

نمونه‌گیری با احتمال متغیر

۰.۳ مقدمه

آنچه در فصل دوم تحت عنوان نمونه‌گیری تصادفی ساده شرح دادیم مبتنی بر این بود که احتمال انتخاب همه نمونه‌های ممکن به حجم n از جامعه به حجم N ، و در نتیجه احتمال انتخاب همه واحدهای جامعه برای شرکت در یک نمونه یکسان باشد. در عمل، این یکسان بودن احتمال انتخاب واحدها، گاهی به حصول برآوردهایی منجر می‌شود که با وجود نااریب بودن، گاهی به دلیل بزرگ بودن واریانس جامعه، از واقعیت دورند. برای روشن شدن مطلب به ذکر مثالی می‌پردازیم. فرض کنید بخواهیم جمعیت شهرنشین استانی را برآورد کنیم. اگر استان شامل ده شهر باشد که جمعیت آنها برحسب ۱۰۰۰ نفر به ترتیب ۱۰۰، ۱۱۰، ۱۵۰، ۱۳۰، ۲۱۰، ۷۰، ۵۰۰، ۵۴۰، ۷۰۰ و ۸۰ است، کل جمعیت شهرنشین استان برابر با ۲۵۹۰ است. اگر بخواهند با نمونه‌ای تصادفی به حجم ۴، جمعیت کل را برآورد کنند ممکن است نمونه حاصل ۱۰۰، ۸۰، ۱۳۰، ۷۰ باشد که چون میانگین نمونه برابر با ۹۵ است جمعیت کل برابر با ۹۵۰ برآورد می‌شود. می‌بینید که این برآورد، با اینکه نااریب است به دلیل بزرگی واریانس از واقعیت، یعنی از ۲۵۹۰ بسیار دور است. علت این اختلاف بارز این است که در روش نمونه‌گیری تصادفی، شانس انتخاب شهری با جمعیت ۵۰۰ برای شرکت در نمونه، با شانس انتخاب شهری با جمعیت مثلاً ۷۰ برای شرکت در نمونه یکی است که ممکن است به انتخاب نمونه‌ای مثل نمونه بالا منجر شود که برآوردی غیرمنطقی را به دست دهد. عقل سلیم می‌گوید که در انتخاب واحدهای نمونه، به شهری که جمعیت بیشتری

دارد باید شانس انتخاب بیشتری داده شود. دادن احتمال به واحدها برای شرکت در نمونه، در واقع دادن وزنهایی است که عمدتاً واریانس جامعه را تقلیل می‌دهد. همین ایده، نمونه‌گیری جدیدی را به نام نمونه‌گیری با احتمال متغیر القا می‌کند.

۱.۳ تعریف نمونه‌گیری با احتمال متغیر

یک جامعه با N واحد را در نظر می‌گیریم و فرض می‌کنیم Y_i معرف مقدار صفت تحت مطالعه برای i امین واحد جامعه ($i = 1, 2, \dots, N$) باشد. به علاوه فرض می‌کنیم که p_i ($i = 1, 2, \dots, N$) معرف احتمال انتخاب i امین واحد جامعه برای عضویت در نمونه‌ای مفروض باشد. بدیهی است $\sum_{i=1}^N p_i = 1$. در این نمونه‌گیری احتمال انتخاب واحدهای جامعه برای شرکت در نمونه ثابت نیست و معمولاً از واحدی به واحد دیگر تغییر می‌کند. اگر در حالت خاص، p_i ها متناسب با اندازه Y_i ها باشند نمونه‌گیری را تصادفی با احتمال متناسب با اندازه می‌نامند و آن را نمونه‌گیری PPS* می‌گویند.

نمونه‌گیری تصادفی با احتمال متغیر را می‌توان به دو روش با جایگذاری و بدون جایگذاری انجام داد. در روش با جایگذاری احتمال انتخاب واحد Y_i برای شرکت در نمونه به همان مقداری است که از قبل به صورت p_i مشخص شده است، زیرا در هر انتخاب با کل جامعه سروکار داریم، اما در نمونه‌گیری بدون جایگذاری با وضعی کمی پیچیده روبه‌رو هستیم. بدیهی است که نحوه انتخاب واحدها با روش انتخاب نمونه تصادفی ساده متفاوت است. قبل از ورود به بحث برآورد پارامترهای جامعه با استفاده از نمونه تصادفی با احتمال متغیر، شیوه‌هایی را برای انتخاب نمونه با احتمال متغیر که خاص روش با جایگذاری است شرح می‌دهیم.

۲.۳ انتخاب نمونه با احتمال متغیر و با جایگذاری به شیوه مجموع

تراکمی

در نمونه‌گیری تصادفی ساده برای انتخاب واحدهای نمونه، از جدول اعداد تصادفی استفاده می‌کنیم، زیرا احتمال انتخاب همه واحدها برای شرکت در نمونه یکسان است ولی در نمونه‌گیری با احتمال متغیر، چون احتمال تخصیص یافته به هر واحد در حالت کلی با احتمالهای منسوب به واحدهای دیگر برای عضویت در نمونه فرق می‌کند استفاده مستقیم از جدول اعداد تصادفی میسر نیست. برای انتخاب نمونه اعداد صحیح X_1, X_2, \dots, X_N را که متناسب با احتمالهای انتخاب واحدهای Y_1, Y_2, \dots, Y_N جامعه هستند در نظر می‌گیریم. این اعداد صحیح وقتی p_i ها اعدادی گویا باشند همیشه وجود دارند. مثلاً اگر $N = 5$ و احتمالهای متناظر با ۵ واحد جامعه به ترتیب $\frac{1}{6}, \frac{1}{4}, \frac{1}{5}, \frac{1}{3}, \frac{1}{6}$ باشند، وقتی احتمالها را هم مخرج کنیم هم‌ارز با $\frac{1}{6}, \frac{2}{6}, \frac{2}{6}, \frac{1}{6}, \frac{1}{6}$ و $\frac{15}{6}$ و $\frac{2}{6}$

خواهند بود. لذا این احتمالات با اعداد صحیح ۱۰، ۲۰، ۱۲، ۱۵، و ۳ متناسبند و برای این مثال

$$X_1 = 10, X_2 = 20, X_3 = 12, X_4 = 15, X_5 = 3$$

اعداد صحیح متناظر با Y_1, Y_2, \dots, Y_5 هستند. ما حالتی را که احتمالات گنگاند نادیده می‌گیریم. اصطلاحاً دنباله X_1, X_2, \dots, X_N را دنباله صفتهای کمکی متناظر با Y_1, Y_2, \dots, Y_N می‌نامیم، لذا

$$\begin{array}{l} \text{صفت اصلی} : Y_1 \quad Y_2 \quad \dots \quad Y_i \quad \dots \quad Y_N \\ \text{صفت کمکی} : X_1 \quad X_2 \quad \dots \quad X_i \quad \dots \quad X_N \end{array}$$

حال برای انتخاب نمونه با احتمال متغیر، به اولین واحد جامعه صفت اصلی، عدد $T_1 = X_1$ و به دومین واحد جامعه عدد $T_2 = X_1 + X_2$ و در حالت کلی به i امین واحد جامعه عدد $T_i = X_1 + X_2 + \dots + X_i$ ($i = 1, \dots, N$) را نسبت می‌دهیم. بدیهی است عدد منسوب به Y_N برابر $T_N = \sum_{i=1}^N X_i$ است. در این صورت جدول زیر را به دست می‌آوریم.

جدول ۱.۳ نمایش احتمالات انتخاب متناسب با اندازه واحدها

| صفت اصلی | صفت کمکی | T_i | p_i |
|----------|----------|---------------------------|-----------|
| Y_1 | X_1 | X_1 | X_1/T_N |
| Y_2 | X_2 | $X_1 + X_2$ | X_2/T_N |
| \vdots | \vdots | \vdots | \vdots |
| Y_i | X_i | $X_1 + X_2 + \dots + X_i$ | X_i/T_N |
| \vdots | \vdots | \vdots | \vdots |
| Y_N | X_N | $X_1 + X_2 + \dots + X_N$ | X_N/T_N |

با رجوع به جدول اعداد تصادفی، عددی از ۱ تا T_N انتخاب می‌کنیم. این عدد را R می‌نامیم. عدد R را با T_i ها در ستون سوم، مقایسه می‌کنیم. عدد R یا با یکی از T_i ها برابر است و یا بین دو T_i متوالی، یعنی مثلاً بین T_{i-1} و T_i ، قرار دارد. پس $T_{i-1} < R \leq T_i$. در این صورت واحد i ام با مقدار Y_i را به عنوان واحد نمونه انتخاب می‌کنیم. باید ثابت کنیم که انتخاب واحد i ام متناسب با X_i صورت گرفته است. تعداد حالت‌های مساعد برای اینکه R در رابطه $T_{i-1} < R \leq T_i$ قرار گیرد برابر است با $T_i - T_{i-1}$ ، یعنی

$$T_i - T_{i-1} = (X_1 + X_2 + \dots + X_i) - (X_1 + X_2 + \dots + X_{i-1}) = X_i$$

و چون R را از بین اعداد ۱ تا T_N انتخاب کرده‌ایم، لذا تعداد حالت‌های ممکن انتخاب R برابر با T_N است. پس

$$\begin{aligned} P(X_{i-1} < R \leq X_i) &= P(\text{انتخاب } Y_i \text{ برای شرکت در نمونه}) \\ &= \frac{X_i}{T_N} = \frac{X_i}{X_1 + \dots + X_N} = p_i \end{aligned}$$

یعنی واحد i ام که با شیوهٔ بالا انتخاب می‌کنیم و واحدی از نمونه است با احتمال p_i برگزیده می‌شود. پس از انتخاب واحد i ام، آن را به جامعه برمی‌گردانیم و این فرایند را n بار تکرار می‌کنیم تا از N واحد جامعه، نمونه‌ای به حجم n انتخاب شود. بدیهی است که امکان دارد برخی از واحدهای نمونهٔ منتخب، تکراری باشند.

مثال ۱.۳ صفتی که تحت بررسی است میزان محصول گندم ۱° روستاست. می‌خواهیم به‌کمک نمونه‌ای به حجم ۴، میزان محصول ۱° روستا را برآورد کنیم. مساحت زمین زیرکشت هر یک از ده روستا را می‌دانیم. در این مثال، چهار واحد نمونه را به روش تصادفی با احتمال متغیر، که احتمالها متناسب با صفت کمکی مساحت‌های زیرکشت‌اند، مشخص می‌کنیم.

میزان محصول ۱° روستا را به‌ترتیب Y_1, Y_2, \dots, Y_{10} می‌نامیم. این مقادیر در دست نیستند. مساحت زمین زیرکشت متناظر با Y_i ها را می‌دانیم و آنها را با X_1, X_2, \dots, X_{10} نشان می‌دهیم. مقادیر X_i در جدول زیر آمده‌اند. با توجه به این جدول، به‌شیوهٔ مجموع تراکمی ۴ واحد نمونه را تعیین می‌کنیم. در ستون سوم جدول، مقادیر T_i را مطابق شرحی که در بیان شیوهٔ انتخاب نمونه ذکر شد، مشخص کرده‌ایم. حال به جدول اعداد تصادفی رجوع می‌کنیم و از اعداد ۱ تا ۱۸۵° چهار عدد انتخاب می‌کنیم. فرض کنید اعداد $۷۲, ۵۵۱, ۱۷۳^\circ$ و ۹۸° به‌دست آیند. به‌ترتیب داریم

$$۷۲ < T_1$$

$$T_2 < ۵۵۱ < T_3$$

$$۱۷۳^\circ = T_8$$

$$T_9 < ۹۸^\circ < T_{10}$$

| T_i | مساحت زمین زیرکشت | محصول | شمارهٔ واحد |
|-------|-------------------|----------|-------------|
| ۱۰۰ | ۱۰۰ | Y_1 | ۱ |
| ۱۹۰ | ۹۰ | Y_2 | ۲ |
| ۴۹۰ | ۳۰۰ | Y_3 | ۳ |
| ۶۰۰ | ۱۱۰ | Y_4 | ۴ |
| ۱۰۰۰ | ۴۰۰ | Y_5 | ۵ |
| ۱۱۵۰ | ۱۵۰ | Y_6 | ۶ |
| ۱۶۵۰ | ۵۰۰ | Y_7 | ۷ |
| ۱۷۳۰ | ۸۰ | Y_8 | ۸ |
| ۱۸۰۰ | ۷۰ | Y_9 | ۹ |
| ۱۸۵۰ | ۵۰ | Y_{10} | ۱۰ |

لذا با توجه به استدلالی که در شیوهٔ انتخاب واحدهای نمونه ارائه شد. نمونهٔ تصادفی مورد نظر

عبارت است از واحدهای ۱، ۴، ۸، و ۵ جامعه با مقادیر

$$Y_1, Y_4, Y_8, Y_5$$

پس باید به روستاهای شماره ۱، ۴، ۸، و ۵ مراجعه و میزان محصول گندم آنها را ثبت کنیم. ممکن بود چهار عدد تصادفی منتخب، مثلاً اعداد ۱۵۰، ۱۷۱، ۱۵۰۱، و ۱۸۰۲ باشند. در این صورت

$$T_1 < 150 < T_2, \quad T_1 < 171 < T_2$$

$$T_6 < 1501 < T_7, \quad T_9 < 1802 < T_{10}$$

و لذا، نمونه منتخب به صورت زیر است: واحدهای ۲، ۲، ۷، و ۱۰ با مقادیر

$$Y_2, Y_2, Y_7, Y_{10}$$

می بینید که در این حالت واحد Y_2 یا روستای شماره ۲، دوبار در نمونه شرکت می کند. Δ

شیوه مجموع تراکمی، مستلزم تشکیل ستون T_i ها برای همه واحدهای مقادیر X_i هاست که مسلماً وقتی جامعه بزرگ است انجام آن حتی با نرم افزارهای کامپیوتری هم پرهزینه و وقتگیر است. روش زیر را که لاهیری برای انتخاب نمونه با احتمال متغیر پیشنهاد کرده است و روشی با هزینه کمتر است، عرضه می کنیم.

۳.۳ انتخاب نمونه با احتمال متغیر به شیوه لاهیری

به ستون اندازه های X_i مراجعه می کنیم و بزرگترین مقدار X_i را مشخص و آن را M می نامیم. اگر حجم جامعه باشد، یک زوج عدد تصادفی (i, j) را با شرایط $1 \leq i \leq N$ و $1 \leq j \leq M$ انتخاب می کنیم. وقتی i انتخاب شد X_i را در نظر می گیریم، در این صورت اگر داشته باشیم $X_i \leq j$ ، واحد i ام را به عنوان واحدی از نمونه مطلوب می کنیم. اگر داشته باشیم $X_i > j$ زوج انتخاب شده (i, j) را نادیده می گیریم و زوج تصادفی دیگری را انتخاب و فرایند قبلی را تکرار می کنیم. این فرایند را تا انتخاب n واحد نمونه ادامه می دهیم. قبل از اثبات اینکه استفاده از این شیوه، به هر واحد نمونه، احتمال انتخابی متناسب با X_i را تخصیص می دهد، به ذکر مثالی می پردازیم.

مثال ۲.۳ به مثال ۱.۳ برمی گردیم و این بار به روش لاهیری نمونه ای به حجم ۴ با احتمال متناسب با X_i ها به دست می آوریم. چون بزرگترین X_i برابر با ۵۰۰ است، پس $M = 500$ و از طرفی $N = 10$. زوج (i, j) از اعداد تصادفی را به قسمی انتخاب می کنیم که $1 \leq i \leq 10$ و $1 \leq j \leq 500$ ، یعنی عدد اول را با استفاده از یک ستون جدول اعداد تصادفی و عدد دوم را با استفاده از سه ستون انتخاب می کنیم. فرض می کنیم زوج $(5, 201)$ به دست آید. به ازای $i = 5$

با مراجعه به جدول داده‌ها، مشاهده می‌شود که $X_5 = 400$ و چون $X_5 < 201 = z$ ، لذا واحد $i = 5$ با مقدار Y_5 اولین واحد نمونه است. مجدداً زوجی تصادفی با شرایط مورد نظر انتخاب می‌کنیم. فرض می‌کنیم زوج $(6, 403)$ نتیجه شود، چون

$$z = 203 > X_6 = 150$$

این زوج تصادفی را نادیده می‌گیریم و زوج جدیدی از اعداد تصادفی را برمی‌گزینیم. ممکن است به دفعاتی متعدد موفق به انتخاب واحد جدیدی برای نمونه نشویم. اگر به‌عنوان مثال زوجهای تصادفی بعدی به ترتیب $(2, 83)$ ، $(8, 302)$ ، $(5, 97)$ ، و $(3, 25)$ باشند به‌سهولت می‌توانیم تحقیق کنیم که با زوج اول، Y_2 به‌عنوان دومین واحد نمونه انتخاب می‌شود. زوج بعدی یعنی $(8, 302)$ را باید نادیده گرفت و با زوج $(5, 97)$ ، واحد 5 را به‌عنوان سومین واحد نمونه برمی‌گزینیم. از زوج آخر، یعنی $(3, 25)$ ، واحد سه، چهارمین واحد نمونه خواهد بود. لذا مقادیر واحدهای نمونه مورد نظر به‌روش لاهیری به‌صورت

$$Y_5, Y_2, Y_5, Y_3$$

است، که در آن واحد پنجم، تکراری است.

▲

با توجه به مثال بالا می‌توان حالتی را تصور کرد که انتخابهای متوالی زوجهای (i, j) هرگز به انتخاب واحدی منجر نشود. در ادامه مطلب خواهیم دید که احتمال وقوع چنین حالتی صفر است. اینک که با ذکر مثال، شیوه لاهیری برای انتخاب واحدهای نمونه مشخص شد ثابت می‌کنیم که با این شیوه، احتمال انتخاب واحد i متناسب با X_i است.

