس با استفاده از فنون نمونه‌گیری تصادفی ساده، از هر یک از گروههای سه‌گانه تعداد معینی از بیمارستانها را انتخاب نماییم.

سپس برآورد تعداد کل تختها را می‌توان از ترکیب نتایج سه طبقه به دست آورد. این، اساس نمونه‌گیری طبقه‌بندی شده است.

برای نشان دادن ایده‌ها و مزایای طبقه‌بندی به یک مثال نگاهی می‌اندازیم.

مثال تشریحی: فرض کنید جاده‌ای با ۲۴ مایل طول از مناطقی عبور می‌کند که می‌توانند به صورت شهری و روستایی طبقه‌بندی شوند و این جاده به هشت قطعه تقسیم شده است که هر یک دارای طولی برابر ۳ مایل است. نمونه‌ای شامل سه قطعه گرفته شده و در هر قطعه نمونه‌گیری شده تجهیزات ویژه‌ای نصب گردیده است تا مجموع مایل‌های (مسافتهای) طی شده توسط وسایل نقلیه، اتومبیل و کامیون، در آن بخش طی یک سال خاص شمارش شود. به علاوه، گزارشی از تمام حوادثی که در هر قطعه نمونه اتفاق می‌افتد نیز ثبت شود.

مقدار مسافت طی شده توسط کامیونها و تعداد حوادثی که در آن کامیون دخالت داشته است طی یک دوره معین برای هر یک از قطعه‌های هشتگانه جامعه در جدول ۱.۵ ارائه شده است. فرض کنید یک نمونه تصادفی ساده متشکل از سه قطعه را به منظور برآورد کل تعداد کامیون - مایل طی شده در جاده گرفته‌ایم. برای یک جامعه متشکل از هشت قطعه، ۵۶ نمونه ممکن سه قطعه‌ای وجود دارند. توزیع نمونه‌گیری تعداد کامیون - مایل طی شده در جاده در جدول ۲.۵ ارائه شده است.

جدول ۱.۵ کامیون - مایل و تعداد حوادثی که در آن کامیون مداخله داشته است

بر حسب قطعه‌های جاده

قطعه	نوع	تعداد کامیون - مایل ($\times 100$)	تعداد حوادث با مداخله کامیون
۱	شهری	۶۳۲۷	۸
۲	روستایی	۲۵۵۵	۵
۳	شهری	۸۶۹۱	۹
۴	شهری	۷۸۳۴	۹
۵	روستایی	۱۵۸۶	۵
۶	روستای	۲۰۳۴	۱
۷	روستایی	۲۰۱۵	۹
۸	روستایی	۳۰۱۲	۴

جدول ۲.۵ توزیع نمونه‌گیری x' برای ۵۶ نمونه ممکن سه قطعه‌ای

x'	قطعه‌ها در نمونه	x'	قطعه‌ها در نمونه
۳۳۰۷۷/۳۳	(۲, ۴, ۷)	۴۶۸۶۱/۳۳	(۱, ۲, ۳)
۳۵۷۳۶	(۲, ۴, ۸)	۴۴۵۷۶	(۱, ۲, ۴)
۱۶۴۶۶/۶۷	(۲, ۵, ۶)	۲۷۹۱۴/۶۷	(۱, ۲, ۵)
۱۶۴۱۶	(۲, ۵, ۷)	۲۹۱۰۹/۳۳	(۱, ۲, ۶)
۱۹۰۷۴/۶۷	(۲, ۵, ۸)	۲۹۰۵۸/۶۷	(۱, ۲, ۷)
۱۷۶۱۰/۶۷	(۲, ۶, ۷)	۳۱۷۱۷/۳۳	(۱, ۲, ۸)
۲۰۲۶۹/۳۳	(۲, ۶, ۸)	۶۰۹۳۸/۶۷	(۱, ۳, ۴)
۲۰۲۱۸/۶۷	(۲, ۷, ۸)	۴۴۲۷۷/۳۳	(۱, ۳, ۵)
۴۸۲۹۶	(۳, ۴, ۵)	۴۵۴۷۲	(۱, ۳, ۶)
۴۹۴۹۰/۶۷	(۳, ۴, ۶)	۴۵۴۲۱/۳۳	(۱, ۳, ۷)
۴۹۴۴۰	(۳, ۴, ۷)	۴۸۰۸۰	(۱, ۳, ۸)
۵۲۰۹۸/۶۷	(۳, ۴, ۸)	۴۱۹۹۲	(۱, ۴, ۵)
۳۲۸۲۹/۳۳	(۳, ۵, ۶)	۴۳۱۸۶/۶۷	(۱, ۴, ۶)
۳۲۷۷۸/۶۷	(۳, ۵, ۷)	۴۳۱۳۶	(۱, ۴, ۷)
۳۵۴۳۷/۳۳	(۳, ۵, ۸)	۴۵۷۹۴/۶۷	(۱, ۴, ۸)
۳۳۹۷۳/۳۳	(۳, ۶, ۷)	۲۶۵۲۵/۳۳	(۱, ۵, ۶)
۳۶۶۳۲	(۳, ۶, ۸)	۲۶۴۷۴/۶۷	(۱, ۵, ۷)
۳۶۵۸۱/۳۳	(۳, ۷, ۸)	۲۹۱۳۳/۳۳	(۱, ۵, ۸)
۳۰۵۴۴	(۴, ۵, ۶)	۲۷۶۶۹/۳۳	(۱, ۶, ۷)
۳۰۴۹۳/۳۳	(۴, ۵, ۷)	۳۰۳۲۸	(۱, ۶, ۸)
۳۳۱۵۲	(۴, ۵, ۸)	۳۰۲۷۷/۳۳	(۱, ۷, ۸)
۳۱۶۸۸	(۴, ۶, ۷)	۵۰۸۸۰	(۲, ۳, ۴)
۳۴۳۴۶/۶۷	(۴, ۶, ۸)	۳۴۲۱۸/۶۷	(۲, ۳, ۵)
۳۴۲۹۶	(۴, ۷, ۸)	۳۵۴۱۳/۳۳	(۲, ۳, ۶)
۱۵۰۲۶/۶۷	(۵, ۶, ۷)	۳۵۳۶۲/۶۷	(۲, ۳, ۷)
۱۷۶۸۵/۳۳	(۵, ۶, ۸)	۳۸۰۲۱/۳۳	(۲, ۳, ۸)
۱۷۶۳۴/۶۷	(۵, ۷, ۸)	۳۱۹۳۳/۳۳	(۲, ۴, ۵)
۱۸۸۲۹/۳۳	(۶, ۷, ۸)	۳۳۱۲۸	(۲, ۴, ۶)

جدول ۳.۵ دو طبقه برای داده‌های جدول ۱.۵

طبقه ۲ (قطعات روستایی)		طبقه ۱ (قطعات شهری)	
قطعه	تعداد کامیون - مایل $\times 1000$	قطعه	تعداد کامیون - مایل $\times 1000$
۱	۲۵۵۵	۲	۶۳۲۷
۳	۱۵۸۶	۵	۸۶۹۱
۴	۲۰۳۴	۶	۷۸۳۴
	۲۰۱۵	۷	
	۳۰۱۲	۸	

برآوردهای کل تعداد کامیون - مایل که از این طریق به دست می‌آید از ۱۵۰۲۶/۶۷ تا ۶۰۹۳۸/۶۷ تغییر می‌کنند و میانگین توزیع نمونه‌گیری x' برابر است با ۳۴۰۵۴ یعنی مجموع جامعه، و خطای معیار x' برابر است با ۱۰۵۳۶/۹.

حالا فرض کنید که به جای گرفتن یک نمونه تصادفی ساده متشکل از سه قطعه از جامعه‌ای با هشت قطعه، اول قطعه‌ها را به دو طبقه تقسیم کنیم، یکی شامل قطعات شهری، دیگری شامل قطعات روستایی، به صورتی که در جدول ۳.۵ نشان داده شده است.

حالا می‌توانیم یک نمونه متشکل از یک قطعه از طبقه ۱ و دو قطعه از طبقه ۲ انتخاب کنیم و کل تعداد کامیون - مایل را به وسیله برآورد x'_{str} به صورت زیر برآورد کنیم:

$$x'_{str} = x'_1 + x'_2$$

که در آن

$$x'_1 = \text{تعداد کامیون - مایل برآورد شده در سه قطعه تشکیل دهنده طبقه ۱}$$

$$x'_2 = \text{تعداد کامیون - مایل برآورد شده در پنج قطعه تشکیل دهنده طبقه ۲}$$

از نظر مفهومی، هر طبقه را یک زیرجامعه در نظر می‌گیریم، در داخل هر زیرجامعه به طور مستقل نمونه می‌گیریم، و برآوردی برای کل جامعه با انبوه سازی برآورد تک تک طبقه‌ها در همه زیرجامعه‌ها یا طبقات به دست می‌آوریم. بنابراین، اگر یک نمونه تصادفی ساده از یک قطعه از طبقه ۱ و دو قطعه از طبقه ۲ بگیریم در مجموع $\binom{5}{2} \times \binom{3}{1} = 30$ نمونه ممکن به دست می‌آوریم. پس از آن اگر از شیوه برآورد که در بالا با معادله مربوط به x'_{str} ارائه شد استفاده کنیم توزیع نمونه‌گیری x'_{str} را به دست می‌آوریم که در جدول ۴.۵ نشان داده شده است.

$E(x'_{str})$ ، میانگین توزیع نمونه‌گیری x'_{str} برابر است با ۳۴۰۵۴ (تابلوی ۳.۲ را ببینید) که همان مجموع واقعی جامعه است. $SE(x'_{str})$ ، خطای معیار کل برآورد شده x'_{str} برابر است با ۳۲۹۷/۶ که

جدول ۴.۵ توزیع نمونه‌گیری x'_{str} برای ۳۰ نمونه ممکن از سه قطعه

$x'_{str} = (x'_1 + x'_2)$	$x'_1 (= 5\bar{x}_1)$	$x'_2 (= 3\bar{x}_2)$	طبقه ۲	طبقه ۱
۲۹۳۳۳/۵	۱۰۳۵۲/۵	۱۸۹۸۱	(۲,۵)	۱
۳۰۴۵۳/۵	۱۱۴۷۲/۵	۱۸۹۸۱	(۲,۶)	۱
۳۰۴۰۶	۱۱۴۲۵	۱۸۹۸۱	(۲,۷)	۱
۳۲۸۹۸/۵	۱۳۹۱۷/۵	۱۸۹۸۱	(۲,۸)	۱
۲۸۰۳۱	۹۰۵۰	۱۸۹۸۱	(۵,۶)	۱
۲۷۹۸۳/۵	۹۰۰۲/۵	۱۸۹۸۱	(۵,۷)	۱
۳۰۴۷۶	۱۱۴۹۵	۱۸۹۸۱	(۵,۸)	۱
۲۹۱۰۳/۵	۱۰۱۲۲/۵	۱۸۹۸۱	(۶,۷)	۱
۳۱۵۹۶	۱۲۶۱۵	۱۸۹۸۱	(۶,۸)	۱
۳۱۵۴۸/۵	۱۲۵۶۷/۵	۱۸۹۸۱	(۷,۸)	۱
۳۶۴۲۵/۵	۱۰۳۵۲/۵	۲۶۰۷۳	(۲,۵)	۳
۳۷۵۴۵/۵	۱۱۴۷۲/۵	۲۶۰۷۳	(۲,۶)	۳
۳۷۴۹۸	۱۱۴۲۵	۲۶۰۷۳	(۲,۷)	۳
۳۹۹۹۰/۵	۱۳۹۱۷/۵	۲۶۰۷۳	(۲,۸)	۳
۳۵۱۲۳	۹۰۵۰	۲۶۰۷۳	(۵,۶)	۳
۳۵۰۷۵/۵	۹۰۰۲/۵	۲۶۰۷۳	(۵,۷)	۳
۳۷۵۶۸	۱۱۴۹۵	۲۶۰۷۳	(۵,۸)	۳
۳۶۱۹۵/۵	۱۰۱۲۲/۵	۲۶۰۷۳	(۶,۷)	۳
۳۸۶۸۸	۱۲۶۱۵	۲۶۰۷۳	(۶,۸)	۳
۳۸۶۴۰/۵	۱۲۵۶۷/۵	۲۶۰۷۳	(۷,۸)	۳
۳۳۸۵۴/۵	۱۰۳۵۲/۵	۲۳۵۰۲	(۲,۵)	۴
۳۴۹۷۴/۵	۱۱۴۷۲/۵	۲۳۵۰۲	(۲,۶)	۴
۳۴۹۲۷	۱۱۴۲۵	۲۳۵۰۲	(۲,۷)	۴
۳۷۴۱۹/۵	۱۳۹۱۷/۵	۲۳۵۰۲	(۲,۸)	۴
۳۲۵۵۲	۹۰۵۰	۲۳۵۰۲	(۵,۶)	۴
۳۲۵۰۴/۵	۹۰۰۲/۵	۲۳۵۰۲	(۵,۷)	۴
۳۴۹۹۷	۱۱۴۹۵	۲۳۵۰۲	(۵,۸)	۴
۳۳۶۲۴/۵	۱۰۱۲۲/۵	۲۳۵۰۲	(۶,۷)	۴
۳۶۱۱۷	۱۲۶۱۵	۲۳۵۰۲	(۶,۸)	۴
۳۶۰۶۹/۵	۱۲۵۶۷/۵	۲۳۵۰۲	(۷,۸)	۴

خیلی کمتر از خطای معیار x' ، یعنی برآورد معادل کل جامعه تحت نمونه‌گیری تصادفی ساده است. این نتایج در جدول ۵.۵ خلاصه شده‌اند.

به صورت شهودی می‌توان دید که چرا برنامه نمونه‌گیری مبتنی بر طبقه‌بندی برآوردی را به دست می‌دهد که خطای معیار آن کمتر از برآورد نظیر بر مبنای نمونه‌گیری تصادفی ساده است. برای x' ، برآورد مجموع، با استفاده از نمونه‌گیری تصادفی ساده ۵۶ مقدار ممکن به دست می‌آید در حالی که از طبقه‌بندی ۳۰ مقدار برای x'_{str} به دست می‌آید. همچنین، امتحان کردن دامنه تغییرات دو توزیع نشان می‌دهد که طبقه‌بندی، آن دسته از نمونه‌ها را که منجر به برآوردهای بسیار زیاد یا بسیار کم برای مجموع می‌شد حذف کرده است. سه قطعه شهری دارای مقادیر زیاد و پنج قطعه روستایی دارای مقادیر کم مشخصه‌های مورد اندازه‌گیری (کامیون - مایل) بودند. طبقه‌بندی اطمینان می‌دهد که حداقل یک قطعه شهری و یک قطعه روستایی در نمونه انتخاب می‌شوند و به این ترتیب امکان برآوردهای بسیار زیاد و بسیار کم مجموع را از میان می‌برد.