برهان شیوه لاهیری. قبلاً متذکر می‌شویم که هر زوج (i, j) که به انتخاب Y_i منجر نشود به زوج نامؤثر موسوم است و اگر زوجی به انتخاب Y_i منجر شود آن را زوج مؤثر می‌نامند. حال زوج (Y_1, X_1) را در نظر می‌گیریم و احتمال انتخاب Y_i را در بار اول انتخاب زوج (i, j) با $P_1(Y_i)$ نشان می‌دهیم. به‌همین ترتیب احتمال انتخاب Y_i را در بار دوم انتخاب زوج (i, j) با $P_2(Y_i)$ نمایش می‌دهیم و الی آخر. احتمال $P_1(Y_i)$ را می‌توان به‌صورت زیر نوشت

$$P_1(Y_i) = P(j \leq X_i, 1 \leq j \leq M) \text{ تا } N \text{ انتخاب شود و سپس برای } 1 \leq j \leq M$$

چون انتخاب i و j مستقل از هم هستند، پس داریم

$$P_1(Y_i) = P(\text{انتخاب } i \text{ از اعداد } 1 \text{ تا } N) \cdot P(j \leq X_i | 1 \leq j \leq M) \quad (۱.۳)$$

اما

$$P(\text{انتخاب } i \text{ از اعداد } 1 \text{ تا } N) = \frac{1}{N} \quad (۲.۳)$$

از طرفی تعداد حالت‌های ممکن برای انتخاب j برابر با M است و حالت‌های مساعد برای وقوع پیشامد $X_i \leq j$ برابر با X_i است، پس

$$\begin{cases} P(j \leq X_i) = \frac{X_i}{M} \\ P(j > X_i) = 1 - \frac{X_i}{M} = \frac{M - X_i}{M} \end{cases} \quad (3.3)$$

اگر رابطه‌های (2.3) و (3.3) را در (1.3) منظور کنیم، نتیجه می‌شود

$$P_1(Y_i) = \frac{1}{N} \cdot \frac{X_i}{M} \quad (4.3)$$

درواقع، این احتمال، احتمال مؤثر بودن انتخاب زوج (i, j) است.

برای اینکه زوج (i, j) در اولین انتخاب مؤثر نباشد باید یکی از پیشامدهای زیر رخ دهد

$$i = 1, \text{ آن‌گاه } j > X_1$$

$$i = 2, \text{ آن‌گاه } j > X_2$$

...

...

$$i = N, \text{ آن‌گاه } j > X_N$$

لذا،

$P((i, j) \text{ مؤثر نبودن زوج})$

$$= P(i = 1, j > X_1) + P(i = 2, j > X_2) + \dots + P(i = N, j > X_N) \quad (5.3)$$

اما

$$P(i = 1, j > X_1) = P(i = 1) \cdot P(j > X_1) = \frac{1}{N} \cdot \frac{M - X_1}{M}$$

$$P(i = 2, j > X_2) = P(i = 2) \cdot P(j > X_2) = \frac{1}{N} \cdot \frac{M - X_2}{M}$$

.....

.....

.....

$$P(i = N, j > X_N) = P(i = N) \cdot P(j > X_N) = \frac{1}{N} \cdot \frac{M - X_N}{M}$$

اگر این رابطه‌ها را در (۵.۳) قرار دهیم، نتیجه می‌شود که

$$\begin{aligned} P((i, j) \text{ مؤثر نبودن زوج}) &= \sum_{i=1}^N \frac{1}{N} \cdot \frac{M - X_i}{M} = \frac{1}{NM} \sum_{i=1}^N (M - X_i) \\ &= \frac{1}{NM} (MN - N\bar{X}) = 1 - \frac{\bar{X}_N}{M} \end{aligned} \quad (۶.۳)$$

که در آن، \bar{X}_N میانگین X_i ‌هاست. حال $P_r(Y_i)$ را محاسبه می‌کنیم

$$\begin{aligned} P_r(Y_i) &= P((i, j) \text{ انتخاب در دومین انتخاب}) \\ &= P((i, j) \text{ مؤثر بودن انتخاب در بار اول} \cap (i, j) \text{ مؤثر نبودن انتخاب در بار اول}) \\ &= P(\text{مؤثر بودن انتخاب در بار اول}) \cdot P(\text{مؤثر نبودن انتخاب در بار اول}) \end{aligned}$$

که با توجه به رابطه‌های (۴.۳) و (۶.۳)، داریم

$$P_r(Y_i) = \left(1 - \frac{\bar{X}_N}{M}\right) \frac{1}{N} (X_i/M) \quad (۷.۳)$$

به همین ترتیب

$$\begin{aligned} P_r(Y_i) &= P\left(\text{مؤثر بودن انتخاب در بار سوم} \cap (i, j) \text{ مؤثر نبودن انتخاب در بار دوم} \cap (i, j) \text{ مؤثر نبودن انتخاب در بار اول}\right) \\ &= \left(1 - \frac{\bar{X}_N}{M}\right) \left(1 - \frac{\bar{X}_N}{M}\right) \frac{1}{N} (X_i/M) \end{aligned}$$

ممکن است متوالیاً، به دفعات بسیار زیاد، انتخاب (i, j) مؤثر نبوده و سپس انتخاب بعدی مؤثر باشد. حتی می‌توان تصور کرد که تعداد مؤثر نبودن‌ها به بینهایت هم بگراید. حال اگر هدف تعیین احتمال انتخاب واحد i با مقدار Y_i به‌عنوان واحدی از نمونه باشد باید در یکی از پیشامدهای متوالی که احتمال متناظر با آنها را در بالا مشخص کردیم Y_i انتخاب شود. پس

$$\begin{aligned} P(\text{انتخاب } Y_i \text{ در نمونه}) &= \frac{1}{N} \cdot \frac{X_i}{M} + \left(1 - \frac{\bar{X}_N}{M}\right) \frac{1}{N} \frac{X_i}{M} \\ &\quad + \left(1 - \frac{\bar{X}_N}{M}\right)^2 \frac{1}{N} \frac{X_i}{M} + \dots \end{aligned}$$

طرف دوم این برابری، یک سری هندسی است که قدر نسبت آن $(1 - \frac{\bar{X}_N}{M})$ است. چون M بزرگترین X_i است، $\frac{\bar{X}_N}{M}$ کوچکتر از واحد است و لذا $(1 - \frac{\bar{X}_N}{M} < 1)$ و مثبت است مگر اینکه همه X_i ها با هم برابر باشند که نمونه‌گیری به نمونه‌گیری تصادفی ساده تبدیل می‌شود. با این توضیح، سری طرف دوم رابطه بالا یک سری همگراست و وقتی تعداد انتخابهای ناموثر متوالی به بینهایت بگراید، داریم

$$\lim P(\text{انتخاب } Y_i \text{ در نمونه}) = \frac{\frac{1}{N} \frac{X_i}{M}}{1 - (1 - \frac{\bar{X}_N}{M})} = \frac{X_i}{N \bar{X}_N} = \frac{X_i}{T_N}$$

یعنی احتمال انتخاب Y_i به روش لاهیری برای شرکت در نمونه متناسب با X_i است. در این روش برای انتخاب واحدها نیازی به داشتن مجموع تراکمی X_i ها نداریم و انتخابها با محاسباتی کمتر صورت می‌گیرند.

تبصره. از رابطه (۶.۳) نتیجه می‌شود که انتخابهای متوالی (i, j) نمی‌توانند همیشه ناموثر باشند زیرا احتمال ناموثر بودن (i, j) برابر با $1 - \frac{\bar{X}_N}{M}$ است که همیشه از ۱ کمتر است و لذا پیشامدی حتمی نخواهد بود.

۴.۳ روش خرد کردن برای اصلاح شیوه لاهیری

دیدیم که احتمال موثر نبودن انتخاب زوج (i, j) برابر $(1 - \frac{\bar{X}_N}{M})$ است. هرچه این مقدار کوچکتر باشد یعنی هرچه کسر $\frac{\bar{X}_N}{M}$ به یک نزدیکتر باشد احتمال بالا کوچکتر بوده و در نتیجه تعداد انتخابهای ناموثر که موجب صرف هزینه و وقت می‌شود کمتر است. مقدار $\frac{\bar{X}_N}{M}$ وقتی به یک نزدیک است که M به میانگین X_i ها نزدیک باشد. اگر مقادیر X_i ها پراکندگی زیاد نداشته باشند، یعنی M به \bar{X}_N نزدیک باشد این منظور حاصل می‌شود. وقتی M ، یعنی بزرگترین X_i با سایر X_i ها و در نتیجه با \bar{X}_N فاصله زیاد داشته باشد می‌توان با روش خرد کردن که ذیلاً شرح می‌دهیم این فاصله را کم و تعداد انتخابهای ناموثر را نیز کم کرد. این روش همان‌طور که توضیح دادیم برای کوچکتر کردن $1 - \frac{\bar{X}_N}{M}$ ابداع شده است. در این روش آن واحدی را که مقدار X آن متناظر با M است به دو یا چند واحد تفکیک می‌کنیم. با این عمل تعداد واحدهای جامعه یک یا چند واحد اضافه می‌شود ولی مقدار M ، یعنی ماکسیمم X_i ها هم کوچک می‌شود. البته با این عمل \bar{X}_N هم کوچک می‌شود ولی تفکیک را به صورتی انجام می‌دهیم که $\frac{\bar{X}_N}{M}$ نهایتاً بزرگ و احتمال $1 - \frac{\bar{X}_N}{M}$ کوچک و در نتیجه تعداد انتخابهای ناموثر کم شود. برای توضیح مطلب به مثال زیر توجه کنید.

مثال ۳.۳ در جدول زیر میزان محصول گندم ۸ ده را با مساحت زمین زیرکشت آنها درج کرده‌ایم. هدف، انتخاب نمونه‌ای از این جامعه با احتمال متناسب با مساحت زیرکشت است.

| شماره واحد | میزان محصول Y_i | مساحت زمین زیرکشت X_i |
|------------|-------------------|-------------------------|
| ۱ | Y_1 | ۴۵ |
| ۲ | Y_2 | ۵۰ |
| ۳ | Y_3 | ۷۵ |
| ۴ | Y_4 | ۳۰ |
| ۵ | Y_5 | ۱۲۰ |
| ۶ | Y_6 | ۶۵ |
| ۷ | Y_7 | ۳۵ |
| ۸ | Y_8 | ۶۰ |

با روش لاهیری احتمال نامؤثر بودن انتخاب زوج (i, j) ، با توجه به مقادیر $N = 8$ و $\bar{X}_N = 60$ ، $M = 120$ ، برابر است با $0.5 - \frac{\bar{X}_N}{M} = 0.1$. یعنی باید انتظار داشته باشیم که نصف زوجهای (i, j) در انتخاب هر نمونه، نامؤثر باشند که مسلماً با اتلاف وقت و هزینه همراهاند. حال با استفاده از روش خرد کردن، روستایی که مساحت زیرکشت آن 120 است به دو روستا با مساحتهای زیرکشت 75 و 45 تقسیم می‌کنیم. در نتیجه حجم جامعه $N = 9$ و $M = 75$ و $\bar{X}_N = \frac{160}{3}$ خواهد شد، لذا

$$1 - \frac{\bar{X}_N}{M} = 1 - \frac{160/3}{75} = \frac{13}{45} \approx 0.29$$

ملاحظه می‌کنید که با این روش احتمال نامؤثر بودن زوج (i, j) از 0.5 به 0.29 کاهش می‌یابد. ▲

۵.۳ برآورد میانگین در نمونه‌گیری با احتمال متغیر و با جایگذاری

همان‌طور که در ابتدای فصل متذکر شدیم، حجم جامعه را N و مقدار صفت اصلی واحد i ام جامعه را Y_i و احتمال انتخاب این واحد در نمونه به حجم n و با جایگذاری را با p_i نشان می‌دهیم که $\sum_{i=1}^N p_i = 1$. فرض می‌کنیم با یکی از شیوه‌های قبل نمونه‌ای به حجم n انتخاب کرده‌ایم. حال می‌خواهیم برآوردکننده‌ای برای \bar{Y}_N ، میانگین جامعه، به دست آوریم.

قضیه ۱.۳ در نمونه‌گیری تصادفی با جایگذاری و با احتمال متغیر و به حجم n ، آماره $\frac{1}{n} \sum_{i=1}^n \frac{Y_i}{np_i}$ برآوردکننده ناریب میانگین جامعه است.

برهان. فرض کنیم نمونه‌ای که با احتمال متغیر و با جایگذاری انتخاب شده است دارای مقادیر و احتمالهای متناظر زیر باشد:

$$\begin{array}{ccccccc} Y_1 & Y_2 & \cdots & Y_i & \cdots & Y_n \\ p_1 & p_2 & \cdots & p_i & \cdots & p_n \end{array}$$

بدیهی است $\sum_{i=1}^n p_i < 1$. اگر میانگین جامعه \bar{Y}_N باشد می‌خواهیم از روی نمونه برآوردکننده‌ای

برای \bar{Y}_N بیابیم. روی مقادیر واحدهای جامعه تبدیل متغیر

$$Z_i = \frac{Y_i}{Np_i} \quad i = 1, 2, \dots, N$$

را اعمال می‌کنیم تا جامعه Z_i ها به وجود آید. میانگین این جامعه به صورت زیر به دست می‌آید

$$\bar{Z}_N = \sum_{i=1}^N Z_i p_i = \sum_{i=1}^N \frac{Y_i}{Np_i} \cdot p_i = \frac{1}{N} \sum_{i=1}^N Y_i = \bar{Y}_N$$

یعنی میانگین جامعه Z_i ها برابر با میانگین جامعه Y_i هاست، پس به جای تعیین برآوردکننده‌ای برای \bar{Y}_N باید برآوردکننده‌ای برای \bar{Z}_N بیابیم. ولی با توجه به آنچه در آمار مقدماتی دیده‌اید، اگر از جامعه متغیر تصادفی Z که تابع جرم احتمال آن $P(Z = Z_i) = p_i$ ، $i = 1, \dots, N$ است نمونه‌ای با جایگذاری، که استقلال واحدها را تأمین می‌کند، به تصادف انتخاب کنیم، و نمونه به صورت Z_1, Z_2, \dots, Z_n باشد، آن‌گاه $\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$ برآوردکننده ناریب \bar{Z}_N جامعه است. پس

$$\bar{Z}_n = \hat{Z}_N = \hat{Y}_N$$

اما با توجه به $Z_i = \frac{Y_i}{Np_i}$ ، برآوردکننده ناریب \bar{Y}_N ، یعنی \bar{Z}_n به صورت زیر است

$$\hat{Y}_N = \bar{Z}_n = \frac{1}{n} \sum_{i=1}^n \frac{Y_i}{Np_i}$$

و قضیه ثابت می‌شود. این برآوردکننده منسوب به هنس و هورویتس است. \square

مثال ۴.۳ از جامعه‌ای با ۵۰ واحد، نمونه‌ای به حجم ۶ و با احتمال متغیر و با جایگذاری انتخاب کرده‌ایم. مقادیر واحدها و احتمالهای متناظر با آنها در جدول زیر آمده‌اند

| | | | | | | |
|-------|----------------|----------------|----------------|----------------|----------------|----------------|
| Y_i | ۱۰ | ۱۲ | ۲۰ | ۸ | ۱۴ | ۱۶ |
| p_i | $\frac{1}{50}$ | $\frac{1}{45}$ | $\frac{1}{60}$ | $\frac{1}{40}$ | $\frac{1}{45}$ | $\frac{1}{70}$ |

برآورد ناریب میانگین جامعه را بیابید.

ابتدا تبدیل متغیر $Z_i = \frac{Y_i}{Np_i}$ را برای جدول داده‌ها اجرا می‌کنیم تا Z_i های متناظر به دست آیند

$$\begin{aligned} Z_1 &= \frac{Y_1}{Np_1} = \frac{10}{50 \left(\frac{1}{50}\right)} = 10 & , & \quad Z_2 = \frac{12}{50 \left(\frac{1}{45}\right)} = 10,8 \\ Z_3 &= \frac{20}{50 \left(\frac{1}{60}\right)} = 24 & , & \quad Z_4 = \frac{8}{50 \left(\frac{1}{40}\right)} = 6,4 \\ Z_5 &= \frac{14}{50 \left(\frac{1}{45}\right)} = 25,2 & , & \quad Z_6 = \frac{16}{50 \left(\frac{1}{70}\right)} = 22,4 \end{aligned}$$

پس نمونه Z_i ها به صورت زیر است

$$۱۰, ۱۰۸, ۲۴, ۶۴, ۲۵۲, ۲۲۴$$

بنابراین

$$\hat{Y}_N = \bar{Z}_n = \frac{1}{6}[۱۰ + ۱۰۸ + \dots + ۲۲۴] \simeq ۱۶۴۶$$

▲ که برآورد نااریب میانگین جامعه به حجم ۵۰ Y_i ها است.