جدول ۵.۵ مقایسه نتایج نمونه‌گیری تصادفی ساده و طبقه‌بندی

طرح طبقه‌بندی		
نمونه‌گیری تصادفی ساده	طبقه‌بندی	
۳	*۳	تعداد عناصر در نمونه
۵۶	۳۰	تعداد نمونه‌های ممکن
۳۴۰۵۴ ⁺	۳۴۰۵۴ ⁺	میانگین توزیع برآورد مجموع
۱۰۵۳۶/۹	۳۲۹۷/۶	خطای معیار برآورد مجموع
۴۵۹۱۲	۱۲۰۰۷	دامنه تغییرات توزیعهای برآورد مجموعها

* یک عنصر از طبقه ۱ و دو عنصر از طبقه ۲.

⁺ این مقدار همان مجموع جامعه است.

□

طبقه‌بندی نسبت به نمونه‌گیری تصادفی ساده، سه امتیاز عمده دارد.

۱. با شرایط معین مفروض، دقت آن نسبت به نمونه‌گیری تصادفی ساده ممکن است افزایش پیدا کند (یعنی ممکن است خطاهای معیار به دست آمده از این شیوه برآورد، کمتر باشند).
۲. امکان به دست آوردن برآوردهایی با دقت مشخص برای هر یک از طبقه‌ها وجود دارد.
۳. ممکن است جمع‌آوری اطلاعات برای نمونه طبقه‌بندی شده نیز به دلایل سیاسی یا اداری به همان آسانی جمع‌آوری اطلاعات برای یک نمونه تصادفی ساده باشد. در این صورت، با

گرفتن نمونه طبقه‌بندی شده چیزی را از دست نمی‌دهیم زیرا، خطاهای معیار حاصل از این طرح به ندرت ممکن است از خطاهای معیار نمونه‌گیری تصادفی ساده بیشتر شود.

راهبرد مورد استفاده برای ساختن طبقات مستلزم دو گام است. ابتدا پارامتر مورد نظر جامعه برای برآورد کردن را تعیین می‌کنیم. سپس جامعه را نسبت به متغیر دیگری که تصور می‌شود با متغیر مورد نظر مربوط باشد طبقه‌بندی می‌کنیم. اگر فرض ما راجع به این پیوند درست باشد این گام دوم تضمین می‌کند که طبقه‌ها نسبت به متغیر تحت بررسی همگن‌اند.

برای مثال، اگر بخواهیم تعداد تختها را در تمام بیمارستانهای ایالت برآورد کنیم ممکن است بخواهیم بیمارستانها را از لحاظ زیربنا (که اطلاعی است که به خاطر مالیات به راحتی قابل دسترسی است) با این استدلال طبقه‌بندی کنیم که بیمارستانهایی که فضای بیشتری دارند تختهای بیشتری خواهند داشت. از سوی دیگر، اگر به متوسط هزینه روزانه هر تخت بیمارستانی به ازای هر بیمار علاقه‌مندیم می‌توانیم طبقه‌بندی بیمارستانها را بر مبنای محل جغرافیایی آنان با این استدلال انجام دهیم که بیمارستانهای مناطقی از ایالت که دارای اقتصاد پر رونقی‌اند می‌توانند هزینه‌های سنگینتری از بیماران خود مطالبه کنند تا بیمارستانهایی که در مناطق با اقتصاد کساد واقع شده‌اند.

در بیشتر موقعیتهای عملی، طبقه‌بندی جامعه نسبت به متغیر تحت بررسی صرفاً به دلایل مربوط به هزینه و عملی بودن، کاری مشکل است. در بسیاری از موارد جامعه به راحتی‌ترین طریق با استفاده از معیارهای اداری (مثل نواحی رأی‌گیری)، معیارهای جغرافیایی (مثل شمال، جنوب، شرق یا غرب) یا سایر معیارهای طبیعی (مثل جنس یا سن) طبقه‌بندی می‌شود. طبقه‌بندی بر اساس راحتی غیرمنطقی نیست زیرا برآورد یک تک پارامتر در آمارگیری مدرن متداول نیست. در عوض، اطلاعات زیادی در مورد هر واحد نمونه‌گیری جمع‌آوری می‌شود و پارامترهای زیادی مورد توجه‌اند. واضح است که آن چه برای یک متغیر می‌تواند راهبرد طبقه‌بندی بهینه‌ای برای تهیه طبقات نسبتاً همگن باشد ممکن است نسبت به یک متغیر دیگر، طبقات بسیار ناهمگن فراهم نماید. اهمیت دارد که آماردان قبل از تصمیم‌گیری در مورد ملاک مناسب برای طبقه‌بندی، دامنه داده‌هایی را که قرار است جمع‌آوری شوند مورد توجه قرار دهد. این قبیل مطالب گاهی اوقات در گزارشهای آمارگیریهای بهداشتی که از طبقه‌بندی استفاده کرده‌اند مورد بحث قرار می‌گیرند.

عیب عمده نمونه‌گیری طبقه‌بندی شده آن است که نیازمند شناسایی یکایک واحدهای شمارش بر حسب طبقه پیش از اجرای نمونه‌گیری است. اگر چنین اطلاعاتی به سهولت فراهم نباشد این روش به ندرت ممکن است عملی باشد.

ممکن است بخواهید طبقه‌بندی را با راهبرد نمونه‌گیری خوشه‌ای (فصلهای ۸ تا ۱۱) مقایسه کنید که در آن ممکن است صرفه‌جویی‌هایی عمده در وقت و هزینه امکان‌پذیر شود. مع‌هذا، به دلیل مزایایی که در بالا شرح داده شد، طبقه‌بندی فن بسیار توانایی است که در سطح وسیعی مورد استفاده قرار می‌گیرد.

۴.۵ پارامترهای جامعه برای طبقات

پارامترهای جامعه برای طبقات را می‌توان با همان نمادهایی تعریف کرد که برای تعریف پارامترهای جامعه به طور کلی مورد استفاده قرار دادیم. در این بخش نمادهایی را معرفی می‌کنیم که در سراسر بحث از برنامه‌های نمونه‌گیری بر مبنای طبقه‌بندی به کار خواهند رفت.

جامعه‌ای را در نظر می‌گیریم که دارای N واحد اولیه است که در L طبقه دو به دو ناسازگار و جامع به گونه‌ای گروه‌بندی شده‌اند که طبقه ۱ دارای N_1 واحد اولیه، طبقه ۲ دارای N_2 واحد اولیه، ...، و طبقه L دارای N_L واحد اولیه است. به عبارت دیگر، $N = \sum_{h=1}^L N_h$ اندازه جامعه است. فرض کنید یک متغیر یا مشخصه X را در جامعه در نظر گرفته‌ایم. از نماد $X_{h,i}$ برای نشان دادن مقدار مشخصه X برای i امین واحد اولیه در داخل طبقه h استفاده می‌کنیم. به عبارت دیگر، واحدهای اولیه داخل هر طبقه ویژه h از ۱ تا N_h شماره‌گذاری می‌شوند.