۶.۳ واریانس برآوردکننده میانگین جامعه در نمونه‌گیری با احتمال متغیر و با جایگذاری

مسئله نمونه‌های متعددی به حجم m ، با روشهای مذکور می‌توان از جامعه مورد نظر انتخاب کرد. لذا پس از تبدیل Y_i های نمونه به Z_i ها، \bar{Z}_n های متعددی خواهیم داشت. هرچه \bar{Z}_n ها پراکندگی کمتری داشته باشند، برآمد تصادفی یکی از آنها به عنوان برآورد نااریب \bar{Y}_N با دقتی بیشتر همراه است. برای داشتن ایده‌ای از پراکندگی \bar{Z}_n ها، واریانس \bar{Z}_n تصادفی را محاسبه می‌کنیم. می‌دانیم

$$V(\bar{Z}_n) = E(\bar{Z}_n^2) - [E(\bar{Z}_n)]^2$$

اما $E(\bar{Z}_n) = \bar{Z}_N$ و $\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$ پس

$$\begin{aligned} V(\bar{Z}_n) &= E \left[\left(\frac{1}{n} \sum_{i=1}^n Z_i \right)^2 \right] - \bar{Z}_N^2 \\ &= \frac{1}{n^2} E \left[\left(\sum_{i=1}^n Z_i \right)^2 \right] - \bar{Z}_N^2 \\ &= \frac{1}{n^2} E \left[\sum_{i=1}^n Z_i^2 + \sum_{(i \neq j)=1}^n Z_i Z_j \right] - \bar{Z}_N^2 \\ &= \frac{1}{n^2} \left[\sum_{i=1}^n E(Z_i^2) + \sum_{(i \neq j)=1}^n E(Z_i Z_j) \right] - \bar{Z}_N^2 \end{aligned} \quad (۸.۳)$$

اما Z_i ها متناظر با احتمالهای p_i ها هستند، لذا

$$E(Z_i^2) = \sum_{i=1}^N p_i Z_i^2 \quad (۹.۳)$$

از طرفی Z_i مستقل از Z_j است، پس

$$E(Z_i Z_j) = E(Z_i)E(Z_j) = \bar{Z}_N \cdot \bar{Z}_N = \bar{Z}_N^2 \quad (۱۰.۳)$$

اگر رابطه‌های (۹.۳) و (۱۰.۳) را در (۸.۳) منظور کنیم، نتیجه می‌شود

$$\begin{aligned} V(\bar{Z}_n) &= \frac{1}{n^2} \left[\sum_{i=1}^n \sum_{i=1}^N p_i Z_i^2 + \sum_{(i \neq j)=1}^n \bar{Z}_N^2 \right] - \bar{Z}_N^2 \\ &= \frac{1}{n^2} \left[n \sum_{i=1}^N p_i Z_i^2 + n(n-1) \bar{Z}_N^2 \right] - \bar{Z}_N^2 \\ &= \frac{1}{n} \sum_{i=1}^N p_i Z_i^2 + \frac{n-1}{n} \bar{Z}_N^2 - \bar{Z}_N^2 \end{aligned}$$

و سرانجام

$$V(\bar{Z}_n) = \frac{1}{n} \left(\sum_{i=1}^N p_i Z_i^2 - \bar{Z}_N^2 \right)$$

عبارت داخل پرانتز برابر با واریانس جامعه Z_i است. اگر واریانس این جامعه را با σ_Z^2 نشان دهیم، داریم

$$V(\bar{Z}_n) = \frac{\sigma_Z^2}{n}$$

یعنی با توجه به اینکه \bar{Z}_n برآوردکننده \bar{Y}_N است

$$V(\hat{Y}_N) = \frac{\sigma_Z^2}{n}$$

از طرفی با توجه به تبدیل $Z_i = \frac{Y_i}{N p_i}$

$$\begin{aligned} \sigma_Z^2 &= \sum_{i=1}^N p_i Z_i^2 - \bar{Z}_N^2 = \sum_{i=1}^N p_i \cdot \left(\frac{Y_i}{N p_i} \right)^2 - \bar{Y}_N^2 \\ &= \sum_{i=1}^N \frac{Y_i^2}{N^2 p_i} - \bar{Y}_N^2 \end{aligned}$$

پس

$$V(\hat{Y}_N) = V(\bar{Z}_n) = \frac{1}{n} \left[\sum_{i=1}^N \frac{Y_i^2}{N^2 p_i} - \bar{Y}_N^2 \right] \quad (۱۱.۳)$$

تبصره ۱. اگر قرار دهیم $p_i = \frac{1}{N}$ ، نمونه‌گیری با احتمال متغیر به نمونه‌گیری تصادفی ساده با جایگذاری تبدیل می‌شود. رابطه $Z_i = \frac{Y_i}{Np_i}$ به $Z_i = Y_i$ تبدیل می‌شود، و

$$V(\hat{Y}_N) = V(\bar{Z}_n) = V(\bar{Y}_n) = \frac{\sigma_Y^2}{n}$$

تبصره ۲. اگر احتمال انتخاب Y_i متناسب با مقدار Y_i باشد، آنگاه

$$p_i = \frac{Y_i}{\sum_{i=1}^N Y_i}$$

و در نتیجه،

$$Y_i = p_i \sum_{i=1}^N Y_i$$

که اگر آن را در رابطه $Z_i = \frac{Y_i}{Np_i}$ منظور کنیم، نتیجه می‌شود

$$Z_i = \frac{p_i \sum_{i=1}^N Y_i}{Np_i} = \bar{Y}_N$$

یعنی Z_i به ازای هر i همیشه برابر با \bar{Y}_N است و لاجرم $\sigma_Z^2 = 0$ و این نمونه‌گیری کاراتر از نمونه‌گیری تصادفی ساده با جایگذاری است.

۷.۳ برآورد واریانس برآوردکننده میانگین در نمونه‌گیری با احتمال متغیر و با جایگذاری

در رابطه $V(\hat{Y}_N) = \frac{\sigma_Z^2}{n}$ ، چون σ_Z^2 مجهول است نمی‌توانیم واریانس برآوردکننده میانگین جامعه را به دست آوریم، لذا باید از روی نمونه برآوردی برای این واریانس بیابیم. می‌توان نشان داد که s_Z^2 ، تغییرات نمونه، برآوردکننده‌ای ناریب برای σ_Z^2 جامعه است (نظیر مورد نمونه‌گیری تصادفی ساده با جایگذاری). می‌دانیم

$$s_Z^2 = \frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2$$

پس

$$\begin{aligned} E(s_Z^2) &= \frac{1}{n-1} E \left[\sum_{i=1}^n (Z_i - \bar{Z}_n)^2 \right] = \frac{1}{n-1} E \left[\sum_{i=1}^n Z_i^2 - n\bar{Z}_n^2 \right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n E(Z_i^2) - nE(\bar{Z}_n^2) \right] \end{aligned} \quad (12.3)$$

از طرفی

$$V(\bar{Z}_n) = E(\bar{Z}_n^2) - \bar{Z}_N^2$$

پس

$$E(\bar{Z}_n^2) = \frac{\sigma_Z^2}{n} + \bar{Z}_N^2 \quad (۱۳.۳)$$

ضمناً می‌دانیم که

$$E(Z_i^2) = \sum_{i=1}^N p_i Z_i^2 \quad (۱۴.۳)$$

اگر (۱۳.۳) و (۱۴.۳) را در (۱۲.۳) قرار دهیم، نتیجه می‌شود

$$\begin{aligned} E(s_Z^2) &= \frac{1}{n-1} \left[\sum_{i=1}^n \sum_{i=1}^N p_i Z_i^2 - n \left(\frac{\sigma_Z^2}{n} + \bar{Z}_N^2 \right) \right] \\ &= \frac{1}{n-1} \left[n \sum_{i=1}^N p_i Z_i^2 - n \bar{Z}_N^2 - \sigma_Z^2 \right] \\ &= \frac{1}{n-1} \left[n \left(\sum_{i=1}^N p_i Z_i^2 - \bar{Z}_N^2 \right) - \sigma_Z^2 \right] \end{aligned}$$

اما عبارت داخل پرانتز برابر با σ_Z^2 است، پس

$$E(s_Z^2) = \frac{1}{n-1} (n\sigma_Z^2 - \sigma_Z^2) = \sigma_Z^2$$

یعنی s_Z^2 برآوردکننده نااریب σ_Z^2 است. لذا با توجه به برابری

$$V(\hat{Y}_N) = V(\bar{Z}_n) = \frac{\sigma_Z^2}{n}$$

برآوردکننده نااریب این واریانس به صورت زیر است

$$\hat{V}(\hat{Y}_N) = \frac{s_Z^2}{n}$$

از برآوردکننده هنس-هورویتس، برآوردکننده نااریب مجموع واحدهای جامعه به صورت زیر به دست می‌آید

$$\hat{T} = N\hat{Y}_N = \frac{1}{n} \sum_{i=1}^n \frac{Y_i}{p_i}$$

که با توجه به (۱۱.۳)، داریم

$$V(\hat{T}) = V(N\hat{Y}_N) = \frac{1}{n} \sum_{i=1}^N p_i \left(\frac{Y_i}{p_i} - T \right)^2$$

برآوردکنندهٔ ناریب این واریانس به‌وضوح عبارت است از

$$\hat{V}(\hat{T}) = \frac{1}{n(n-1)} \sum_{i=1}^n \left(\frac{Y_i}{p_i} - \hat{T} \right)^2$$

یک بازهٔ اطمینان $(1 - \alpha)$ درصد برای مجموع کل واحدهای جامعه، مبتنی بر نرمال بودن توزیع بزرگ نمونه‌ای برای \hat{T} ، به‌صورت زیر است

$$\hat{T} \pm z \sqrt{\hat{V}(\hat{T})}$$

که z نقطهٔ $\alpha/2$ بالایی توزیع نرمال استاندارد است. برای حجمهای نمونه‌ای کمتر از 50 ، باید از توزیع t با $n - 1$ درجه آزادی استفاده کرد.

مثال ۵.۳ با توجه به داده‌های مثال ۴.۳ برآوردی ناریب برای واریانس برآوردکنندهٔ میانگین جامعه بیابید.

در مثال ۴.۳، پس از استفاده از تبدیل $Z_i = \frac{Y_i}{Np_i}$ ، به‌جای نمونهٔ به‌حجم ۶ از جامعهٔ Y_i ها، نمونهٔ به‌حجم ۶ از جامعهٔ Z_i ها را به‌صورت زیر به‌دست آوردیم

$$10, \quad 108, \quad 24, \quad 64, \quad 252, \quad 224$$

مقدار \hat{Y}_N حاصل از این نمونه برابر 1646 است. این مقدار رخدادی از متغیر تصادفی \hat{Y}_N است. حال می‌خواهیم واریانس این متغیر تصادفی را برآورد کنیم. ابتدا s_Z^2 را به‌دست می‌آوریم

$$\begin{aligned} s_Z^2 &= \frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 = \frac{1}{n-1} \left(\sum_{i=1}^n Z_i^2 - n\bar{Z}_n^2 \right) \\ &= \frac{1}{5} [10^2 + \dots + 224^2 - 6(1646)^2] = \frac{1}{5} (197040 - 162559) \\ &\approx 6896 \end{aligned}$$

پس

$$\hat{V}(\hat{Y}_N) = \frac{1}{6} (6896) \approx 1149$$

تبصره. برای کاربردها می‌توان $\hat{V}(\hat{Y}_N)$ را به صورت زیر نوشت

$$\begin{aligned}\hat{V}(\hat{Y}_N) &= \frac{s_Z^2}{n} = \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 \\ &= \frac{1}{n(n-1)} \sum_{i=1}^n \left(\frac{Y_i}{Np_i} - \hat{Y}_N \right)^2 \\ &= \frac{1}{n(n-1)N^2} \sum_{i=1}^n \left(\frac{Y_i}{p_i} - N\hat{Y}_N \right)^2\end{aligned}\quad (۱۵.۳)$$

۸.۳ مقایسه نمونه‌گیری تصادفی با احتمال متغیر و نمونه‌گیری

تصادفی ساده (حالت با جایگذاری) در برآورد میانگین جامعه

از جامعه به حجم N نمونه‌ای تصادفی به حجم n ، با احتمال متغیر و با جایگذاری به دست می‌آوریم. با این داده‌ها میانگین جامعه را برآورد می‌کنیم. حال تصور می‌کنیم این داده‌ها حاصل یک نمونه‌گیری تصادفی ساده با جایگذاری است. یک بار دیگر با این تصور، میانگین جامعه را برآورد می‌نماییم. اگر این دو برآوردکننده را به ترتیب با \bar{Z}_n و \bar{Y}_{Ran} نشان دهیم، برای مقایسه دقت این دو برآوردکننده باید واریانس آنها را با هم مقایسه کنیم. اگر مقادیر واحدهای نمونه Y_1, Y_2, \dots, Y_n باشند، مطابق (۱۱.۳) داریم

$$V(\bar{Z}_n) = \frac{1}{n} \left[\sum_{i=1}^N \frac{Y_i^2}{N^2 p_i} - \bar{Y}_N^2 \right]$$

و

$$\begin{aligned}V(\hat{Y}_{\text{Ran}}) &= \frac{\sigma^2}{n} = \frac{1}{nN} \sum_{i=1}^N (Y_i - \bar{Y}_N)^2 = \frac{1}{nN} \left[\sum_{i=1}^N Y_i^2 - N\bar{Y}_N^2 \right] \\ &= \frac{1}{n} \left[\sum_{i=1}^N \frac{Y_i^2}{N} - \bar{Y}_N^2 \right]\end{aligned}$$

برای مقایسه دقت دو روش تفاضل دو واریانس را حساب می‌کنیم

$$\begin{aligned}d = V(\bar{Z}_n) - V(\bar{Y}_{\text{Ran}}) &= \left[\frac{1}{n} \sum_{i=1}^N \frac{Y_i^2}{N^2 p_i} - \frac{1}{n} \bar{Y}_N^2 \right] - \left[\frac{1}{n} \sum_{i=1}^N \frac{Y_i^2}{N} - \frac{1}{n} \bar{Y}_N^2 \right] \\ &= \frac{1}{nN} \sum_{i=1}^N \left(\frac{Y_i^2}{Np_i} - Y_i^2 \right) \\ &= \frac{1}{nN} \sum_{i=1}^N \left(\frac{1}{Np_i} - 1 \right) Y_i^2\end{aligned}$$

حال اگر به جای Y_i مقدار آن $\sum \frac{X_i}{N X_i}$ یا $\frac{X_i}{N X_i}$ را قرار دهیم نتیجه می‌شود که

$$d = \frac{1}{nN} \sum_{i=1}^N \left(\frac{\sum X_i}{N X_i} - 1 \right) Y_i^2 = \frac{1}{nN} \sum_{i=1}^N \left(\frac{\bar{X}_N}{X_i} - 1 \right) Y_i^2$$

پس

$$d = \frac{1}{nN} \left(\bar{X}_N \sum_{i=1}^N \frac{Y_i^2}{X_i} - \sum_{i=1}^N Y_i^2 \right) \quad (۱۶.۳)$$

اگر داشته باشیم

$$\bar{X}_N \sum_{i=1}^N \frac{Y_i^2}{X_i} < \sum_{i=1}^N Y_i^2$$

در نتیجه، $d < 0$ و نمونه‌گیری با احتمال متغیر و با جایگذاری دقیقتر از نمونه‌گیری تصادفی ساده با جایگذاری است. نابرابری بالا را می‌توان به صورت زیر نوشت

$$\sum_{i=1}^N \frac{Y_i^2}{X_i} < \sum_{i=1}^N \frac{Y_i^2}{\bar{X}_N}$$

این نابرابری برحسب پارامترهای جامعه اصلی و جامعه کمکی بیان شده است، لذا کاربرد عملی در تعیین اینکه کدامیک از دو نمونه‌گیری در جامعه‌ای مفروض مناسبتر است ندارد. نابرابری اخیر را گاهی به صورت زیر بیان می‌کنند

$$\sum_{i=1}^N (X_i - \bar{X}_N) \frac{Y_i^2}{X_i} > 0 \quad (۱۷.۳)$$

راج نشان داده است که این رابطه وقتی X_i ها و $\frac{Y_i^2}{X_i}$ ها به صورت مثبت وابسته‌اند برقرار است و در چنین حالتی نمونه‌گیری با احتمال متغیر و با جایگذاری بهتر از نمونه‌گیری تصادفی ساده با جایگذاری است. تبصره. اگر رابطه بین صفت اصلی و صفت کمکی به صورت خطی $Y = a + bX$ ، یعنی همبستگی بین Y و X کامل باشد با استفاده از نابرابری نتیجه می‌شود که نمونه‌گیری با احتمال متغیر وقتی دقت کمتری دارد که

$$\frac{\bar{X} - \tilde{X}}{\bar{X} \sigma_X^2} > \frac{b^2}{a^2}$$

که در آن $\tilde{X} = \frac{\sum (X_i^2)}{\sum X_i}$

۹.۳ نمونه‌گیری با احتمال متغیر و بدون جایگذاری

جامعه‌ای N واحدی با مقادیر $Y_1, Y_2, \dots, Y_i, \dots, Y_N$ را در نظر می‌گیریم. فرض می‌کنیم $p_1, p_2, \dots, p_i, \dots, p_N$ احتمالهای متناظر با واحدهای جامعه باشند. می‌خواهیم با روش بدون جایگذاری نمونه‌ای به حجم n از این جامعه انتخاب کنیم به قسمی که احتمال انتخاب واحد i ام جامعه، $i = 1, \dots, N$ در هر استخراج برابر p_i باشد. واحد اول را نظیر حالت با جایگذاری به یکی از دو روشی که متذکر شدیم به دست می‌آوریم. اگر مثلاً واحد i ام در انتخاب اول برای عضویت در نمونه انتخاب شود همان‌طور که دیدیم

$$P(\text{استخراج } Y_i \text{ در بار اول}) = p_i, \quad i = 1, \dots, N$$

حال این واحد را به جامعه برمی‌گردانیم. جامعه‌ای که فعلاً در اختیار داریم دارای مقادیر

$$Y_1, Y_2, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_N$$

است که حجم آن $N - 1$ است. احتمالهای متناظر با این واحدها در جامعه اصلی با دنباله زیر مشخص می‌شود

$$p_1, p_2, \dots, p_{i-1}, p_{i+1}, \dots, p_N$$

در مورد این جامعه جدید و احتمالهای متناظر با واحدهای آن، نمی‌توانیم به دو روش مذکور، واحد دوم نمونه را انتخاب کنیم زیرا مجموع این احتمالها برابر با یک نیست. برای تهیه واحدهای نمونه در حالت بدون جایگذاری پیشنهادهایی متعدد که کم و بیش با پیچیدگی همراه‌اند وجود دارند. این روشها را می‌توان عمدتاً در مراجع [۸] و [۱۳] یافت. از جمله این روشها، روش راتو و همکاران، روش دوربین^۱، روش ساتتر^۲، و روش هدایت-لین است. ما برای تعیین واحدهای دوم، سوم، ... و n ام نمونه بدون جایگذاری به اثبات و ذکر جزئیات مطالب نمی‌پردازیم، و تنها روشی را که به برآورد مرتب موسوم است به صورتی مشروح مطرح می‌کنیم. در پایان نیز به اختصار شیوه نمونه‌گیری راتو-هارتلی-کوکران را بررسی می‌نماییم.

۱.۹.۳ برآورد مرتب

پس از انتخاب اولین واحد نمونه که فرض می‌کنیم واحد i ام جامعه است، جامعه باقی‌مانده را که به حجم $N - 1$ است با مقادیر و احتمالهای متناظری که در زیر آورده‌ایم در نظر می‌گیریم

$$\begin{array}{cccccc} Y_1 & Y_2 & \dots & Y_{i-1} & Y_{i+1} & \dots & Y_N \\ \frac{p_1}{1-p_i} & \frac{p_2}{1-p_i} & \dots & \frac{p_{i-1}}{1-p_i} & \frac{p_{i+1}}{1-p_i} & \dots & \frac{p_N}{1-p_i} \end{array} \quad (18.3)$$

بدین طریق، احتمالهای متناظر با واحدها در جامعه اصلی را با تقسیم بر $1 - p_i$ به احتمالهایی متناسب که مجموع آنها یک است تبدیل کرده‌ایم. واضح است که

$$\sum_{\substack{j=1 \\ j \neq i}}^n \frac{p_j}{1 - p_i} = \frac{1}{1 - p_i} \sum_{\substack{j=1 \\ j \neq i}}^N p_j = \frac{1 - p_i}{1 - p_i} = 1$$

حال در این جامعه جدید به حجم $N - 1$ ، با توجه به احتمالهای متناظر با واحدها به یکی از دو روش مجموع تراکمی یا لاهیری، واحدی را که اولین انتخاب از این جامعه است به‌عنوان دومین واحد نمونه بدون جایگذاری از جامعه اصلی برمی‌گزینیم. با همین فرایند جامعه باقی‌مانده را که به حجم $N - 2$ است در نظر می‌گیریم، و اگر Y_i و Y_j دو واحد منتخب قبلی باشند، این جامعه و احتمالهای متناظر با واحدهای آن را به‌صورت زیر می‌نویسیم

$$\begin{array}{cccccccc} Y_1 & \dots & Y_{j-1} & Y_{j+1} & \dots & Y_{i-1} & Y_{i+1} & \dots & Y_N \\ \frac{p_1}{1 - p_i - p_j} & \dots & \frac{p_{j-1}}{1 - p_i - p_j} & \frac{p_{j+1}}{1 - p_i - p_j} & \dots & \frac{p_{i-1}}{1 - p_i - p_j} & \frac{p_{i+1}}{1 - p_i - p_j} & \dots & \frac{p_N}{1 - p_i - p_j} \end{array}$$

باز هم احتمالها متناسب‌اند و مجموع آنها یک است. به یکی از دو روش نمونه‌گیری با جایگذاری، یک واحد از این جامعه را به‌عنوان سومین واحد نمونه اصلی انتخاب می‌کنیم. این فرایند را تا انتخاب n واحد نمونه ادامه می‌دهیم. بدیهی است انجام این فرایند به‌کمک کامپیوتر چندان مشکل نخواهد بود.

حال نمونه حاصل از این فرایند را در نظر می‌گیریم، اولین واحد نمونه، واحد شماره i با مقدار Y_i از جامعه اصلی است و احتمال متناظر با آن p_i است. اگر قرار دهیم

$$Z_1 = \frac{Y_i}{N p_i} \quad (19.3)$$

داریم

$$E(Z_1) = \sum_{i=1}^N \frac{Y_i}{N p_i} \cdot p_i = \bar{Y}_N$$

دومین واحد نمونه، واحد شماره j با مقدار Y_j از جامعه است، و احتمال متناظر با آن p_j است. اگر قرار دهیم

$$Z_2 = \frac{1}{N} \left[Y_i + Y_j \frac{1 - p_i}{p_j} \right] \quad (20.3)$$

آنگاه برای Z_i ثابت، با توجه به (۱۸.۳)

$$\begin{aligned} E \left[Y_j \frac{1-p_i}{p_j} \middle| Y_i \right] &= \sum_{\substack{j=1 \\ j \neq i}}^N Y_j \frac{1-p_i}{p_j} \cdot \frac{p_j}{1-p_i} \\ &= \sum_{\substack{j=1 \\ j \neq i}}^N Y_j = N\bar{Y}_N - Y_i \end{aligned}$$

لذا

$$E(Z_i) = \frac{1}{N} \left[E(Y_i) + E \left(Y_j \frac{1-p_i}{p_j} \right) \right]$$

چون i ثابت است

$$= \frac{1}{N} [Y_i + N\bar{Y}_N - Y_i] = \bar{Y}_N$$

اگر روابطی مشابه با (۲۰.۳) برای واحدهای سوم و چهارم و ... و n ام نمونه بنویسیم، میانگین هر Z_i در دنباله Z_1, Z_2, \dots, Z_n برابر با \bar{Y}_N است. لذا میانگین

$$\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$$

نیز برابر با \bar{Y}_N است. یعنی \bar{Z}_n برآوردکننده نااریب \bar{Y}_N است.

در مورد روابط مشابه با (۲۰.۳)، می‌توان مطلب را به صورت کلی زیر مطرح کرد:

اگر Y_1, Y_2, \dots, Y_n مقادیر واحدهای نمونه باشند که با توضیح بالا، بدون جایگذاری انتخاب شده‌اند، و اگر p_1, \dots, p_n احتمالهای متناظر با آنها فرض شوند، متغیر Z_i را که تعمیمی از (۲۰.۳) است به صورت زیر می‌نویسیم

$$Z_i = \frac{1}{N} \left[Y_1 + Y_2 + \dots + Y_{i-1} + Y_i \frac{1-(p_1+p_2+\dots+p_{i-1})}{p_i} \right] \quad (21.3)$$

$i = 1, 2, \dots, n$

در این صورت \bar{Z}_n برای دنباله Z_1, Z_2, \dots, Z_n برآوردکننده نااریب \bar{Y}_N است. می‌توان ثابت

$$\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$$

$$\hat{V}(\bar{Z}_n) = \hat{V}(\bar{Y}_N) = \frac{1}{n(n-1)} \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 \quad (22.3)$$

که امکان می‌دهد واریانس برآوردکننده \bar{Y}_N را برآورد کنیم.

مثال ۶.۳ جدول زیر درآمد روزانه ۱۰ خانوار و تعداد افراد هر خانوار را نشان می‌دهد. Y_i درآمد روزانه و X_i تعداد افراد خانواده نام است که اولی صفت اصلی و دومی صفت کمکی است.

| | | | | | | | | | | |
|------------------|---|---|----|-----|----|----|----|----|----|-----|
| شماره خانوار | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ | ۱۰ |
| Y_i برحسب ۱۰۰۰ | ۶ | ۸ | ۷ | ۶٫۵ | ۵ | ۴ | ۹ | ۷ | ۸ | ۷٫۵ |
| X_i | ۳ | ۴ | ۴ | ۳ | ۲ | ۳ | ۵ | ۳ | ۴ | ۳ |
| T_i | ۳ | ۷ | ۱۱ | ۱۴ | ۱۶ | ۱۹ | ۲۴ | ۲۷ | ۳۱ | ۳۴ |

می‌خواهیم با نمونه‌ای به حجم ۴ بدون جایگذاری، برآوردی برای میانگین درآمد روزانه به دست آوریم. ابتدا با توجه به شیوه مجموع تراکمی، با مراجعه به جدول اعداد تصادفی، عددی از اعداد ۱ تا ۳۴ اختیار می‌کنیم. فرض کنید ۲۲ به دست آید. در این صورت $Y_7 = 9$ اولین واحد نمونه است. احتمالهای متناظر با Y_i ها در جامعه اصلی به صورت زیرند

$$\frac{3}{34}, \frac{4}{34}, \frac{4}{34}, \frac{3}{34}, \frac{2}{34}, \frac{3}{34}, \frac{5}{34}, \frac{3}{34}, \frac{4}{34}, \frac{3}{34}$$

چون Y_7 را به عنوان اولین واحد نمونه انتخاب کرده‌ایم و نمونه‌گیری بدون جایگذاری است، جامعه‌ای جدید به حجم ۹ تشکیل می‌شود که احتمالهای متناظر با واحدهای آن، طبق آنچه گفتیم از $\frac{p_i}{1-p_i}$ به دست می‌آیند. چون $p_7 = \frac{5}{34}$ داریم

| | | | | | | | | | |
|------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| واحدهای جامعه | Y_1 | Y_2 | Y_3 | Y_4 | Y_5 | Y_6 | Y_8 | Y_9 | Y_{10} |
| احتمالهای متناظر | $\frac{3}{29}$ | $\frac{4}{29}$ | $\frac{4}{29}$ | $\frac{3}{29}$ | $\frac{2}{29}$ | $\frac{3}{29}$ | $\frac{3}{29}$ | $\frac{4}{29}$ | $\frac{3}{29}$ |

بنابراین، جدول زیر را داریم

| | | | | | | | | | |
|------------------|---|---|----|-----|----|----|----|----|-----|
| شماره خانوار | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۸ | ۹ | ۱۰ |
| Y_i برحسب ۱۰۰۰ | ۶ | ۸ | ۷ | ۶٫۵ | ۵ | ۴ | ۷ | ۸ | ۷٫۵ |
| X_i | ۳ | ۴ | ۴ | ۳ | ۲ | ۳ | ۳ | ۴ | ۳ |
| T_i | ۳ | ۷ | ۱۱ | ۱۴ | ۱۶ | ۱۹ | ۲۲ | ۲۶ | ۲۹ |

ملاحظه می‌کنید سه سطر اول، همان سه سطر اول جدول قبل هستند که Y_7 آنها کنار گذاشته شده است. به شیوه مجموع تراکمی، با مراجعه به جدول اعداد تصادفی، عددی از اعداد ۱ تا ۲۹ اختیار می‌کنیم. فرض کنید ۹ به دست آید، لذا $Y_3 = 7$ دومین واحد نمونه است. از حذف Y_7 در جدول بالا سطر T_i به صورت زیر درمی‌آید

$$Y_i: Y_1 \ Y_2 \ Y_3 \ Y_4 \ Y_5 \ Y_6 \ Y_8 \ Y_9 \ Y_{10}$$

$$T_i: 3 \ 7 \ 10 \ 12 \ 15 \ 18 \ 22 \ 25$$

از اعداد ۱ تا ۲۵ عددی به تصادف برمی‌گزینیم. فرض کنید ۲۴ به دست آید. پس $Y_{10} = 7.5$ سومین واحد نمونه است. پس از حذف این واحد، داریم

$$Y_i : Y_1 \ Y_2 \ Y_4 \ Y_5 \ Y_6 \ Y_8 \ Y_9$$

$$T_i : 3 \ 7 \ 10 \ 12 \ 15 \ 18 \ 22$$

از اعداد ۱ تا ۲۲، عددی به تصادف اختیار می‌کنیم. فرض کنید ۱۱ به دست آید، پس $Y_5 = 5$ چهارمین واحد نمونه است. بنابراین نمونه منتخب به صورت زیر است

$$i : \text{شماره واحد منتخب} \quad 7 \quad 3 \quad 10 \quad 5$$

$$Y_i \text{ مقادیر} : \quad 9 \quad 7 \quad 7,5 \quad 5$$

$$\text{احتمالهای متناظر} : \quad \frac{5}{34} \quad \frac{4}{34} \quad \frac{3}{34} \quad \frac{2}{34}$$

Z_i ها را مطابق فرمول (۲۱.۳) محاسبه می‌کنیم

$$Z_1 = \frac{Y_1}{Np_1} = \frac{9}{10 \left(\frac{5}{34}\right)} = \frac{306}{50}$$

$$Z_2 = \frac{1}{N} \left(Y_1 + Y_2 \cdot \frac{1-p_1}{p_2} \right) = \frac{1}{10} \left[9 + 7 \left(\frac{\frac{21}{34}}{\frac{4}{34}} \right) \right] = \frac{1}{10} \left[9 + \frac{203}{4} \right] = \frac{121}{20}$$

$$Z_3 = \frac{1}{N} \left(Y_1 + Y_2 + Y_3 \cdot \frac{1-p_1-p_2}{p_3} \right) = \frac{1}{10} \left[9 + 7 + 7,5 \left(\frac{\frac{25}{34}}{\frac{3}{34}} \right) \right] = \frac{157}{20}$$

$$Z_4 = \frac{1}{N} \left(Y_1 + Y_2 + Y_3 + Y_4 \cdot \frac{1-p_1-p_2-p_3}{p_4} \right) = \frac{1}{10} \left[23,5 + 5 \left(\frac{\frac{22}{34}}{\frac{2}{34}} \right) \right] = \frac{157}{20}$$

$$\hat{Y}_N = \bar{Z}_n = \frac{1}{4} \left(\frac{306}{50} + \frac{121}{20} + \frac{157}{20} + \frac{157}{20} \right) \simeq 6,96$$

از رابطه (۲۲.۳) می‌توان برآورد واریانس \hat{Y}_N را که همان برآورد واریانس \hat{Z}_n است محاسبه کرد. \blacktriangle

۲.۹.۳ طرح نمونه‌گیری راثو-هارتلی-کوکران و برآوردکننده وابسته به آن

در این بخش رهبرد نمونه‌گیری معروف و راحت دیگری را که راثو و همکاران او پیشنهاد کرده‌اند و مبتنی بر صفت کمکی است متذکر می‌شویم. موضوع را با گروهبندی تصادفی جامعه تحت بررسی آغاز می‌کنیم.

جامعه P به حجم N و مجموعه اعداد صحیح N_1, N_2, \dots, N_n را در نظر می‌گیریم، به قسمی که برای $1 \leq i \leq n$ ، $N_i \geq 1$ و $\sum_{i=1}^n N_i = N$. در این صورت می‌گوییم $\{G_1, G_2, \dots, G_n\}$ یک گروهبندی تصادفی $(N, n; N_1, N_2, \dots, N_n)$ از جامعه P است اگر

$$P = \bigcup_{i=1}^n G_i$$

و برای هر $1 \leq i \leq n$ ، تعداد عناصر موجود در G_i عبارت باشد از $|G_i| = N_i$ و به علاوه برای هر i در تشکیل G_i ، هر N_i تایی جامعه P شانس برابر برای انتخاب شدن داشته باشد. یک راه

تشکیل چنین گروهبندی تصادفی این است که فرض کنیم G_1 یک نمونه تصادفی از P است که به روش تصادفی ساده به حجم N_1 انتخاب شده است، G_2 نمونه‌ای تصادفی از $P - G_1$ است که از جامعه به حجم $N - N_1$ انتخاب شده است و نظایر آن. در مرحله $(n-1)$ ام، نمونه G_{n-1} به حجم N_{n-1} از جامعه به حجم $N - \sum_{i=1}^{n-2} N_i$ به روش تصادفی ساده انتخاب می‌شود، و بالاخره در مرحله n ام، G_n از باقیمانده واحدهای جامعه تشکیل می‌شود.

حال طرح نمونه‌گیری راثو-هارتلی-کوکران (RHC) را از جامعه P به حجم N با احتمالهای متناظر p_1, p_2, \dots, p_N برای ورود به نمونه شرح می‌دهیم.

روش، شامل دو مرحله است. در زیر N_1, N_2, \dots, N_n را به عنوان مقادیر دلخواه ولی ثابت در نظر می‌گیریم که $1 \leq i \leq n, N_i \geq 1$ و $\sum_{i=1}^n N_i = N$.

مرحله ۱. ابتدا یک گروهبندی تصادفی $(N, n; N_1, N_2, \dots, N_n)$ جامعه P را که مثلاً به صورت $\{G_1, G_2, \dots, G_n\}$ است تشکیل می‌دهیم.

مرحله ۲. با استفاده از روش PPS نمونه‌گیری مستقل در داخل هر گروه، یک واحد از هر گروه استخراج می‌کنیم. بنابر روش PPS، با فرض داشتن گروهبندی تصادفی، شانس انتخاب i امین واحد در نمونه، اگر $i \in G_K$ برابر با $\frac{p_i}{P_K}$ است که در آن $P_K = \sum_{j \in G_K} p_j$ با شرط $1 \leq i \leq N, 1 \leq K \leq n$.

برآوردکننده مجموع جامعه که با طرح RHC پیشنهاد می‌شود به صورت زیر است

$$\hat{T}_{\text{RHC}} = \sum_{K=1}^n Y_{i_K} P_K / p_{i_K}$$

که در آن $i_1, i_2, \dots, i_K, \dots, i_n$ نمونه منتخبی است که واحد i_K از G_K ، $1 \leq K \leq n$ انتخاب شده است. در زمینه طرح RHC قضیه زیر را بدون اثبات ذکر می‌کنیم.

قضیه ۲.۳ تحت طرح نمونه‌گیری RHC، برآوردکننده \hat{T}_{RHC} برای مجموع واحدهای جامعه ناریب است و واریانس آن به صورت زیر است

$$V(\hat{T}_{\text{RHC}}) = \left\{ \left(\sum_{K=1}^n N_K^2 - N \right) / (N^2 - N) \right\} \left\{ \sum_{i=1}^N \left(\frac{Y_i}{p_i} - T \right)^2 p_i \right\}$$

و به علاوه، برآورد ناریب این واریانس به صورت زیر است

$$\hat{V}(\hat{T}_{\text{RHC}}) = \left\{ \left(\sum_{K=1}^n N_K^2 - N \right) / \left(N^2 - \sum_{K=1}^n N_K^2 \right) \right\} \left\{ \sum_{K=1}^n \left(\frac{Y_{i_K}}{p_{i_K}} - \hat{T}_{\text{RHC}} \right)^2 P_K \right\}$$

تبصره. تا اینجا فرض کردیم N_i ها متفاوت، دلخواه ولی تثبیت شده‌اند و $N_K \geq 1$ ، $\sum_{K=1}^n N_K = N$ ، $1 \leq K \leq n$ ، با بررسی فرمول $V(\hat{T}_{RHC})$ می‌توان نتیجه گرفت که انتخاب N_1, N_2, \dots, N_n که این واریانس را مینیمم می‌کند متناظر است با

$$N_1 = N_2 = \dots = N_n = \frac{N}{n} \quad \text{اگر } N \text{ بر } n \text{ تقسیم‌پذیر باشد}$$

و

$$\begin{cases} N_1 = N_2 = \dots = N_u = m + 1 \\ N_{u+1} = N_{u+2} = \dots = N_n = m \end{cases} \quad \text{اگر } N = mn + u, \quad 1 \leq u \leq n - 1$$

به عبارت دیگر برای اینکه طرح نمونه‌گیری RHC را به مفیدترین صورت درآوریم باید گروهبندی تصادفی را به نحوی انجام دهیم که حجمهای گروهی تا حد ممکن به هم نزدیک باشند. چنگ^۱ و لی^۲ (۱۹۸۳، ۱۹۸۷) و داس‌گوپتا و سینها (۱۹۸۳) بررسیهایی در زمینه ویژگیهای نظیر مینیماکس بودن نمونه‌گیری RHC با انتخاب اپتیم حجمهای گروهی انجام داده‌اند.