مثال تشریحی: فرض کنید می‌خواهیم متوسط هزینه روزانه دارو به ازای هر بیمار را در یک بیمارستان برآورد کنیم. تصمیم می‌گیریم بیمارستان را بر حسب خدمات (درمانی، جراحی، زنان و زایمان، و همه خدمات دیگر با هم) طبقه‌بندی کنیم و واحدهای اولیه را به صورت بیماران در یک روز معین تعریف می‌کنیم.

فرض کنید که در یک روز تعیین شده ۲۵۰ بیمار در بیمارستان است که ۱۰۰ نفر برای درمان، ۷۵ نفر برای جراحی، ۵۰ نفر زنان و زایمان و ۲۵ نفر برای سایر خدمات مراجعه کرده‌اند. پس با استفاده از نمادهایی که در بالا معرفی شدند داریم:

$$N = 250 \quad N_1 = 100 \quad N_2 = 75 \quad N_3 = 50 \quad N_4 = 25$$

اگر $X_{h,i}$ نشان‌دهنده مقدار متغیر X برای i امین واحد اولیه در داخل طبقه h باشد، آنگاه مثلاً برای طبقه ۲ داریم:

$$X_{2,1} = \text{مقدار متغیر } X \text{ برای عنصر } 1 \text{ در داخل طبقه } 2$$

$$X_{2,2} = \text{مقدار متغیر } X \text{ برای عنصر } 2 \text{ در داخل طبقه } 2$$

و به همین ترتیب تا می‌رسیم به

$$X_{r,75} = \text{مقدار متغیر } x \text{ برای عنصر } 75 \text{ در داخل طبقه } 2$$

عناصر داخل سایر طبقه‌ها نیز به همین ترتیب تعریف می‌شوند.

□

پارامترهای جامعه برای طبقات را به روشی تعریف می‌کنیم که پارامترهای جامعه کل را تعریف کردیم.

مجموع یا مقدار تجمعی متغیر x در داخل یک طبقه h را با نماد X_{h+} نشان می‌دهیم که از فرمول زیر به دست می‌آید:

$$X_{h+} = \sum_{i=1}^{N_h} X_{h,i}$$

مجموع کل جامعه با جمع کردن مجموع طبقه‌ها به دست می‌آید، یا

$$X = \sum_{h=1}^L \sum_{i=1}^{N_h} X_{h,i} = \sum_{h=1}^L X_{h+}$$

سطح میانگین یک مشخصه x برای طبقه h با نماد \bar{X}_h نشان داده می‌شود و از فرمول زیر به دست می‌آید:

$$\bar{X}_h = \frac{\sum_{i=1}^{N_h} X_{h,i}}{N_h} = \frac{X_{h+}}{N_h}$$

میانگین \bar{X} متغیر x برای تمام جامعه از فرمول زیر

$$\bar{X} = \frac{X}{N}$$

یا معادل جبری آن،

$$\bar{X} = \frac{\sum_{h=1}^L N_h \bar{X}_h}{N} = \sum_{h=1}^L W_h \bar{X}_h$$

به دست می‌آید، که در آن،

$$W_h = \frac{N_h}{N}$$

به عبارت دیگر، میانگین \bar{X} برای تمام جامعه عبارت است از متوسط موزون میانگینهای تکی طبقه‌ها،

$$\bar{X}_h, \text{ با وزنه‌ای } (W_h = \frac{N_h}{N}) \text{ متناسب با تعداد عناصر در هر طبقه.}$$

واریانس σ_{hx}^2 توزیع متغیر x در داخل یک طبقه خاص h به صورت متوسط توان دوم انحرافها

حول میانگین طبقه تعریف و با فرمول زیر نشان داده می‌شود:

$$\sigma_{hx}^2 = \frac{\sum_{i=1}^{N_h} (X_{h,i} - \bar{X}_h)^2}{N_h}$$

ضرایب تغییرات و واریانسهای نسبی برای هر طبقه به همان صورتی تعریف می‌شوند که برای جامعه‌ای که به طبقات گروه‌بندی نشده است تعریف شد. ضریب تغییرات برای توزیع در داخل یک طبقه خاص از فرمول زیر به دست می‌آید:

$$V_{hx} = \frac{\sigma_{hx}}{\bar{X}_h}$$

و واریانس نسبی صرفاً عبارت است از توان دوم ضریب تغییرات.

برای راحتی کار در ارجاعهای بعدی، پارامترهای جامعه برای نمونه‌گیری طبقه‌بندی شده در تابلوی ۱.۵ خلاصه شده‌اند. حالا ببینیم این فرمولها در عمل چگونه به کار می‌روند. مثال تشریحی: جامعه‌ای متشکل از ۱۴ خانواده را در نظر بگیرید که در سه بلوک شهری زندگی می‌کنند. اگر خانواده‌ها را واحدهای اولیه، بلوکها را به عنوان طبقه‌ها، و اندازه خانواده‌ها را مشخصه X در نظر بگیریم، وضعیتی خواهیم داشت که در جدول ۶.۵ نشان داده شده است.

تابلوی ۱.۵ پارامترهای طبقه‌ها و جامعه برای نمونه‌گیری طبقه‌بندی شده		
در داخل طبقه	تمام جامعه	
		مجموع
$X_{h+} = \sum_{i=1}^{N_h} X_{h,i}$	$X = \sum_{h=1}^L \sum_{i=1}^{N_h} X_{h,i} = \sum_{h=1}^L X_{h+}$	(۱.۵)
		میانگین
$\bar{X}_h = \frac{\sum_{i=1}^{N_h} X_{h,i}}{N_h} = \frac{X_{h+}}{N_h}$	$\bar{X} = \frac{\sum_{h=1}^L N_h \bar{X}_h}{N} = \sum_{h=1}^L W_h \bar{X}_h = \frac{X}{N}$	(۲.۵)
		نسبت
$P_{hy} = \frac{\sum_{i=1}^{N_h} Y_{h,i}}{N_h}$	$P_y = \frac{\sum_{h=1}^L N_h P_{hy}}{N} = \sum_{h=1}^L W_h P_{hy}$	(۳.۵)
		واریانس
$\sigma_{hx}^2 = \frac{\sum_{i=1}^{N_h} (X_{h,i} - \bar{X}_h)^2}{N_h}$		(۴.۵)
		واریانس نسبی
$V_{hx}^2 = \frac{\sigma_{hx}^2}{\bar{X}_h^2}$		(۵.۵)
<p>در این تعریفها L تعداد طبقه‌ها، N_h تعداد عناصر در طبقه h، N تعداد کل عناصر، $X_{h,i}$ مقدار X برای iامین عنصر در طبقه h، $Y_{h,i}$ نشانه بود یا نبود یک صفت کیفی دو حالتی Y برای iامین عنصر در طبقه h، و $W_h = \frac{N_h}{N}$ نسبتی از کل جامعه که متعلق به طبقه h است.</p>		

جدول ۶.۵ طبقه‌ها برای جامعه متشکل از ۱۴ خانواده

اندازه خانواده	خانواده	بلوک
۴	۱	۱
۳	۲	
۴	۳	
۴	۱	۲
۶	۲	
۴	۳	
۸	۴	
۸	۵	
۲	۱	۳
۳	۲	
۲	۳	
۲	۴	
۲	۵	
۳	۶	