مثال ۷.۳ جامعه‌ای متشکل از ۱۵ واحد است. در زیر مقادیر صفت‌های کمکی متناظر با آنها را آورده‌ایم

| | | | | | | | | |
|------------------|-----|----|----|----|----|----|----|----|
| i شماره واحدها | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ |
| X_i صفت کمکی | ۲۳ | ۳ | ۴۱ | ۲۶ | ۵۳ | ۶ | ۲۸ | ۵۲ |
| i شماره واحدها | ۹ | ۱۰ | ۱۱ | ۱۲ | ۱۳ | ۱۴ | ۱۵ | |
| X_i صفت کمکی | ۱۱۳ | ۷۲ | ۸ | ۳۵ | ۴۲ | ۳۸ | ۵۲ | |

می‌خواهیم مجموع واحدهای جامعه را بر مبنای نمونه‌ای به حجم ۳ با استفاده از شیوه RHC برآورد کنیم. قرار می‌دهیم $N_1 = N_2 = N_3 = 5$ که متناظر با انتخاب اپتیم N_i هاست.

مرحله ۱. به جدول اعداد تصادفی رجوع می‌کنیم. فرض می‌کنیم گروهبندی تصادفی زیر حاصل شود

$$G_1 = \{7, 3, 12, 8, 1\} \quad G_2 = \{9, 14, 6, 2, 10\}$$

$$G_3 = \{4, 5, 11, 13, 15\}$$

مرحله ۲. با توجه به داده‌ها $\sum X_i = 745$ ، لذا از $p_i = \frac{X_i}{\sum X_i}$ همه مقادیر p_i ، $1 \leq i \leq 15$ مشخص می‌شوند. حال با واحدهای G_1, G_2, G_3 ، مقادیر p_i متناظر را همراه می‌کنیم و با استفاده از روش PPS از هر کدام یک واحد انتخاب می‌کنیم. می‌توان فرض کرد سه واحد منتخب عبارت‌اند از

$$i_1 = 7 \quad i_2 = 2 \quad i_3 = 13$$

حال به واحد شماره ۷، ۲، و ۱۳ جامعه مراجعه کرده مشخصه مورد نظر را مشاهده می‌کنیم. اگر اندازه مشخصه‌ها به صورت زیر باشند

$$Y_V = 30 \quad Y_r = 32 \quad Y_{13} = 45$$

آن‌گاه محاسبات زیر را انجام می‌دهیم

$$P_1 = \sum_{i \in G_1} p_i = \frac{X_V + X_r + X_{13} + X_8 + X_1}{\sum_i X_i} = \frac{278 + 41 + 35 + 52 + 23}{745}$$

$$= 0.2403$$

$$P_r = \sum_{i \in G_r} p_i = \frac{X_9 + X_{14} + X_6 + X_r + X_{10}}{\sum_i X_i} = \frac{113 + 38 + 6 + 3 + 72}{745}$$

$$= 0.4201$$

$$P_r = \sum_{i \in G_r} p_i = 0.3396$$

$$p_{i_1} = p_V, \quad \frac{p_V}{P_1} = \frac{\frac{278}{745}}{\frac{179}{745}} = \frac{278}{179} = 0.1564$$

$$p_{i_r} = p_r, \quad \frac{p_r}{P_r} = \frac{3}{313} = 0.00958$$

$$p_{i_r} = p_{13}, \quad \frac{p_{13}}{p_r} = \frac{42}{253} = 0.1660$$

لذا

$$\hat{T}_{RHC} = \sum_{K=1}^n Y_{i_K} P_K / p_{i_K}$$

$$= 30 \left(\frac{1}{0.1564} \right) + 32 \left(\frac{1}{0.00958} \right) + 45 \left(\frac{1}{0.1660} \right)$$

$$\# 191,8159 + 334,0292 + 271,0843$$

$$= 796,9294$$

▲ از رابطه $\hat{V}(\hat{T}_{RHC})$ برآورد واریانس \hat{T}_{RHC} برابر با 20814 به دست می‌آید.

تمرینها

۱. مشخصه مورد بررسی، میزان پنبه تولیدشده در ۱۰ روستای گرگان است. مساحت زمین زیر کشت پنبه در این ۱۰ روستا به ترتیب ۱۴، ۱۲، ۱۰، ۱۶، ۱۸، ۲۰، ۱۰، ۱۴، ۱۷، ۲۰ هکتار است.

ای بررسی میزان پنبه، ابتدا باید نمونه‌ای به حجم ۴ با احتمال متغیر و متناسب با مساحت‌های متناظر انتخاب کنیم. این نمونه را یک بار با روش جایگذاری و بار دیگر با روش بدون جایگذاری مشخص کنید. ۲. از جامعه‌ای که به حجم ۵۰ است، نمونه‌ای به حجم ۸، با احتمال متغیر و به روش تصادفی با جایگذاری انتخاب کرده، مشخصه واحدهای منتخب را همراه با احتمال متناظر با آنها در جدول زیر ثبت کرده‌ایم

| | | | | | | | | |
|-------|-----------------|-----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| Y_i | ۱۰۰ | ۱۰۴ | ۱۱۵ | ۱۱۲ | ۱۲۰ | ۱۲۵ | ۱۲۲ | ۱۵۰ |
| P_i | $\frac{۳}{۱۰۰}$ | $\frac{۲}{۱۰۰}$ | $\frac{۱}{۵۰}$ | $\frac{۱}{۵۰}$ | $\frac{۱}{۴۰}$ | $\frac{۱}{۶۰}$ | $\frac{۱}{۴۵}$ | $\frac{۲}{۹۰}$ |

الف) برآورد ناریبی برای میانگین جامعه بیابید.

ب) برآورد ناریبی برای واریانس برآوردکننده میانگین جامعه به دست آورید.

۳. جامعه‌ای مرکب از ۱۲ واحد است می‌خواهیم نمونه‌ای به حجم ۴، با احتمال متغیر، و به روش تصادفی با جایگذاری از این جامعه انتخاب کنیم. مقدار و احتمالهای متناظر با واحدها به صورت زیرند:

| | | | | | | | | | | | | |
|-------|----------------|----------------|---------------|----------------|---------------|----------------|---------------|----------------|----------------|----------------|----------------|----------------|
| Y_i | Y_1 | Y_2 | Y_3 | Y_4 | Y_5 | Y_6 | Y_7 | Y_8 | Y_9 | Y_{10} | Y_{11} | Y_{12} |
| P_i | $\frac{۱}{۱۲}$ | $\frac{۱}{۱۸}$ | $\frac{۱}{۹}$ | $\frac{۱}{۱۲}$ | $\frac{۱}{۶}$ | $\frac{۱}{۱۸}$ | $\frac{۱}{۹}$ | $\frac{۱}{۱۲}$ | $\frac{۱}{۱۵}$ | $\frac{۱}{۱۸}$ | $\frac{۱}{۱۲}$ | $\frac{۲}{۴۵}$ |

نمونه را یک بار با روش لاهییری و یک بار با روش مجموع تراکمی به دست آورید.

۴. فرض کنید از جامعه بالا نمونه‌ای به حجم ۴، بدون جایگذاری و با احتمال متغیر انتخاب کرده‌ایم. مشخصه ۴ واحد منتخب را اندازه‌گیری کرده و با احتمالهای متناظر در جدول زیر آورده‌ایم

| | | | | |
|-------|----------------|---------------|---------------|----------------|
| Y_i | ۱۰ | ۱۱ | ۸ | ۱۳ |
| P_i | $\frac{۲}{۴۵}$ | $\frac{۱}{۹}$ | $\frac{۱}{۶}$ | $\frac{۱}{۱۲}$ |

الف) میانگین جامعه را برآورد کنید.

ب) برآورد واریانس برآوردکننده میانگین جامعه را به دست آورید.

۵. از جامعه مذکور در تمرین ۳، نمونه‌ای به حجم ۴، بدون جایگذاری، انتخاب کنید و میانگین جامعه را برآورد نمایید.

۶. جامعه‌ای متشکل از ۱۲ واحد است. در زیر مقادیر صفت‌های کمکی متناظر با آنها را آورده‌ایم

| | | | | | | | | | | | | |
|------------------|---|---|---|---|---|----|----|----|----|----|----|----|
| i : شماره واحد | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ | ۱۰ | ۱۱ | ۱۲ |
| X_i : صفت کمکی | ۲ | ۵ | ۷ | ۳ | ۱ | ۱۲ | ۱۸ | ۲۳ | ۳۷ | ۴ | ۳۱ | ۱۹ |

مجموع واحدهای جامعه را بر مبنای نمونه‌ای به حجم ۴ به شیوه RHC برآورد کنید.

تمرینهای چهارگزینه‌ای

۱. در جامعه‌ای به حجم ۱۰۰ اگر بزرگترین واحد جامعه کمکی ۶۰۰ باشد و بخواهیم با روش

۱۳۴ نمونه‌گیری با احتمال متغیر

لاهیری نمونه‌ای با جایگذاری و با احتمال متغیر انتخاب کنیم، احتمال مؤثر بودن انتخاب زوج $(j, 5)$ برابر است با

الف) $\frac{X_5}{600}$ ب) $\frac{X_5}{100}$ ج) $\frac{X_5}{60000}$ د) $\frac{X_j}{5}$

۲. در یک نمونه‌گیری تصادفی با احتمال متغیر، میانگین صفت کمکی ۴۸ است. اگر بزرگترین واحد صفت کمکی ۶۰ باشد و نمونه‌گیری با جایگذاری و با روش لاهیری انجام شود، احتمال نامؤثر بودن هر انتخاب برابر است با

الف) $\frac{4}{5}$ ب) $\frac{1}{48}$ ج) $\frac{1}{60}$ د) $\frac{1}{5}$

۳. در یک نمونه‌گیری تصادفی با احتمال متغیر و با جایگذاری نمونه‌ای به حجم ۵ گرفته‌ایم. اگر حجم جامعه ۵۰ و مجموع $\frac{Y_i}{p_i}$ های نمونه ۵۰۰ باشد برآورد ناریب میانگین جامعه برابر است با

الف) ۲ ب) ۱۰ ج) $\frac{1}{4}$ د) ۱۰۰

۴. در جامعه‌ای به حجم ۱۰ نمونه‌ای به حجم ۲ با احتمال متغیر و بدون جایگذاری گرفته‌ایم. اندازه این ۲ واحد به ترتیب ۱۲ و ۱۶ و احتمالهای متناظر با آنها به ترتیب $\frac{1}{8}$ و $\frac{1}{10}$ است. در این صورت برآورد میانگین جامعه برابر است با

الف) ۱۵ ب) ۱۵٫۵ ج) ۱۴ د) ۱۲٫۴

۵. اگر نمونه‌ای به روش تصادفی با احتمال متغیر و با جایگذاری از جامعه‌ای به حجم ۵۰ بگیریم و اگر در جامعه $\sum \frac{Y_i^2}{X_i} = 2$ و $\sum Y_i^2 = 100$ و $\sum X_i = 4000$ ، آنگاه دقت نمونه‌گیری مزبور از نمونه‌گیری تصادفی ساده با جایگذاری

الف) بیشتر است ب) کمتر است ج) دقتها یکسان‌اند د) نمی‌توان قضاوت کرد.

نمونه‌گیری تصادفی با طبقه‌بندی

۰.۴ مقدمه

(در نمونه‌گیری تصادفی ساده، هر واحد بانس برآورد میانگین جامعه علاوه بر اینکه به حجم نمونه بستگی دارد به تغییرپذیری مشخصه تحت بررسی هم وابسته است.) اگر جامعه خیلی ناهمگن باشد و به دلیل محدودیت‌های اقتصادی نتوانیم حجم نمونه را بزرگ آخبار کنیم تقریباً غیرممکن است که با روش نمونه‌گیری تصادفی ساده، برآوردی به حد کافی دقیق برای پارامتر مورد نظر جامعه بیابیم. مثلاً در بررسی‌های مربوط به مشخصه‌های کارکنان شرکتها در شهری بزرگ، تعداد کارکنان در برخی از شرکتها بیش از ۱۰۰۰ نفر و در بعضی محدود به ده نفر یا کمتر است. در این صورت انتخاب نمونه تصادفی ساده از جامعه شرکتها، با توجه به ناهمگنی تعداد کارکنان، از نمونه‌ای به نمونه دیگر با نوسانات شدید همراه است. اما اگر ممکن باشد که شرکتها را از نظر تعداد کارکنان به چند طبقه با حجم خیلی کم، کم، متوسط، زیاد، و خیلی زیاد تقسیم کنیم و از هر طبقه برای بررسی مشخصه مورد نظر به انتخاب نمونه بپردازیم از نوسانات شدید در نمونه‌های مختلف جلوگیری خواهد شد. این رهیافت، نمونه‌گیری دیگری را تحت عنوان نمونه‌گیری با طبقه‌بندی مطرح می‌کند. این روش، تکنیکی بسیار متداول است که به دلایل زیاد انجام می‌شود. عمده‌ترین این دلایل به شرح زیرند:

۱. اگر برای بعضی از زیرجامعه‌های یک جامعه، داده‌ها و اطلاعاتی با دقت معلوم خواهند لازم است که هر زیرجامعه را یک طبقه به حساب آورند.
۲. تشکیلاتی که در یک کشور، مسئول انجام نمونه‌گیری برای ارائه نتایج به سازمانهای ذی‌ربط

است، در هر یک از نواحی مختلف کشور واحدهایی دارد. کارکنان هر واحد دربارهٔ ویژگیهای مورد نظر در ناحیهٔ خود، اطلاعاتی دقیقتر از سایرین دارند و لذا اگر نمونه‌گیری در هر ناحیه، به‌عنوان یک طبقه، به‌صورتی مستقل از نواحی دیگر صورت گیرد با دقتی بیشتر همراه بوده، به‌علاوه از لحاظ میزان هزینه و نحوهٔ سازماندهی کار نمونه‌گیری، تسهیلاتی بیشتر فراهم می‌شوند. همگنی تقریبی بعضی از صفات تحت نمونه‌گیری در یک ناحیه نیز، به گونه‌ای که خواهیم دید به بالا بردن کارایی نمونه‌گیری با طبقه‌بندی کمک عمده‌ای خواهد کرد. مشکلات نمونه‌گیری در بخشهای مختلف یک جامعهٔ بزرگ، به‌صورتی بارز، متفاوت‌اند، که لاجرم نمی‌توان رفتار و سیاستی یکسان در همهٔ جامعه داشت و طبیعتاً تقسیم جامعه به طبقه‌ها امری منطقی است.

۳. با طبقه‌بندی کردن جامعه می‌توان دقت برآورد مجموع جامعه را کنترل کرد. زیرا می‌توان یک جامعهٔ ناهمگن را به زیرجامعه‌های تقریباً همگن تقسیم کرد. منظور از زیرجامعه یا طبقهٔ همگن، طبقه‌ای است که اندازه‌ها از واحدی به واحد دیگر تغییر عمده‌ای ندارند و می‌توان در چنین طبقه‌ای با نمونه‌ای به حجم اندک، برآوردی دقیق از صفت تحت بررسی به‌دست آورد. برآوردهایی که جداگانه از طبقه‌های مختلف فراهم می‌شوند سرانجام ترکیب‌شده و برآوردی دقیق برای صفت مورد نظر در کل جامعه به‌دست می‌آید. چون انتخاب نمونه‌ها از طبقه‌های مختلف مستقل از هم انجام می‌شود واریانس هر برآوردکننده در طبقه‌ها با هم جمع می‌شوند تا واریانس برآوردکننده در کل جامعه با احتساب ضرایبی که بعداً از آن گفتگو می‌کنیم به‌دست آید. لذا اصل طبقه‌بندی مبتنی بر آن چنان افزایی از جامعه است که واریانس برآوردکننده در هر طبقه تا حد ممکن کوچک باشد. برای این هدف، لاجرم باید در هر طبقه واحدهایی همگن قرار گیرند.

در این مقدمه لازم است تذکر دهیم که در بعضی از جامعه‌ها تقسیم جامعه به طبقات را باید با توجه به ویژگیهای قسمتهای مختلف جامعه انجام داد ولی در اکثر جامعه‌ها طبقه‌ها به‌صورت طبیعی و با ساختار جامعه از قبل مشخص شده‌اند. مثالهایی نوعی از جامعه‌هایی که به‌طور طبیعی طبقه‌بندی شده‌اند به‌شرح زیرند:

الف) جامعهٔ دانش‌آموزان مدارس خودبه‌خود به‌وسیلهٔ کلاس‌بندی مدارس و جنس دانش‌آموز به‌صورت طبقه‌ها درمی‌آید. یک طبقه در این جامعه، مثلاً «کلاس چهارم دختران» است.

ب) جامعهٔ افراد مالیات‌دهنده، به‌وسیلهٔ شهر، جنس، و دامنهٔ درآمد گزارش‌شده طبقه‌بندی می‌شود. تعیین دامنهٔ درآمدها دلخواه است. یک طبقه در این جامعه، مثلاً «مردان شهر تهران با درآمد سالیانه‌ای در دامنهٔ [دومیلیون-یک میلیون] تومان» است.

ج) برای جامعهٔ خانوارها در یک کشور، هر استان کشور یک طبقه است. بدیهی است مرز هر استان باید دقیقاً تعریف شود.

د) جامعهٔ مؤسسات خرده‌فروشی به‌وسیلهٔ نوع جنسی که عرضه می‌شود (ادویه، گوشت، ...)

و مکان مؤسسه، خودبه‌خود طبقه‌بندی می‌شود. بدیهی است بعضی از ویژگیهایی که طبقه‌ها را به وجود می‌آورند جنبهٔ کیفی و برخی جنبهٔ

کمی دارند. مثلاً جنس دانش آموز، یا نام خرده‌فروشی جنبه کیفی و میزان درآمد جنبه کمی دارد. معمولاً اگر یک جامعه به صورت طبیعی طبقه‌بندی نشده باشد و بخواهیم خود ما به کار طبقه‌بندی بپردازیم باید ابتدا نمونه‌ای نسبتاً بزرگ به روش تصادفی ساده از جامعه بگیریم و سپس با در نظر گرفتن ناهمگنی و پراکندگی این نمونه، طبقاتی را با کرانه‌های مناسب انتخاب کنیم، به طوری که اگر واحدهای این نمونه را در طبقه‌های ساخته شده توزیع کنیم در هر طبقه همگنی موجود باشد. بعداً درباره تعیین کرانه‌های طبقات گفتگو خواهیم کرد.