بر حسب نمادهایی که در بالا معرفی شدند، مقادیر X برای هر طبقه به شرح زیرند:

طبقه ۳ ($N_3 = 6$)	طبقه ۲ ($N_2 = 5$)	طبقه ۱ ($N_1 = 3$)
$X_{3,1} = 2$	$X_{2,1} = 4$	$X_{1,1} = 4$
$X_{3,2} = 3$	$X_{2,2} = 6$	$X_{1,2} = 3$
$X_{3,3} = 2$	$X_{2,3} = 4$	$X_{1,3} = 4$
$X_{3,4} = 2$	$X_{2,4} = 7$	
$X_{3,5} = 2$	$X_{2,5} = 8$	
$X_{3,6} = 3$		

مجموع متغیر X در داخل هر طبقه h با استفاده از معادله (۱.۵) به شرح زیر است:

$$X_{1+} = 4 + 3 + 4 = 11$$

$$X_{2+} = 4 + 6 + 4 + 7 + 8 = 29$$

$$X_{3+} = 2 + 3 + 2 + 2 + 2 + 3 = 14$$

مجموع برای تمام جامعه دوباره با استفاده از معادله (۱.۵) عبارت است از:

$$X = 11 + 29 + 14 = 54$$

میانگینهای جامعه‌ای برای طبقات با استفاده از معادله (۲.۵) به صورت زیرند:

$$\bar{X}_1 = \frac{11}{3} = 3/67 \quad \bar{X}_2 = \frac{29}{5} = 5/8 \quad \bar{X}_3 = \frac{14}{6} = 2/33$$

میانگین جامعه کل از معادله (۲.۵) به صورت زیر به دست می‌آید:

$$\bar{X} = \frac{3}{14} \times 3/67 + \frac{5}{14} \times 5/8 + \frac{6}{14} \times 2/33 = 3/857$$

واریانسهای طبقات با استفاده از معادله (۴.۵) عبارت‌اند از:

$$\sigma_{1x}^2 = \frac{(4 - 3/67)^2 + (3 - 3/67)^2 + (4 - 3/67)^2}{3} = 0/222$$

$$\sigma_{2x}^2 = \frac{(4 - 5/8)^2 + (6 - 5/8)^2 + \dots + (8 - 5/8)^2}{5} = 2/56$$

$$\sigma_{3x}^2 = \frac{(2 - 2/33)^2 + (3 - 2/33)^2 + \dots + (3 - 2/33)^2}{6} = 0/222$$

واریانسهای نسبی برای طبقات از معادله (۵.۵) به شرح زیر به دست می‌آیند:

$$V_{1x}^2 = \frac{0/222}{(3/67)^2} = 0/165$$

$$V_{2x}^2 = \frac{2/56}{(5/8)^2} = 0/761$$

$$V_{3x}^2 = \frac{0/222}{(2/33)^2} = 0/409$$

□

۵.۵ آماره‌های نمونه برای طبقات

فرض کنید که در داخل یک طبقه خاص h به نحوی یک نمونه n_h عنصری از N_h عنصر موجود در طبقه انتخاب و هر عنصر نمونه را نسبت به نوعی متغیر X اندازه‌گیری کرده‌ایم. برای راحتی کار، عناصر

نمونه را از ۱ تا n_h شماره‌گذاری می‌کنیم و قرار می‌گذاریم که $x_{h,1}$ مقدار متغیر X برای عنصر شماره ۱ «۱» نمونه و $x_{h,2}$ مقدار متغیر X برای عنصر شماره «۲» نمونه و x_{h,n_h} مقدار X برای عنصر شماره « n_h » نمونه باشد. اگر نمونه‌ای متشکل از n_h عنصر در داخل هر طبقه h گرفته شود، در آن صورت کل اندازه نمونه n از راه زیر به دست می‌آید:

$$n = \sum_{h=1}^L n_h$$

که در آن L ، تعداد طبقاتی است که عناصر جامعه در آنها گروه‌بندی شده‌اند.

مثلاً، در مثال قبل اگر یک عنصر از طبقه ۱، دو عنصر از طبقه ۲ و چهار عنصر از طبقه ۳، انتخاب کنیم خواهیم داشت

$$n_1 = 1, \quad n_2 = 2, \quad n_3 = 4$$

و

$$n = 1 + 2 + 4 = 7$$

اگر مقادیر چهار عنصر انتخاب شده از طبقه ۳ عبارت باشند از $X_{3,1}, X_{3,2}, X_{3,3}, X_{3,4}$ و $X_{3,5}, X_{3,6}$ در آن صورت خواهیم داشت:

$$\begin{aligned} x_{3,1} = X_{3,1} = 2 & \quad x_{3,3} = X_{3,3} = 2 \\ x_{3,2} = X_{3,2} = 2 & \quad x_{3,4} = X_{3,4} = 3 \end{aligned}$$

مجموع نمونه‌ای و میانگین نمونه‌ای برای یک طبقه خاص از فرمول زیر به دست می‌آیند:

$$x_{h+} = \sum_{i=1}^{n_h} x_{h,i} \quad \text{و} \quad \bar{x}_h = \frac{\sum_{i=1}^{n_h} x_{h,i}}{n_h} = \frac{x_{h+}}{n_h}$$

۶.۵ برآورد پارامترهای جامعه از روی نمونه‌گیری تصادفی طبقه‌بندی شده

اگر انتخاب عناصر نمونه به طور مستقل در داخل هر طبقه با استفاده از نمونه‌گیری تصادفی ساده انجام پذیرفته باشد، در آن صورت این طرح را طرح نمونه‌گیری تصادفی طبقه‌بندی شده می‌نامند و برآوردهای مناسب پارامترهای جامعه همراه با خطاهای معیار آنها در تابلوی ۶.۵ نشان داده شده‌اند.

در این تعاریف، L تعداد طبقات، N_h تعداد عناصر در طبقه h ، N کل تعداد عناصر، $x_{h,i}$ مقدار X برای i امین عنصر در طبقه h و $y_{h,i}$ معرف بود یا نبود یک صفت کیفی دو حالتی Y برای i امین عنصر در طبقه h است.

همچنین S_{hx}^2 ، برآورد واریانس توزیع X برای طبقه h است که از فرمول زیر به دست می‌آید:

$$S_{hx}^2 = \frac{\sum_{i=1}^{n_h} (x_{h,i} - \bar{x}_h)^2}{n_h - 1} \quad (9.5)$$

تابلوی ۲.۵ برآوردهای پارامترهای جامعه و خطاهای معیار این برآوردها در نمونه‌گیری طبقه‌بندی شده

پارامتر	خطای معیار برآورد برای تمام جامعه	برآورد برای تمام جامعه	برآورد برای طبقه
مجموع			
	$\hat{SE}(x'_{str}) = \sqrt{\sum_{h=1}^L \frac{N_h^2 S_{hx}^2}{n_h} \left(\frac{N_h - n_h}{N_h} \right)}$	$x'_{str} = \sum_{h=1}^L x'_h$	$x'_h = N_h \bar{x}_h$
میانگین			
	$\hat{SE}(\bar{x}_{str}) = \sqrt{\sum_{h=1}^L \left(\frac{N_h}{N} \right)^2 \frac{S_{hx}^2}{n_h} \left(\frac{N_h - n_h}{N_h} \right)}$	$\bar{x}_{str} = \frac{\sum_{h=1}^L N_h \bar{x}_h}{N}$	$\bar{x}_h = \frac{\sum_{i=1}^{n_h} x_{h,i}}{n_h}$
نسبت			
	$\hat{SE}(p_{y, str}) = \sqrt{\sum_{h=1}^L \left(\frac{N_h}{N} \right)^2 \frac{p_{hy}(1-p_{hy})}{n_h - 1} \left(\frac{N_h - n_h}{N_h} \right)}$	$p_{y, str} = \frac{\sum_{h=1}^L y_{hi}}{N}$	$p_{hy} = \frac{\sum_{i=1}^{n_h} y_{h,i}}{n_h}$

مثال تشریحی: با استفاده از STATA یکصد و پنجاه و هشت (۱۵۸) بیمارستان در ناحیه‌ای بر حسب سطح خدمات زایمان در سه طبقه گروه‌بندی شده‌اند که با متغیر *oblevel* نشان داده می‌شود. از این ۱۵۸ بیمارستان ناحیه، ۴۲ بیمارستان در طبقه ۱ (تخصصی)، ۹۹ بیمارستان در طبقه ۲ (متوسط) و ۱۷ بیمارستان در طبقه ۳ (درجه سه) قرار دارند. یک نمونه تصادفی طبقه‌بندی شده متشکل از ۱۵ بیمارستان از این جامعه گرفته شده است: ۴ بیمارستان در طبقه ۱؛ ۵ بیمارستان در طبقه ۲؛ و ۶ بیمارستان در طبقه ۳. این نمونه خاص متشکل از ۱۵ بیمارستان همراه با متغیرهای مربوط برای برآورد، در پرونده اطلاعاتی STATA به صورت *hospsamp.dta* ثبت شده‌اند. این پرونده در صفحه بعد نشان داده شده است.

Hospno	oblevel	weighta	tothosp	births
15	1	10.50	42	480
80	1	10.50	42	426
86	1	10.50	42	342
136	1	10.50	42	174
7	2	19.80	99	2022
26	2	19.80	99	576
62	2	19.80	99	1999
90	2	19.80	99	482
101	2	19.80	99	836
28	3	2.83	17	3108
34	3	2.83	17	4674
39	3	2.83	17	2539
102	3	2.83	17	1610
119	3	2.83	17	4618
149	3	2.83	17	1781

معانی متغیرها به شرح زیرند:

hospno = کد شناسایی برای بیمارستان خاص.

oblevel = سطح خدمات زایمان برای هر بیمارستان.

weighta = وزن نمونه‌گیری $\frac{N_h}{n_h}$ برای بیمارستان خاص.

tothosp = تعداد کل جامعه بیمارستانها در طبقه‌ای که بیمارستان نمونه قرار دارد.

births = تعداد رویدادهای تولد در بیمارستان در طی سال قبل.

از فرمانهای زیر می‌توان برای به دست آوردن برآوردهای کل تعداد تولدها و خطاهای معیار

برآوردها استفاده کرد.

```
. use "a:\hospsamp.dta", clear
. svyset pweight weighta
. svyset strata oblevel
. svyset fpc tothosp
. svytotal births
. svytotal births, by(oblevel)
```

فرمان اول نشان می‌دهد که داده‌ها در فایل اطلاعاتی STATA به نام *hospsamp.dta* قرار دارند.

فرمان دوم نشان می‌دهد که وزن نمونه‌گیری برای بیمارستان در متغیر *weighta* قرار دارد.

فرمان سوم نشان می‌دهد که متغیر نشانگر طبقه در متغیر *oblevel* قرار دارد.

فرمان چهارم نشان می‌دهد که N_h ، تعداد بیمارستانهای طبقه در متغیر *tothosp* قرار دارد. در اینجا نیز این فرمان برای گزارشهای متعلق به یک طبقه یکسان است ولی برای گزارشهای مربوط به بیمارستانها در طبقه‌های مختلف یکسان نیست. تصحیح جامعه متناهی از روی این متغیر محاسبه می‌شود. فرمان پنجم نشان می‌دهد که کل تعداد تولدها باید از روی داده‌های نمونه برآورد شود. فرمان ششم نشان می‌دهد که کل تعداد تولدها برای هر طبقه باید از این داده‌ها برآورد شود. خروجی STATA که از فرمانهای بالا تولید شده به شرح زیر است:

Survey total estimation					
pweight:	weighta		Number of obs	=	15
Strata:	oblevel		Number of strata	=	3
PSU:	<observations>		Number of PSUs	=	15
FPC:	tothosp		Population size	=	158
Total	Estimate	Std. Err.	[95% Conf. Interval]		Deff
births	183983	34014.35	109872.1	258093.9	.7035476
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.					
Survey total estimation					
pweight:	weighta		Number of obs	=	15
Strata:	oblevel		Number of strata	=	3
PSU:	<observations>		Number of PSUs	=	15
FPC:	tothosp		Population size	=	158
Total Subpop	Estimate	Std. Err.	[95% Conf. Interval]		Deff
births					
oblevel == 1	14931	2669.857	9113.882	20748.12	.1564799
oblevel == 2	117117	33067.68	45068.7	189165.3	1.089406
oblevel == 3	51935	7508.403	35575.6	68294.4	.0330073
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.					

مثال تشریحی با استفاده از SUDAAN: دوباره مثالی را که در آن ۱۵۸ بیمارستان در یک ناحیه بر حسب سطح خدمات زایمان در سه طبقه گروه‌بندی شده‌اند در نظر می‌گیریم، که در آن یک نمونه تصادفی طبقه‌بندی شده گرفتیم و ۴ بیمارستان از طبقه ۱، ۵ بیمارستان از طبقه ۲ و ۶ بیمارستان از طبقه ۳ انتخاب کردیم. این نمونه خاص متشکل از ۱۵ بیمارستان در بحث استفاده از STATA برای تحلیل این داده‌ها در بالا نشان داده شد. حالا فرض می‌کنیم که پرونده اطلاعاتی *hospsamp.dta* یک

پرونده اطلاعاتی SAS PC است (*hospsamp.ssd*). در زیر نشان خواهیم داد که چگونه نرم‌افزار آماری SUDAAN می‌تواند برآوردها و خطاهای معیار را برای این طرح نمونه‌گیری تصادفی طبقه‌بندی شده تولید کند.

دوباره با فرض این که *hospsamp* به پرونده اطلاعاتی SAS ارسال شده است، می‌توانیم از فرمانهای زیر در دست آوردن برآوردهای کل تولدها برای هر طبقه و برای ۱۵۸ بیمارستانی که جامعه را تشکیل می‌دهند، استفاده کنیم.

```
PROC DESCRIPT DATA = HOSPSAMP FILETYPE = SAS DESIGN = WOR TOTALS;
  NEST OBLEVEL;
  WEIGHT WEIGHTA;
  TOTCNT TOTHOSP;
  VAR BIRTHS;
  SUBGROUP OBLEVEL;
  LEVELS 3;
  SETENV COLWIDTH=20;
  SETENV DESWIDTH=3;
```

سطر اول اطلاعاتی درباره پرونده اطلاعاتی ارائه می‌دهد و حاکی از آن است که نمونه‌گیری بدون جایگذاری است و قرار است برآورد مجموعها را به دست آوریم.