۱.۴ تعریف نمونه تصادفی با طبقه‌بندی

در نمونه‌گیری با طبقه‌بندی، همان‌طور که توضیح دادیم، ابتدا جامعه به حجم N را به L زیرجامعه به حجم‌های $N_1, N_2, \dots, N_h, \dots, N_L$ تقسیم می‌کنیم، به قسمی که این زیرجامعه‌ها متداخل نباشند و هر واحد جامعه به یک و تنها به یک زیرجامعه متعلق باشد. در واقع زیرجامعه‌ها افزایی از جامعه را به وجود می‌آورند. پس

$$N = \sum_{h=1}^L N_h$$

جامعه ←

هر زیرجامعه را یک طبقه می‌نامیم. این طبقه‌ها همان‌گونه که در بالا اشاره کردیم به‌گونه‌ای تعیین می‌شوند که واحدهای با حداکثر همگنی در هر طبقه قرار بگیرند. برای به دست آوردن برآوردهای دقیق باید مقادیر N_h ($h = 1, 2, \dots, L$) را بدانیم. وقتی طبقه‌ها و حجم آنها مشخص شد از هر طبقه نمونه‌ای انتخاب می‌شود. این نمونه را می‌توان با هر روشی انتخاب کرد. بعداً که با نمونه‌گیریهای سیستماتیک، خوشه‌ای، برآورد رگرسیونی، برآورد نسبتی آشنا شدید، می‌توانید در هر طبقه نمونه را با یکی از این روشها هم به دست آورید. انتخاب نمونه در هر طبقه مستقل از سایر طبقه‌هاست. حجمهای نمونه‌ها در طبقه‌ها را به ترتیب با $n_1, n_2, \dots, n_h, \dots, n_L$ نشان می‌دهیم به طوری که

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

معرف حجم نمونه از کل جامعه است. اگر نمونه‌ای که از هر طبقه انتخاب می‌شود، به روش نمونه‌گیری تصادفی ساده بدون جایگذاری باشد، نمونه‌گیری را نمونه‌گیری تصادفی با طبقه‌بندی می‌نامند و نمونه متشکل از کل نمونه‌های طبقات را نمونه تصادفی با طبقه‌بندی می‌گویند. در این فصل، این نوع نمونه‌گیری را بررسی می‌کنیم.

نظریه نمونه‌گیری با طبقه‌بندی، از ویژگیهای برآوردهای نمونه با طبقه‌بندی و از انتخاب بهترین حجم نمونه در هر طبقه برای تأمین دقت ماکسیمم بحث می‌کند. ابتدا فرض می‌کنیم طبقه‌ها پیشاپیش ساخته شده‌اند، و مطالب مربوط به برآوردها را مطرح می‌کنیم، سپس در گام بعد از چگونگی ساختن طبقه‌ها گفتگو می‌نماییم.

۲.۴ نمادها و برخی از تعریفها

تعداد طبقه‌ها را با L نشان می‌دهیم. زیرنویس h ، معرف شماره طبقه است که از ۱ تا L تغییر می‌کند. زیرنویس i برای تعیین شماره واحدها در داخل هر طبقه به‌کار می‌رود. N تعداد کل افراد جامعه است. نمادهای زیر برای طبقه h ، $h = 1, 2, \dots, L$ تعریف می‌شوند

| | |
|---|---|
| \bar{Y}_N | میانگین جامعه |
| N_h | تعداد کل واحدها در طبقه h ام |
| n_h | تعداد واحدهای نمونه از طبقه h ام |
| $Y_{hi}, \quad i = 1, 2, \dots, n_h$ | مقدار صفت واحد i ام در طبقه h ام |
| $W_h = \frac{N_h}{N}$ | نسبت تعداد واحدهای طبقه h ام به تعداد واحدهای کل جامعه یا وزن طبقه h ام |
| $f_h = \frac{n_h}{N_h}$ | کسر نمونه‌گیری برای طبقه h ام |
| $\bar{y}_h = \frac{1}{N_h} \sum_{i=1}^{N_h} Y_{hi}$ | میانگین طبقه h ام |
| $\bar{Y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} Y_{hi}$ | میانگین نمونه طبقه h ام |
| $S_h^2 = \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (Y_{hi} - \bar{y}_h)^2$ | تغییرات طبقه h ام |
| $s_h^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (Y_{hi} - \bar{Y}_h)^2$ | تغییرات نمونه‌ای طبقه h ام |

ار واریانس \leftarrow واریانس نمونه

گاهی اوقات S_h^2 و s_h^2 را به ترتیب واریانس طبقه h و واریانس نمونه این طبقه می‌نامند. توجه کنید که واحدهای هر طبقه و نمونه آن را یکسان نشان داده‌ایم. خواننده از محتوای مطلب در هر جا متوجه خواهد شد که واحد مورد بحث به طبقه تعلق دارد یا به نمونه طبقه. اگر از L طبقه که به حجمهای N_1, N_2, \dots, N_L هستند، L نمونه به ترتیب به حجمهای n_1, n_2, \dots, n_L به روش تصادفی ساده اختیار کنیم، دنباله‌ای با $\sum_{h=1}^L n_h = n$ عضو خواهیم داشت که همان نمونه تصادفی با طبقه‌بندی است. میانگین موزون میانگینهای نمونه‌ای طبقه‌ها را میانگین نمونه با طبقه‌بندی می‌نامیم و آن را با \bar{Y}_{st} نمایش می‌دهیم

$$\boxed{\bar{Y}_{st}} = \frac{1}{N} \sum_{h=1}^L N_h \bar{Y}_h = \sum_{h=1}^L W_h \bar{Y}_h \quad (1.4)$$

توجه کنید که این نمونه، با نمونه تصادفی به حجم n از کل جامعه تفاوت دارد.

قضیه ۱.۴ اگر در طبقه h , $h = 1, \dots, L$ برآوردکننده ناریب میانگین طبقه h ام باشد، آنگاه \bar{Y}_h برآوردکننده ناریب میانگین کل جامعه، یعنی \bar{Y}_N است.

برهان. با توجه به رابطه (۱.۴)، داریم

$$E(\bar{Y}_{st}) = E \left[\sum_{h=1}^L W_h \bar{Y}_h \right] = \sum_{h=1}^L W_h E(\bar{Y}_h)$$

اما چون بنا بر فرض قضیه، \bar{Y}_h برآوردکننده ناریب \bar{Y}_h است، پس

$$E(\bar{Y}_{st}) = \sum_{h=1}^L W_h \bar{Y}_h$$

از طرفی $\bar{Y}_h = \frac{1}{N_h} \sum_{i=1}^{N_h} Y_{hi}$ پس

$$E(\bar{Y}_{st}) = \sum_{h=1}^L \left[\frac{W_h}{N_h} \sum_{i=1}^{N_h} Y_{hi} \right]$$

اما $W_h = \frac{N_h}{N}$ و لذا

$$E(\bar{Y}_{st}) = \frac{1}{N} \sum_{h=1}^L \sum_{i=1}^{N_h} Y_{hi}$$

ضریب $\frac{1}{N}$ در عبارت طرف دوم برابر مجموع همه واحدهای جامعه است، پس

$$E(\bar{Y}_{st}) = \bar{Y}_N \quad (۲.۴)$$

لذا $\hat{Y}_N = \bar{Y}_{st}$ برآوردکننده ای ناریب است. لازم است تذکر دهیم که برای هر نوع نمونه‌گیری در داخل طبقه‌ها، وقتی \bar{Y}_h برآوردکننده ناریب میانگین طبقه h , $h = 1, 2, \dots, L$ باشد، حکم قضیه بالا صحیح است. \square

(فرع) اگر T_N مجموع واحدهای جامعه باشد، آنگاه $N\bar{Y}_{st}$ برآوردکننده ناریب T_N است، زیرا از رابطه $E(\bar{Y}_{st}) = \bar{Y}_N$ نتیجه می‌شود که $E(N\bar{Y}_{st}) = N\bar{Y}_N$ و چون $T_N = N\bar{Y}_N$ پس $E(N\bar{Y}_{st}) = T_N$ یعنی $N\bar{Y}_{st}$ برآوردکننده ناریب T_N است. \square

قضیه ۲.۴ اگر نمونه‌های طبقه‌ها مستقل از هم باشند

$$V(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 V(\bar{Y}_h) \quad (۳.۲)$$

برهان. \bar{Y}_{st} تابعی خطی از \bar{Y}_h ها با ضریبهای ثابت W_h , $h = 1, \dots, L$ است، لذا

$$V(\bar{Y}_{st}) = V(W_1\bar{Y}_1 + W_2\bar{Y}_2 + \dots + W_L\bar{Y}_L)$$

چون نمونه‌ها در طبقه‌ها مستقل از یکدیگرند، \bar{Y}_h ها از هم مستقل‌اند و بنابراین

$$V(\bar{Y}_{st}) = \sum_{h=1}^L V(W_h\bar{Y}_h) = \sum_{h=1}^L W_h^2 V(\bar{Y}_h)$$

□

رابطه بالا نشان می‌دهد که $V(\bar{Y}_{st})$ به \bar{Y}_h ها و وزنهای بستگی دارد. اگر واحدهای هر طبقه همگن باشند واریانس طبقه‌ها کوچک هستند. به عبارت دیگر اگر بتوان جامعه با پراکندگی زیاد را به طبقاتی تقسیم کنیم که واحدهای هر طبقه نزدیک به هم باشند \bar{Y}_N تقریباً با خطای اندک برآورد می‌شود. در واقع استفاده از وزنهای طبقاتی $\frac{N_h}{N}$ است که موجب می‌شود \bar{Y}_{st} تقریباً با خطای اندک \bar{Y}_N را برآورد کند. رابطه (۳.۴) به نوع نمونه‌گیری طبقات بستگی ندارد.

قضیه ۳.۴ در نمونه‌گیری تصادفی با طبقه‌بندی، واریانس \bar{Y}_{st} به صورت زیر است

$$\begin{aligned} V(\bar{Y}_{st}) &= \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{S_h^2}{n_h} = \sum_{h=1}^L W_h^2 (1 - f_h) \frac{S_h^2}{n_h} \\ &= \sum_{h=1}^L \frac{W_h^2 S_h^2}{n_h} - \sum_{h=1}^L \frac{W_h S_h^2}{N} \end{aligned} \quad (4.4)$$

برهان. چون در طبقه h , $h = 1, \dots, L$ ، نمونه‌گیری به روش تصادفی ساده بدون جایگذاری انجام می‌شود، بنابراین

$$V(\bar{Y}_h) = \left(\frac{1}{n_h} - \frac{1}{N_h} \right) S_h^2 = \frac{N_h - n_h}{N_h} \cdot \frac{S_h^2}{n_h}$$

اینک با توجه به (۳.۴)

$$\begin{aligned} V(\bar{Y}_{st}) &= \sum_{h=1}^L W_h^2 V(\bar{Y}_h) = \sum_{h=1}^L W_h^2 \frac{N_h - n_h}{N_h} \frac{S_h^2}{n_h} \\ &= \sum_{h=1}^L \frac{N_h^2}{N^2} \cdot \frac{1}{N_h} (N_h - n_h) \frac{S_h^2}{n_h} \\ &= \frac{1}{N^2} \sum_{h=1}^L N_h (N_h - n_h) \frac{S_h^2}{n_h} \end{aligned}$$

که درستی برابری اول را در حکم قضیه نشان می‌دهد. از طرفی

$$V(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 \frac{N_h - n_h}{N_h} \cdot \frac{S_h^2}{n_h} = \sum_{h=1}^L W_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_h^2}{n_h}$$

چون $\frac{n_h}{N_h} = f_h$ ، پس

$$V(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 (1 - f_h) \frac{S_h^2}{n_h}$$

که درستی برابری دوم را در حکم قضیه نشان می‌دهد. (درستی آخرین برابری را در حکم قضیه تحقیق کنید.)

فرع ۱. اگر در همه طبقه‌ها f_h ها کوچک و قابل اغماض باشند، آنگاه

$$V(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 \frac{S_h^2}{n_h} = \frac{1}{N^2} \sum_{h=1}^L N_h^2 \frac{S_h^2}{n_h} \quad (5.4)$$

فرع ۲. اگر T_N مجموع واحدهای جامعه باشد، آنگاه همان طور که دیدیم $N\bar{Y}_{st}$ برآوردکننده ناریب T_N است و

$$\begin{aligned} V(\hat{T}_N) &= V(N\bar{Y}_{st}) = N^2 V(\bar{Y}_{st}) = N^2 \cdot \frac{1}{N^2} \sum_{h=1}^L N_h (N_h - n_h) \frac{S_h^2}{n_h} \\ &= \sum_{h=1}^L N_h (N_h - n_h) \frac{S_h^2}{n_h} \end{aligned} \quad (6.4)$$

□

۳.۴ نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب

اگر در نمونه‌گیری با طبقه‌بندی حجم نمونه‌های طبقه‌ها با حجم طبقه‌ها متناسب باشند، نمونه‌گیری را با تخصیص متناسب می‌نامند. در واقع در این نوع نمونه‌گیری

$$\frac{n_h}{n} = \frac{N_h}{N} = W_h \quad h = 1, \dots, L$$

یا

$$\frac{n_h}{N_h} = \frac{n}{N} \quad h = 1, \dots, L$$

که نتیجه می‌دهد به‌ازای $L, 1, 2, \dots, L$ ، $f_h = f$ ، به عبارت دیگر در این نوع نمونه‌گیری با طبقه‌بندی، کسر نمونه‌گیری در همه طبقه‌ها یکسان و برابر با $\frac{n}{N}$ است.

قضیه ۴.۴ در نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب

$$V(\bar{Y}_{st}) = \frac{1-f}{n} \sum_{h=1}^L W_h S_h^2 \quad (7.4)$$

برهان. با توجه به (۴.۴)، رابطه کلی

$$V(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 (1-f_h) \frac{S_h^2}{n_h}$$

را در نظر می‌گیریم. چون نمونه‌گیری با تخصیص متناسب است، همه f_h ها برابر f هستند، پس

$$V(\bar{Y}_{st}) = \sum_{h=1}^L \frac{N_h}{N} \cdot W_h (1-f) \frac{S_h^2}{n_h}$$

چون، $\frac{N_h}{N} = \frac{n_h}{n}$ داریم

$$\begin{aligned} V(\bar{Y}_{st}) &= (1-f) \sum_{h=1}^L W_h \cdot \frac{n_h}{n} \cdot \frac{S_h^2}{n_h} \\ &= \frac{1-f}{n} \sum_{h=1}^L W_h S_h^2 \end{aligned}$$

□

فرع. وقتی نمونه‌گیری با تخصیص متناسب بوده و مقدار واریانس در همه طبقه‌ها یکسان باشد، آنگاه اگر این واریانس مشترک را با S_w^2 نشان دهیم همه S_h^2 ها برابر با S_w^2 هستند. پس برابری (۷.۴) به صورت زیر درمی‌آید

$$V(\bar{Y}_{st}) = \frac{1-f}{n} \sum_{h=1}^L W_h S_h^2 = \frac{1-f}{n} S_w^2 \sum_{h=1}^L W_h = (1-f) \frac{S_w^2}{n} \quad (8.4)$$

که در آن، f ، کسر نمونه‌گیری در کل جامعه است.

مثال ۱.۴ جدول زیر معرف تعداد ساکنین ۶۴ شهر بزرگ در سال ۱۹۳۰ برحسب هزار نفر است. در این جدول، جمعیت ۶۴ شهر در سال ۱۹۳۰ و در دو طبقه ثبت شده است. طبقه اول شامل ۱۶ شهر با جمعیت بیشتر و طبقه دوم شامل ۴۸ شهر با جمعیت کمتر است.

نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب ۴۳

جمعیت ۶۴ شهر (برحسب ۱۰۰۰۰) در سال ۱۹۳۰*

Y_{hi} ها

$h :$

| | ۱ | ۲ | |
|------|-----|-----|-----|
| ۹۰۰ | ۳۶۴ | ۲۰۹ | ۱۱۳ |
| ۸۲۲ | ۳۱۷ | ۱۸۳ | ۱۱۵ |
| ۷۸۱ | ۳۲۸ | ۱۶۳ | ۱۲۳ |
| ۸۰۵ | ۳۰۲ | ۲۵۳ | ۱۵۴ |
| ۶۷۰ | ۲۸۸ | ۲۳۲ | ۱۴۰ |
| ۱۲۳۸ | ۲۹۱ | ۲۶۰ | ۱۱۹ |
| ۵۷۳ | ۲۵۳ | ۲۰۱ | ۱۳۰ |
| ۶۳۴ | ۲۹۱ | ۱۴۷ | ۱۲۷ |
| ۵۷۸ | ۳۰۸ | ۳۹۲ | ۱۰۰ |
| ۴۸۷ | ۲۷۲ | ۱۶۴ | ۱۰۷ |
| ۴۴۲ | ۲۸۴ | ۱۴۳ | ۱۱۴ |
| ۴۵۱ | ۲۵۵ | ۱۶۹ | ۱۱۱ |
| ۴۵۹ | ۲۷۰ | ۱۳۹ | ۱۶۳ |
| ۴۶۴ | ۲۱۴ | ۱۷۰ | ۱۱۶ |
| ۴۰۰ | ۱۹۵ | ۱۵۰ | ۱۲۲ |
| ۳۶۶ | ۲۶۰ | ۱۴۳ | ۱۳۴ |

* داده‌ها از تکنیکهای نمونه‌گیری تألیف کوکران اقتباس شده‌اند.

قرار است کل جمعیت ۶۴ شهر در سال ۱۹۳۰ را از روی نمونه‌ای به حجم ۲۴ برآورد کنند. انحراف معیار برآوردکننده کل جمعیت را (۱) به وسیله نمونه تصادفی ساده، (۲) به وسیله نمونه تصادفی با طبقه‌بندی و با تخصیص متناسب، (۳) به وسیله نمونه تصادفی با طبقه‌بندی و استخراج ۱۲ واحد از هر طبقه بیابید.