سطر دوم نشان می‌دهد که گزارشهای نمونه در داخل متغیر طبقه *oblevel* جای داده شده‌اند.

سطر سوم متغیری را که قرار است به عنوان وزن نمونه‌گیری به کار رود مشخص می‌سازد که در این مورد *weighta* نامیده می‌شود و وارون کسر نمونه‌گیری خاص طبقه است.

سطر چهارم نشان می‌دهد که متغیر *tothosp* شامل کل تعداد بیمارستانهای جامعه آماری است، (۴۲ بیمارستان برای طبقه ۱، ۹۹ بیمارستان برای طبقه ۲ و ۱۷ بیمارستان برای طبقه ۳).

سطر پنجم نشان می‌دهد که کل تعداد تولدها قرار است برآورد شود.

سطرهای ششم و هفتم نشان می‌دهند که کل تعداد تولدها برای هر طبقه قرار است برآورد شود.

سطرهای ۸ و ۹ فورمت خروجی را تعیین می‌کنند.

خروجی زیر از فرمانهای فوق‌الذکر SUDAAN تولید شده است.

برآوردهای حاصل از SUDAAN و STATA با یکدیگر مطابقت دارند و خواننده می‌تواند تحقیق

کند که با برآوردهای حاصل از فرمولهای تابلوی ۲.۵ نیز مطابقت دارند.

□

۷.۵ خلاصه

در این فصل مروری بر طبقه‌بندی داشتیم. مفهوم طبقه‌بندی را معرفی و تعریف کردیم که منظور از نمونه تصادفی طبقه‌بندی شده چیست. مواردی را به بحث گذاشتیم که ممکن بود طبقه‌بندی در آنها

مناسب باشد و مزایای طبقه‌بندی را نسبت به نمونه‌گیری تصادفی ساده و نیز معایب عمده آن را توضیح دادیم. عبارتهایی را که برای مشخص کردن پارامترهای جامعه به کار می‌روند و آماره‌های تحت نمونه‌گیری تصادفی ساده را ارائه کردیم. سرانجام نشان دادیم که چگونه می‌توان برآوردهایی از طریق نمونه‌گیری تصادفی طبقه‌بندی شده، همراه با خطاهای معیار آنها، با استفاده از دو نرم‌افزار STATA و SUDAAN به دست آورد.

```

1 PROC DESCRIPT DATA = HOSPSAMP FILETYPE = SAS DESIGN = WOR
  TOTALS;
2 NEST OBLEVEL;
3 WEIGHT WEIGHTA;
4 TOTCNT TOTHOSP;
5 VAR BIRTHS;
6 SUBGROUP OBLEVEL;
7 LEVELS 3;
8 SETENV COLWIDTH=20;
9 SETENV DECWIDTH=3;

```

Number of observations read : 15 Weighted count: 158

Number of observations skipped : 0

(WEIGHT variable nonpositive)

denominator degrees of freedom : 12

by: Variable, OBLEVEL.

Variable		OBLEVEL	
		Total	1
BIRTHS	Sample Size	15.000	4.000
	Weighted Size	158.000	42.000
	Total	183982.904	14931.000
	SE Total	34014.329	2669.857
	Mean	1164.449	355.500
	SE Mean	215.281	63.568

by: Variable. OBLEVEL.

Variable		OBLEVEL	
		2	3
BIRTHS	Sample Size	5.000	6.000
	Weighted Size	99.000	17.000
	Total	117116.928	51934.977
	SE Total	33067.664	7508.399
	Mean	1183.000	3055.000
	SE Mean	334.017	441.671

تمرین

۱.۵ در یک نمونه تصادفی ساده در مورد ۱۰۰۰۰ دونده از میان ۱۰۰۰۰ دونده‌ای که در مسابقهٔ ماراتون ۱۹۹۵ شیکاگو شرکت کرده بودند، نتیجهٔ آزمایش خون ۳۵ نفر برای مصرف استروئید و سایر داروهایی که فعالیت را زیاد می‌کنند مثبت بود. وقتی نتایج بر حسب زمان تکمیل مسابقه دسته‌بندی شدند جدول زیر به دست آمد.

زمان تکمیل مسابقه (ساعت)	تعداد در نمونه	تعداد نتیجهٔ مثبت برای مصرف دارو	درصد مثبت
کمتر از ۲/۵	۱۰۰	۲۵	۲۵/۱۰۰
۲/۵-۴	۵۰۰	۷	۱/۴۰
بیشتر از ۴	۴۰۰	۳	۰/۷۵
جمع	۱۰۰۰	۳۵	۳/۵۰

الف. خطای معیار نسبی که برای مصرف داروی فعالیت‌زا مثبت برآورد شده چقدر است؟
ب. با بررسی (بدون محاسبه) نرخهای ارائه شده در بالا فکر می‌کنید طبقه‌بندی می‌توانست منجر به برآورد اساساً بهتری شود؟ چرا آری یا چرا نه؟

۲.۵ فرض کنید که در تمرین ۱.۵ از ۱۰۰۰۰ نفری که مسابقهٔ ماراتون را به پایان رساندند ۲۰۰۰ نفر در زمانی کمتر از ۲/۵ ساعت، ۶۰۰۰ نفر بین ۲/۵ تا ۴ ساعت و ۲۰۰۰ نفر در زمانی بیشتر از ۴ ساعت مسابقه را تکمیل کرده باشند. در مورد نتایج نمونه‌گیری به صورتی که در تمرین ۱.۵ ارائه شد اظهار نظر کنید.

۳.۵ اگر یک نمونهٔ تصادفی طبقه‌بندی شده متشکل از ۳۳۳ نفر در هر یک از گروههای سه‌گانه گرفته می‌شد، برآورد درصد مثبت که از این طریق به دست می‌آمد احتمالاً چقدر می‌شد؟

۴.۵ تعداد ۲۰۰۰ دوندهٔ دیگر نیز وارد بازی شدند ولی آن را به پایان نرساندند. از میان این عده به تصادف برای ۶۰۰ نفر پرسشنامه‌هایی از طریق پست ارسال شد که ۵۰۰ نفر آن را پر کردند. یکی از اقلام این پرسشنامه از پاسخگو می‌خواست که متوسط تعداد مایلهایی را که در هفته طی مدت ۸ هفته قبل از مسابقهٔ ماراتون پیموده است برآورد کند. میانگین، \bar{x} ، تعداد مایلهای طی شده $32/4$ با انحراف معیار s_x معادل $7/3$ بود. یک نمونهٔ تصادفی ساده متشکل از ۵۰۰ نفر از ۱۰۰۰۰ دونده‌ای که مسابقهٔ ماراتون را تکمیل کرده بودند ۴۰۰ پاسخگو داشت که

متوسط تعداد مایلهای هفتگی آنها $46/8$ با انحراف معیاری معادل $6/2$ مایل بود. متوسط مایلهای طی شده در هفته طی مدت مورد نظر برای همه کسانی که وارد مسابقهٔ ماراتون شدند چقدر است؟

۵.۵ داده‌های زیر مربوط به ۶ سازمان حفظ بهداشت برای سال ۱۹۸۸ در یک شهر با اندازه متوسط است:

تعداد کارکنان تأمین‌کنندهٔ مراقبت از بیماران و تعداد رویارویی با بیمار طی سال ۱۹۸۸ توسط ۶ سازمان

تعداد رویارویی با بیمار	تعداد پزشکان تأمین‌کنندهٔ مراقبت از بیماران	HMO سازمان حفظ بهداشت
۲۲۰۰۰	۱۰	۱
۱۴۰۰۰	۶	۲
۱۰۲۰۰	۴	۳
۷۰۰۰۰	۳۰	۴
۱۵۰۰۰	۷	۵
۵۰۰۰	۳	۶

چون این طرح مستلزم استفادهٔ نوبتی از هزینه و زمان است پیشنهاد می‌شود که یک نمونهٔ دوتایی از این سازمانهای حفظ بهداشت برای برآورد کل تعداد رویارویی با بیماران در این شش سازمان طی سال ۱۹۸۸ گرفته شود.

الف. همهٔ نمونه‌های تصادفی سادهٔ دوتایی را شمارش و کل تعداد رویارویی با بیمار در سال ۱۹۸۸ را برای ۶ سازمان برآورد کنید.

ب. خطای معیار برآورد مجموع از روی طرح نمونهٔ تعیین شده در بالا چقدر است؟

پ. سازمانهای شمارهٔ ۱، ۲، ۳، ۵ و ۶ را در یک طبقه و سازمان شمارهٔ ۴ را (به تنهایی) در یک طبقهٔ دیگر گروه‌بندی کنید. از روی این دو طبقه، تمام نمونه‌های تصادفی طبقه‌بندی شدهٔ دوتایی را شمارش کنید. خطای معیار برآورد مجموع حاصل از این برنامهٔ نمونه‌گیری چقدر است؟

ت. در مورد مطلوبیت استفاده از طبقه‌بندی در این وضعیت در مقابل نمونه‌گیری تصادفی ساده اظهار نظر کنید.

۶.۵ فرض کنید که در وضعیت ارائه شده در تمرین ۵.۵، برآورد تعداد رویارویی در سال ۱۹۸۸ به ازای هر پزشک مورد نظر است. قسمت‌های (الف) تا (ت) تمرین ۵.۵ را برای این وضعیت برآورد تکرار کنید. آیا در این وضعیت نمونه‌گیری تصادفی طبقه‌بندی شده بهتر از نمونه‌گیری تصادفی ساده است؟ چرا بهتر است و چرا نیست؟

۷.۵ در درمانگاه بزرگی واقع در یک بیمارستان مرکز شهر، ۵۶ بیماری که دارای عفونت بدون علامت‌های ظاهری ایدز هستند با یک داروی آزمایشی تحت درمان قرار گرفته‌اند که تصور می‌شود توانایی بازدهی برخی از کارکردهای سیستم دفاعی بدن را که با عفونت HIV همراه است دارند. از این بیماران شمارش گلبولی CD4 برای ۱۲ نفر در ویزیت اول کمتر از ۲۵۰ بود، برای ۲۰ نفر بین ۲۵۰ تا ۴۰۰ بود و برای ۲۴ نفر بیشتر از ۴۰۰ بود. (هر چه این شمارش کمتر باشد وضعیت بیمار وخیمتر است). می‌خواهیم یک نمونه تصادفی طبقه‌بندی شده متشکل از ۳۰ بیمار یعنی ۱۰ نفر از هر یک از گروه‌های سه‌گانه بالا بگیریم با این هدف که وقوع پیشامدهای ناشی از ایدز را در یک دوره ۱۲ ماهه در میان این بیماران برآورد کنیم.

الف. چند نمونه متشکل از ۳۰ بیمار به شرحی که در بالا گذشت امکان پذیر است؟
ب. با فرض این که بیمار می‌تواند بیش از یک پیشامد حاکی از ایدز طی این مدت داشته باشد به صورت جبری نشان دهید که چگونه وقوع پیشامدهای حاکی از ایدز را از روی نمونه برای دوره ۱۲ ماهه برآورد می‌کنید.

پ. چند نمونه ۳۰ تایی تحت نمونه‌گیری تصادفی ساده بدون طبقه‌بندی امکان پذیر است؟

۸.۵ در زیر نتایج آمارگیری نمونه‌ای که در تمرین ۷.۵ توصیف شد ارائه شده است:

تعداد افراد

تعداد پیشامدها	$CD4 < 250$	$250 \leq CD4 < 400$	$CD4 \geq 400$
۰	۵	۷	۹
۱	۱	۲	۱
۲	۲	۱	۰
۳	۲	۰	۰

از روی این داده‌ها وقوع پیشامدهای حاکی از ایدز را در جامعه هدف برآورد کنید. خطای معیار برآورد نرخ وقوع چقدر است؟ (باید فرمول خطای معیار را از «اصول اولیه» به دست آورید).

۹.۵ از روی داده‌های تمرین ۸.۵ نسبت بیماران دارای یک یا چند پیشامد حاکی از آیدز را برآورد کنید. خطای معیار این نسبت برآورد شده چقدر است؟

۱۰.۵ یک نمونه تصادفی ساده (بدون طبقه‌بندی) متشکل از ۳۰ نفر از این ۵۶ بیمار، داده‌های زیر را به دست داده است:

تعداد افراد			تعداد پیشامدها
$CD4 \geq 400$	$250 \leq CD4 < 400$	$CD4 < 250$	
۱۲	۸	۳	۰
۰	۲	۱	۱
۱	۱	۰	۲
۰	۰	۲	۳

وقوع پیشامدهای حاکی از آیدز را از روی این داده‌ها برآورد کنید و خطای معیار این برآورد را به دست آورید. این نتایج را با نتایج حاصل از طرح نمونه‌گیری طبقه‌بندی شده مقایسه کنید. کدام طرح، برآورد قابل اعتمادتری را تولید می‌کند؟ چرا؟

کتابشناسی

The following articles present sample surveys in which stratification is used.

1. Hemphill, F. M., A sample survey of home injuries, *Public Health Reports*, 67: 1026, 1952.
2. Horvitz, D. G., Sampling and field procedures of the Pittsburgh morbidity survey, *Public Health Reports*, 67: 1003, 1952.
3. Goldberg, J., Levy, P. S., Mullner, R., Gelfand, H., Iverson, N., Lemeshow, S., and Rothrock, J., Factors affecting trauma center utilization in Illinois, *Medical Care* 19: 547, 1981.
4. Stasny, E. A., Toomey, B. G., and First, R. J., Estimating the rate of rural homelessness: A study of nonurban Ohio, *Survey Methodology*, 20: 87, 1994.
5. Barner, B. M., and Levy, P. S., State-wide shoulder belt usage by type of roadway and posted speed limit: A three year comparison, *38th Annual Proceedings Association for the Advancement of Automotive Medicine*, Association for the Advancement of Automotive Medicine, Des Plaines Ill., 1994.

The following software programs have the capability of taking stratified samples.

6. Frankel, M. R., and Spencer, B. D., *SAMPLE. A Supplementary Module for SYSTAT*, SYSTAT, Inc., Evanston, Ill., 1990.
7. SAS Institute Inc., *SAS Language and Procedures*, Usage 2, Version 6, 1st ed., Cary, N.C., SAS Institute Inc., 1991, pp. 649.

In addition, there is material on stratification in virtually every text on sampling theory and survey methodology, including those listed in the bibliography sections of earlier chapters.