این جامعه به بسیاری از جامعه‌های صنعتی شباهت دارد که بعضی واحدها، مثلاً شهرهای بزرگ، سهمی عمده از کل را دارند و در تغییرپذیری نقشی مهمتر از بقیه واحدها بازی می‌کنند. در جدول زیر، جدول داده‌ها، مجموعهای طبقاتی، و مجموعهای مربعات ارائه شده‌اند. در این مثال فقط از داده‌های سال ۱۹۳۰ استفاده کرده‌ایم.

مجموع و مجموع مربعات

| طبقه | $\sum Y_{hi}$ | $\sum (Y_{hi}^2)$ |
|------|---------------|-------------------|
| ۱ | ۱۰۰۷۰ | ۷۱۴۵۴۵۰ |
| ۲ | ۹۴۹۸ | ۲۱۴۱۷۲۰ |

برای کل جامعه در سال ۱۹۳۰، به دست می‌آوریم

$$T_N = ۱۹۵۶۸ \quad S^2 = ۵۲۴۴۸$$

برای ۳ موردی که در صورت مثال مطرح شده‌اند برآورد کل جامعه را با نمادهای \hat{T}_{prop} ، \hat{T}_{ran} و \hat{T}_{equal} نشان می‌دهیم. طبق فرض $N_1 = ۱۶$ ، $N_2 = ۴۸$.
 ۱. برای نمونه‌گیری تصادفی ساده

$$V(\hat{T}_{ran}) = N^2 \left(1 - \frac{n}{N}\right) \frac{S^2}{n} = (۶۴)^2 \left(1 - \frac{۲۴}{۶۴}\right) \frac{۵۲۴۴۸}{۲۴} = ۵۵۹۴۴۵۳$$

و لذا $\sigma(\hat{T}_{ran}) = ۲۳۶۵$.
 ۲. برای دو طبقه به ترتیب

$$S_1^2 = ۵۳۸۴۳ \quad S_2^2 = ۵۵۸۱$$

توجه کنید واریانس مربوط به طبقه ۱ که شامل شهرهای پرجمعیت است تقریباً ۱۰ برابر واریانس طبقه ۲ است.

در تخصیص متناسب، داریم $n_1 = ۶$ و $n_2 = ۱۸$. بنابر (۷.۴)

$$V(\bar{Y}_{st}) = \frac{1-f}{n} \sum_{h=1}^L W_h S_h^2 = \frac{N-n}{nN^2} \sum_{h=1}^L N_h S_h^2$$

از ضرب طرفین در N^2 ، داریم

$$\begin{aligned} V(N\bar{Y}_{st}) &= V(\hat{T}_{prop}) = \frac{N-n}{n} \sum_{h=1}^L N_h S_h^2 \\ &= \frac{۶۴-۲۴}{۲۴} [(۱۶)(۵۳۸۴۳) + (۴۸)(۵۵۸۱)] = ۱۸۸۲۲۹۳ \end{aligned}$$

و لذا

$$\sigma(\hat{T}_{prop}) = ۱۳۷۲$$

۳. برای $n_1 = n_2 = ۱۲$ ، فرمول کلی

$$V(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h (N_h - n_h) \frac{S_h^2}{n_h}$$

را به‌کار می‌بریم. از ضرب طرفین در N^2 داریم

$$V(N\bar{Y}_{st}) = V(\hat{T}_{equal}) = \sum_{h=1}^L N_h(N_h - n_h) \frac{S_h^2}{n_h}$$

پس

$$V(\hat{T}_{equal}) = \frac{16(16 - 12)(52842)}{12} + \frac{(28)(28 - 12)(5581)}{12}$$

$$= 1090827$$

$$\sigma(\hat{T}_{equal}) = 1044$$

در این مثال، نمونه‌گیری با حجمهای نمونه‌ای برابر در دو طبقه دقیقتر از تخصیص متناسب است، ولی هر دو از نمونه‌گیری تصادفی ساده بهترند. ▲

تبصره. احتمال اینکه واحدی متعلق به طبقه h به‌عنوان عضوی از نمونه با طبقه‌بندی انتخاب شود برابر با $\frac{n_h}{N_h}$ است. احتمال اینکه دو واحد متعلق به طبقه h به‌عنوان اعضای نمونه با طبقه‌بندی انتخاب شوند $n_h(n_h - 1)/N_h(N_h - 1)$ است. اگر دو واحد به دو طبقه مختلف مثلاً طبقه‌های h و k متعلق باشند احتمال اینکه به‌عنوان اعضای نمونه با طبقه‌بندی انتخاب شوند برابر با $n_h \cdot n_k / N_h \cdot N_k$ است. تمام این احتمالها با احتمال انتخاب هر واحد در نمونه‌گیری تصادفی ساده متفاوت‌اند.

۴.۴ برآورد واریانس برآوردکننده میانگین در نمونه‌گیری با طبقه‌بندی

در نمونه‌گیری تصادفی ساده بدون جایگذاری دیدیم که S^2 ی نمونه برآوردکننده ناریب S^2 ی جامعه است. لذا S_h^2 یعنی واریانس نمونه‌ای در طبقه h ام برآوردکننده ناریب S_h^2 ، یعنی واریانس طبقه h ام است. با توجه به این مطلب قضیه زیر را بیان می‌کنیم.

قضیه ۵.۲ در نمونه‌گیری تصادفی ساده با طبقه‌بندی، برآوردکننده واریانس برای برآوردکننده میانگین جامعه برابر است با

$$V(\bar{Y}_{st}) = \hat{\sigma}^2(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{S_h^2}{n_h} \quad (9.4)$$

برهان. برابر (۲.۴) داریم

$$V(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{S_h^2}{n_h}$$

با توجه به برابری

$$E(s_h^2) = S_h^2 \quad h = 1, 2, \dots, L$$

رابطه بالا به صورت زیر درمی‌آید

$$V(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{E(s_h^2)}{n_h}$$

چون، برای مقدار ثابت a ، $E(aX) = aE(X)$ که در آن X ، متغیری تصادفی است، پس رابطه بالا را می‌توان چنین نوشت

$$V(\bar{Y}_{st}) = E \left[\frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{s_h^2}{n_h} \right]$$

این برابری، معرف آن است که عبارت داخل کروشه برآوردکننده نااریب $V(\bar{Y}_{st})$ است، یعنی

$$\hat{V}(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{s_h^2}{n_h}$$

و یا

$$\hat{V}(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 (1 - f_h) \frac{s_h^2}{n_h} \quad (10.4)$$

در عمل، با داشتن تنها یک نمونه در هر طبقه، s_h^2 ها به صورت مقادیری عددی هستند و لذا برآوردی نااریب، برای واریانس \bar{Y}_{st} به دست می‌آید. صورت دیگر (۱۰.۴)، با توجه به آخرین برابری (۴.۴) به صورت زیر است

$$\hat{V}(\bar{Y}_{st}) = \sum_{h=1}^L \frac{W_h^2 s_h^2}{n_h} - \sum_{h=1}^L \frac{W_h s_h^2}{N} \quad (11.4)$$

برای محاسبه $\hat{V}(\bar{Y}_{st})$ باید از هر طبقه حداقل دو واحد انتخاب شود، زیرا در غیر این صورت محاسبه s_h^2 به دلیل صفر شدن $n_h - 1$ در رابطه

$$s_h^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (Y_{hi} - \bar{Y}_h)^2$$

میسر نیست. برای حالت نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب، از (۷.۴) داریم

$$\hat{V}(\bar{Y}_{st}) = \frac{1-f}{n} \sum_{h=1}^L W_h s_h^2$$

که برآوردکننده‌ای نااریب است.

۵.۴ حدود اطمینان میانگین و مجموع واحدهای جامعه در نمونه‌گیری با طبقه‌بندی

برای یک نمونه تصادفی ساده با طبقه‌بندی به حجم n از جامعه‌ای به حجم N ، تعداد نمونه‌های ممکن $\prod_{h=1}^L \binom{N_h}{n_h}$ است. مثلاً اگر ۴ طبقه به حجمهای ۱۰، ۸، ۹، ۱۲ داشته باشیم و از این طبقه‌ها به ترتیب ۴، ۲، ۳، ۵ واحد به تصادف انتخاب کنیم تعداد نمونه‌های ممکن برابر است با

$$\binom{10}{4} \binom{8}{2} \binom{9}{3} \binom{12}{5} = (120)(28)(84)(792) = 39118464$$

بنابراین در حالت کلی حداکثر به همین تعداد هم \bar{Y}_{st} خواهیم داشت. اگر تعداد واحدهایی که \bar{Y}_{st} را تولید می‌کنند یعنی n بزرگ باشد، مثلاً بیش از ۲۰، فرض نرمال بودن توزیع \bar{Y}_{st} را به شرط متناهی بودن واریانسهای طبقات، و با توجه به قضیه حد مرکزی می‌توان پذیرفت. با توجه به این مقدمه، برای تهیه بازه اطمینان، ابتدا توزیع \bar{Y}_{st} را نرمال می‌گیریم. لذا متغیر تصادفی $Z = \frac{\bar{Y}_{st} - \bar{Y}_N}{\sigma(\bar{Y}_{st})}$ دارای توزیع نرمال استاندارد است. اگر z مقدار متغیر نرمال استاندارد متناظر با احتمال تجمعی $1 - \frac{\alpha}{2}$ باشد، بازه اطمینان $(1 - \alpha) 100\%$ درصد \bar{Y}_N به صورت زیر است

$$[\bar{Y}_{st} - z\sigma(\bar{Y}_{st}), \quad \bar{Y}_{st} + z\sigma(\bar{Y}_{st})] \quad (12.4)$$

این بازه، بازه‌ای تصادفی است زیرا \bar{Y}_{st} تصادفی است. برای یک نمونه مشخص، این بازه ثابت است. معمولاً $\sigma(\bar{Y}_{st})$ مجهول است و باید به جای آن، $\hat{\sigma}(\bar{Y}_{st})$ را از رابطه (۹.۴) منظور کرد. لذا برآورد بازه اطمینان به صورت زیر است

$$[\bar{Y}_{st} - z\hat{\sigma}(\bar{Y}_{st}), \quad \bar{Y}_{st} + z\hat{\sigma}(\bar{Y}_{st})] \quad (13.4)$$

اگر جمله‌های بازه‌های (۱۲.۴) (۱۳.۴) را در N ضرب کنیم، بازه اطمینان و برآورد بازه اطمینان برای مجموع واحدهای جامعه، یعنی T_N به دست می‌آید. برآورد بازه اطمینان برای T_N با ضرب اطمینان $1 - \alpha$ به صورت زیر است

$$[N\bar{Y}_{st} - zN\hat{\sigma}(\bar{Y}_{st}), \quad N\bar{Y}_{st} + zN\hat{\sigma}(\bar{Y}_{st})]$$

همان‌طور که اشاره شد می‌پذیریم که \bar{Y}_{st} دارای توزیع نرمال است و z را از جدول توزیع نرمال استاندارد به دست می‌آوریم. اگر تعداد واحدهایی که از طبقه‌ها انتخاب می‌کنیم کم باشد باید مقدار z را به جای توزیع نرمال از توزیع t به دست آوریم. در این بازه اطمینان \bar{Y}_{st} و $\hat{\sigma}(\bar{Y}_{st})$ تصادفی هستند. توزیع $\hat{\sigma}(\bar{Y}_{st})$ در حالت کلی بسیار پیچیده است، ولی یک روش تقریبی برای تخصیص

با توجه به برابری

$$E(s_h^2) = S_h^2 \quad h = 1, 2, \dots, L$$

رابطه بالا به صورت زیر درمی‌آید

$$V(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{E(s_h^2)}{n_h}$$

چون، برای مقدار ثابت a ، $E(aX) = aE(X)$ که در آن X ، متغیری تصادفی است، پس رابطه بالا را می‌توان چنین نوشت

$$V(\bar{Y}_{st}) = E \left[\frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{s_h^2}{n_h} \right]$$

این برابری، معرف آن است که عبارت داخل کروشه برآوردکننده نااریب $V(\bar{Y}_{st})$ است، یعنی

$$\hat{V}(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{s_h^2}{n_h}$$

و یا

$$\hat{V}(\bar{Y}_{st}) = \sum_{h=1}^L W_h^2 (1 - f_h) \frac{s_h^2}{n_h} \quad (10.4)$$

در عمل، با داشتن تنها یک نمونه در هر طبقه، s_h^2 ها به صورت مقادیری عددی هستند و لذا برآوردی نااریب، برای واریانس \bar{Y}_{st} به دست می‌آید. صورت دیگر (۱۰.۴)، با توجه به آخرین برابری (۴.۴) به صورت زیر است

$$\hat{V}(\bar{Y}_{st}) = \sum_{h=1}^L \frac{W_h^2 s_h^2}{n_h} - \sum_{h=1}^L \frac{W_h s_h^2}{N} \quad (11.4)$$

برای محاسبه $\hat{V}(\bar{Y}_{st})$ باید از هر طبقه حداقل دو واحد انتخاب شود، زیرا در غیر این صورت محاسبه s_h^2 به دلیل صفر شدن $n_h - 1$ در رابطه

$$s_h^2 = \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (Y_{hi} - \bar{Y}_h)^2$$

میسر نیست. برای حالت نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب، از (۷.۴) داریم

$$\hat{V}(\bar{Y}_{st}) = \frac{1-f}{n} \sum_{h=1}^L W_h s_h^2$$

□

که برآوردکننده‌ای نااریب است.

۵.۴ حدود اطمینان میانگین و مجموع واحدهای جامعه در نمونه‌گیری با طبقه‌بندی

برای یک نمونه تصادفی ساده با طبقه‌بندی به حجم n از جامعه‌ای به حجم N ، تعداد نمونه‌های ممکن $\prod_{h=1}^L \binom{N_h}{n_h}$ است. مثلاً اگر ۴ طبقه به حجمهای ۱۰، ۸، ۹، ۱۲ داشته باشیم و از این طبقه‌ها به ترتیب ۴، ۲، ۳، ۵ واحد به تصادف انتخاب کنیم تعداد نمونه‌های ممکن برابر است با

$$\binom{10}{4} \binom{8}{2} \binom{9}{3} \binom{12}{5} = (120)(28)(84)(792) = 39118464$$

بنابراین در حالت کلی حداکثر به همین تعداد هم \bar{Y}_{st} خواهیم داشت. اگر تعداد واحدهایی که \bar{Y}_{st} را تولید می‌کنند یعنی n بزرگ باشد، مثلاً بیش از ۲۰، فرض نرمال بودن توزیع \bar{Y}_{st} را به شرط متناهی بودن واریانسهای طبقات، و با توجه به قضیه حد مرکزی می‌توان پذیرفت. با توجه به این مقدمه، برای تهیه بازه اطمینان، ابتدا توزیع \bar{Y}_{st} را نرمال می‌گیریم. لذا متغیر تصادفی $Z = \frac{\bar{Y}_{st} - \bar{Y}_N}{\sigma(\bar{Y}_{st})}$ دارای توزیع نرمال استاندارد است. اگر z مقدار متغیر نرمال استاندارد متناظر با احتمال تجمعی $1 - \frac{\alpha}{2}$ باشد، بازه اطمینان $100(1 - \alpha)$ درصد به صورت زیر است

$$[\bar{Y}_{st} - z\sigma(\bar{Y}_{st}), \quad \bar{Y}_{st} + z\sigma(\bar{Y}_{st})] \quad (12.4)$$

این بازه، بازه‌ای تصادفی است زیرا \bar{Y}_{st} تصادفی است. برای یک نمونه مشخص، این بازه ثابت است. معمولاً $\sigma(\bar{Y}_{st})$ مجهول است و باید به جای آن، $\hat{\sigma}(\bar{Y}_{st})$ را از رابطه (۹.۴) منظور کرد. لذا برآورد بازه اطمینان به صورت زیر است

$$[\bar{Y}_{st} - z\hat{\sigma}(\bar{Y}_{st}), \quad \bar{Y}_{st} + z\hat{\sigma}(\bar{Y}_{st})] \quad (13.4)$$

اگر جمله‌های بازه‌های (۱۲.۴) (۱۳.۴) را در N ضرب کنیم، بازه اطمینان و برآورد بازه اطمینان برای مجموع واحدهای جامعه، یعنی T_N به دست می‌آید. برآورد بازه اطمینان برای T_N با ضریب اطمینان $1 - \alpha$ به صورت زیر است

$$[N\bar{Y}_{st} - zN\hat{\sigma}(\bar{Y}_{st}), \quad N\bar{Y}_{st} + zN\hat{\sigma}(\bar{Y}_{st})]$$

همان‌طور که اشاره شد می‌پذیریم که \bar{Y}_{st} دارای توزیع نرمال است و z را از جدول توزیع نرمال استاندارد به دست می‌آوریم. اگر تعداد واحدهایی که از طبقه‌ها انتخاب می‌کنیم کم باشد باید مقدار z را به جای توزیع نرمال از توزیع t به دست آوریم. در این بازه اطمینان \bar{Y}_{st} و $\hat{\sigma}(\bar{Y}_{st})$ تصادفی هستند. توزیع $\hat{\sigma}(\bar{Y}_{st})$ در حالت کلی بسیار پیچیده است، ولی یک روش تقریبی برای تخصیص

درجه آزادی، وقتی از توزیع t استفاده می‌کنیم، تخصیص ساترتوایت^۱ است. می‌توانیم $\hat{\sigma}^2(\bar{Y}_{st})$ را چنین بنویسیم

$$\hat{\sigma}^2(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L g_h s_h^2 \quad (14.4)$$

که در آن

$$g_h = \frac{N_h(N_h - n_h)}{n_h}$$

درجه آزادی مؤثر را با n_e نشان می‌دهیم که به صورت زیر محاسبه می‌شود

$$n_e = \frac{\left(\sum_{h=1}^L g_h s_h^2 \right)^2}{\sum_{h=1}^L \frac{g_h^2 s_h^2}{n_h - 1}} \quad (15.4)$$

مقدار n_e همیشه بین کوچکترین مقدار از مقادیر $1, 2, \dots, L, n_h - 1$ و مجموع این مقادیر واقع می‌شود. اگر توزیع جامعه دارای چاولگی مثبت باشد، فرمول (۱۵.۴) درجه آزادی را بیشتر برآورد می‌کند.

۶.۴ تخصیص اپتیم

در نمونه‌گیری تصادفی با طبقه‌بندی، معمولاً اگر هزینه انتخاب هر واحد نمونه در همه طبقه‌ها یکسان باشد، نمونه‌گیر از طبقه‌های با حجم بیشتر نمونه‌ای با حجم بزرگتر انتخاب می‌کند، ولی اگر هزینه انتخاب هر واحد متفاوت باشد برای انتخاب حجم نمونه فقط حجم طبقه نمی‌تواند ملاک انتخاب حجم نمونه باشد. در این صورت خط‌مشی نمونه‌گیر باید بر این دو نکته متکی باشد که n_h را طوری انتخاب کند که برای بودجه‌ای معین، $V(\bar{Y}_{st})$ مینیمم شود و یا برای مقدار از پیش تعیین شده $V(\bar{Y}_{st})$ ، هزینه نمونه‌گیری در کل مینیمم باشد.

اگر در طبقه h ام هزینه هر واحد نمونه‌گیری C_h فرض شود، وقتی حجم نمونه در این طبقه n_h باشد، هزینه نمونه‌گیری از این طبقه برابر با $C_h n_h$ است و لذا هزینه نمونه‌گیری از L طبقه برابر با $\sum_{h=1}^L C_h n_h$ خواهد شد. اگر هزینه رفت و آمد بین طبقات و هزینه‌های اداری و غیره را C_0 فرض کنیم، هزینه کل نمونه‌گیری، C ، با معادله زیر بیان می‌شود

$$C = C_0 + \sum_{h=1}^L C_h n_h \quad (16.4)$$

اگر در حالتی خاص هزینه‌های اداری وجود نداشته باشند و C_0 فقط همان هزینه رفت و آمد بین طبقه‌ها باشد نشان داده‌اند که می‌توان C_0 را برابر با $\sum t_h \sqrt{n_h}$ گرفت، که در آن t_h هزینه رفت و آمد به ازای هر واحد است. ما در این کتاب تابع هزینه را به صورت تابع هزینه (۱۶.۴) در نظر گرفته‌ایم که برحسب حجم نمونه، معادله‌ای خطی است.

نفسیه ۶.۲ در نمونه‌گیری تصادفی با طبقه‌بندی که تابع هزینه به صورت (۱۶.۴) است، برآورد واریانس \bar{Y}_{st} برای مقدار مشخص C و یا هزینه C برای واریانس مشخص V ، وقتی مینیمم می‌شود که n_h متناسب با $W_h S_h / \sqrt{C_h}$ باشد.

برهان. داریم

$$C = C_0 + \sum_{h=1}^L C_h n_h$$

از طرفی اگر $V(\bar{Y}_{st})$ را به قصد خلاصه‌نویسی فقط با V نشان دهیم، بنابر (۲.۴)، داریم

$$V = V(\bar{Y}_{st}) = \sum_{h=1}^L \frac{W_h^2 S_h^2}{n_h} - \sum_{h=1}^L \frac{W_h S_h^2}{N} \quad (I)$$

آنچه باید مورد بحث قرار دهیم عبارت‌اند از:

(۱) مقادیر n_h را طوری انتخاب کنیم که وقتی C مشخص است مقدار V مینیمم باشد.

(۲) مقادیر n_h را طوری انتخاب کنیم که وقتی V مشخص است مقدار C مینیمم باشد.

در واقع مسأله عبارت از انتخاب n_h ها برای مینیمم کردن V به ازای C ی ثابت یا مینیمم کردن

C به ازای V ی ثابت است. از رابطه I داریم

$$V + \sum_{h=1}^L \frac{W_h S_h^2}{N} = \sum_{h=1}^L \frac{W_h^2 S_h^2}{n_h} \quad (II)$$

جمله دوم طرف اول (II) ثابت است زیرا فقط به پارامترهای جامعه بستگی دارد.

از طرفی، تابع هزینه را به صورت زیر می‌نویسم

$$C - C_0 = \sum_{h=1}^L C_h n_h \quad (III)$$

اگر قرار دهیم $V' = V + \sum_{h=1}^L \frac{W_h S_h^2}{N}$ و $C' = C - C_0$ ، از ضرب طرفین دو برابری II و III داریم

$$C'V' = \left(\sum_{h=1}^L \frac{W_h^2 S_h^2}{n_h} \right) \left(\sum_{h=1}^L C_h n_h \right) \quad (IV)$$

وقتی V ثابت باشد V' هم ثابت است و باید مقادیر n_h را طوری تعیین کنیم که C و در نتیجه C' مینیمم شود. وقتی C ثابت باشد C' هم ثابت است و باید V و در نتیجه V' مینیمم شود. پس هر دو قسمت مسأله هم‌ارز این هستند که در رابطه (IV) وقتی V' ثابت است C' را مینیمم کنیم و وقتی C' ثابت است V را مینیمم کنیم. لذا در هر دو حالت باید مقادیر n_h را به‌قسمی بیابیم که طرف دوم (IV) مینیمم شود. برای انجام این کار یادآور می‌شویم که از اتحاد جبری

$$\left(\sum a_i^2\right) \left(\sum b_i^2\right) - \left(\sum a_i b_i\right)^2 = \sum_{i < j} \sum_j (a_i b_j - a_j b_i)^2$$

که در آن a_i و b_i مثبت‌اند و حدود \sum برای تمام جملات یکی است، به‌سهولت نابرابری زیر که به نابرابری کوشی-شوارتس موسوم است نتیجه می‌شود

$$\left(\sum a_i^2\right) \left(\sum b_i^2\right) \geq \left(\sum a_i b_i\right)^2$$

برابری وقتی برقرار می‌شود که به‌ازای هر مقدار i داشته باشیم $\frac{b_i}{a_i} = C^{te}$ ، زیرا باید طرف دوم اتحاد که مجموع چندین مربع است برابر با صفر باشد، و ناچار هر یک از جملات باید صفر باشد که به تناسب

$$\frac{b_j}{a_j} = \frac{b_i}{a_i}, \quad j \neq i$$

منجر می‌شود که همان شرط $\frac{b_i}{a_i} = C^{te}$ را به‌دست می‌دهد. حال اگر در نابرابری کوشی-شوارتس قرار دهیم $b_i = \sqrt{C_h n_h}$ و $a_i = \frac{W_h S_h}{\sqrt{n_h}}$ داریم

$$\left(\sum \frac{W_h^2 S_h^2}{n_h}\right) \left(\sum C_h n_h\right) \geq \left(\sum W_h S_h \sqrt{C_h}\right)^2$$

اما طرف اول این نابرابری با $C'V'$ برابر است و وقتی مینیمم می‌شود که با طرف دوم برابر باشد و آن هم وقتی رخ می‌دهد که داشته باشیم

$$\frac{b_h}{a_h} = \frac{\sqrt{C_h n_h}}{\frac{W_h S_h}{\sqrt{n_h}}} = \frac{n_h \sqrt{C_h}}{W_h S_h} = c^{te} = K$$

یعنی

$$n_h = \frac{W_h S_h}{\sqrt{C_h}} \cdot K$$

و در نتیجه

$$n = \sum_{h=1}^L n_h = \sum_{h=1}^L \frac{W_h S_h}{\sqrt{C_h}} \cdot K$$

از تقسیم نظیر به نظیر طرفین دو برابری اخیر بر یکدیگر داریم

$$\frac{n_h}{n} = \frac{W_h S_h / \sqrt{C_h}}{\sum_{h=1}^L (W_h S_h / \sqrt{C_h})}$$

یا

$$n_h = \frac{W_h S_h / \sqrt{C_h}}{\sum_{h=1}^L (W_h S_h / \sqrt{C_h})} \cdot n = \frac{N_h S_h / \sqrt{C_h}}{\sum_{h=1}^L (N_h S_h / \sqrt{C_h})} \cdot n \quad (۱۷.۴)$$

□ یعنی $n_h \propto W_h S_h / \sqrt{C_h}$ یا $n_h \propto N_h S_h / \sqrt{C_h}$

از این قضیه بی‌درنگ می‌توان نتایج زیر را توجیه کرد:

۱. چون n_h نسبت مستقیم با N_h دارد هرچه حجم طبقه بیشتر باشد حجم نمونه باید بزرگتر باشد.

۲. چون n_h با S_h نسبت مستقیم دارد، هرچه تغییرات طبقه، یعنی ناهمگنی واحدها بیشتر

باشد باید حجم نمونه بزرگتر باشد.

۳. چون n_h با $\sqrt{C_h}$ و در نتیجه با C_h نسبت معکوس دارد هرچه C_h کوچکتر باشد، یعنی

نمونه‌گیری در طبقه‌ای ارزاتر تمام شود، باید حجم نمونه آن بیشتر باشد.

از معادله (۱۷.۴) مقادیر $n_h (h = 1, \dots, L)$ برحسب n به دست می‌آیند ولی n مجهول

است. لذا n را برای دو مورد تحت بررسی، یعنی معلوم بودن C یا V ، به دست می‌آوریم.

الف. اگر هزینه کل نمونه‌گیری، C ، از قبل تثبیت شده باشد، معادله (۱۷.۴) مقادیر

$n_h (h = 1, \dots, L)$ را به گونه‌ای می‌دهد که V مینیمم باشد. اگر مقدار n_h را از معادله (۱۷.۴)

در معادله تابع هزینه قرار دهیم، داریم

$$C - C_0 = \sum_{h=1}^L C_h \cdot \frac{N_h S_h / \sqrt{C_h}}{\sum_{h=1}^L (N_h S_h / \sqrt{C_h})} \cdot n$$

از این معادله، مقدار n به صورت زیر نتیجه می‌شود

$$n = \frac{(C - C_0) \sum_{h=1}^L (N_h S_h / \sqrt{C_h})}{\sum_{h=1}^L (N_h S_h \sqrt{C_h})} \quad (۱۸.۴)$$

دستی C_0 ، هزینه کل نمونه‌گیری، مشخص باشد و مقادیر N_h ، S_h ، و C_h به ازای $h = 1, \dots, L$

معلوم باشند، ابتدا فرمول (۱۸.۴) را به کار می‌بریم و مقدار n را می‌یابیم. در گام بعد به کمک فرمول

(۱۷.۴) مقادیر n_h ها را معین می‌کنیم، و به نمونه‌گیری می‌پردازیم.

ب. اگر V از قبل تثبیت شده باشد، مقادیر n_h ها را از فرمول (۱۷.۴) در رابطه (۴.۴) قرار

می‌دهیم

$$V = \sum_{h=1}^L \left[W_h^2 S_h^2 / n_h - \frac{W_h S_h^2}{N} \right]$$

$$V = \sum_{h=1}^L W_h S_h^2 \left[W_h / \frac{W_h S_h / \sqrt{C_h}}{\sum_{h=1}^L W_h S_h / \sqrt{C_h}} n \right] - \sum_{h=1}^L \frac{W_h S_h^2}{N}$$

$$V = \frac{1}{n} \sum_{h=1}^L \left[(W_h S_h \sqrt{C_h}) \left(\sum_{h=1}^L W_h S_h / \sqrt{C_h} \right) \right] - \sum_{h=1}^L \frac{W_h S_h^2}{N}$$

وقتی V معلوم است از معادله بالا n به صورت زیر به دست می‌آید

$$n = \frac{\left[\sum_{h=1}^L W_h S_h \sqrt{C_h} \right] \left[\sum_{h=1}^L W_h S_h / \sqrt{C_h} \right]}{V + \frac{1}{N} \sum_{h=1}^L W_h S_h^2} \quad (۱۹.۴)$$

لذا ابتدا از (۱۹.۴) با داشتن W_h و S_h و C_h مقدار n را به دست می‌آوریم و سپس از (۱۷.۴) مقادیر n_h را مشخص می‌کنیم.

تبصره ۱. یک حالت خاص مهم، حالتی است که هزینه نمونه‌گیری برای هر واحد در تمام طبقات یکسان باشد، یعنی $C_h = c$. در این صورت $C = C_0 + cn$ درمی‌آید و تخصیص ایتیم برای هزینه ثابت به تخصیص ایتیم برای حجم نمونه‌ای ثابت تبدیل می‌شود، زیرا $n = \frac{C - C_0}{c}$. سپس مقادیر n_h از (۱۷.۴) که به صورت زیر خلاصه می‌شود به دست می‌آیند

$$n_h = \frac{W_h S_h / \sqrt{c}}{\sum_{h=1}^L (W_h S_h / \sqrt{c})} \cdot n = \frac{W_h S_h}{\sum_{h=1}^L W_h S_h} n = \frac{N_h S_h}{\sum_{h=1}^L N_h S_h} \cdot n \quad (۲۰.۴)$$

تخصیص (۲۰.۴) را گاهی تخصیص نیمین^۱ یا انتساب نیمین می‌نامند.

تبصره ۲. در تبصره ۱ دیدیم که برای $C_h = c$ و C ثابت واریانس \bar{Y}_{st} به ازای مقادیر n_h حاصل از (۲۰.۴) مینیمم می‌شود. مقدار این مینیمم را می‌توانیم به دست آوریم. می‌دانیم

$$V = V(\bar{Y}_{st}) = \sum_{h=1}^L \frac{W_h^2 S_h^2}{n_h} - \sum_{h=1}^L \frac{W_h S_h^2}{N}$$

به جای n_h مقدار آن را از (۲۰.۴) قرار می دهیم تا V ی اپتیم به دست آید

$$V_{opt} = \sum_{h=1}^L \frac{W_h^T S_h^T}{n W_h S_h / \sum_{h=1}^L W_h S_h} - \sum_{h=1}^L \frac{W_h S_h^T}{N}$$

$$= \frac{1}{n} \left(\sum_{h=1}^L W_h S_h \right)^T - \frac{1}{N} \sum_{h=1}^L W_h S_h^T \quad (21.4)$$

جملة دوم سمت راست، معرف fpc است.

تبصره ۳. برای استفاده از فرمولهای (۱۷.۴) تا (۲۱.۴) باید مقادیر N_h یا مقادیر W_h و S_h^T نیز معلوم باشند. در کاربردها معمولاً مقادیر N_h و در نتیجه W_h معلوم اند. اما باید با نمونه گیری تصادفی مقدماتی از هر طبقه، مثل مورد نمونه گیری تصادفی، ابتدا برآوردهایی ناریب برای S_h^T به دست آوریم و در فرمولهای مذکور به جای مقادیر S_h^T برآوردهای ناریب s_h^T را قرار دهیم. در این صورت n و n_h و V_{opt} به ترتیب به \hat{n} و \hat{n}_h و \hat{V}_{opt} تبدیل می شوند، و حجمهای نمونه ای و واریانس اپتیم با تقریب به دست می آیند.

مثال ۲.۴. مؤسسه ای تحقیقاتی برای تعیین متوسط مدت زمانی که بیماران مبتلا به دیابت در بیمارستانهای شهری بستری می شوند تعداد بیمارانی را که در یک سال در سه بیمارستان شهر بستری بوده اند در نظر می گیرند. این تعداد به ترتیب ۳۰۰، ۱۲۰، و ۱۸۰ است. از روی نمونه گیری مقدماتی در هر بیمارستان واریانس تعداد روزهای بستری بودن بیماران تقریباً برابر $\frac{299}{3}$ ، $\frac{238}{15}$ ، و $\frac{179}{5}$ به دست آمده است. هزینه کسب اطلاعات درباره هر بیمار در سه بیمارستان به ترتیب ۴، ۹، ۱۶ است. الف) مؤسسه ابتدا تصمیم می گیرد که کلاً n فرد از سه بیمارستان را به عنوان نمونه انتخاب کند. با اطلاعات بالا، در هر بیمارستان چه سهمی از n را باید در نظر گرفت تا تخصیص، اپتیم باشد؟ ب) اگر مؤسسه بودجه ای برابر ۸۸۰ واحد پول که شامل هزینه های اداری نیست برای انجام تحقیق تخصیص دهد، با این تخصیص دقیقاً حجم نمونه ای که باید از هر بیمارستان انتخاب کرد چقدر است؟

ج) قبل از اجرای نمونه گیری، اهمیت مطلب ایجاب می کند که برآورد متوسط روزهای بستری بودن را با دقتی خاص معین کنند، به شکلی که واریانس برآورد این متوسط برابر ۴۰ باشد. در این صورت از هر بیمارستان به چه حجمی باید نمونه گرفت؟

د) اگر از سه بیمارستان به ترتیب ۹۸، ۱۱، ۱۸ بیمار را به عنوان نمونه انتخاب کنیم و میانگین مدت بستری بودن بیماران این سه نمونه به ترتیب ۳۰، ۲۶، و ۲۴ روز باشد برآورد ناریبی برای متوسط مدت بستری بودن جامعه بیماران دیابتی به دست آورید. اگر تغییرات این سه نمونه به ترتیب ۸۰، ۱۰، و ۳۰ باشند برآورد واریانس متوسط مدت بستری بودن در هر بیمارستان و برآورد واریانس برآوردکننده متوسط مدت بستری بودن در بیمارستانها را به دست آورید.

ه) یک بازه اطمینان ۹۵ درصدی برای میانگین واقعی مدت بستری بودن بیابید.

الف) بیماران دیابتی به ۳ طبقه تقسیم شده‌اند به طوری که

$$N_1 = 300 \quad N_2 = 120 \quad N_3 = 180 \quad N = 600$$

واریانسهای طبقات به ترتیب عبارت‌اند از $\sigma_1^2 = \frac{299}{3}$ ، $\sigma_2^2 = \frac{238}{15}$ ، $\sigma_3^2 = \frac{179}{5}$ و

$$S_h^2 = \frac{N_h}{N_h - 1} \sigma_h^2 \quad h = 1, 2, 3$$

پس

$$S_1^2 = \frac{300}{299} \cdot \frac{299}{3} = 100$$

$$S_2^2 = \frac{120}{119} \cdot \frac{238}{15} = 16$$

$$S_3^2 = \frac{180}{179} \cdot \frac{179}{5} = 36$$

لذا $S_1 = 10$ ، $S_2 = 4$ و $S_3 = 6$.

حال با توجه به تخصیص ایتیم

$$n_h = \frac{N_h S_h / \sqrt{C_h}}{\sum (N_h S_h / \sqrt{C_h})} \cdot n$$

و داریم $C_1 = 4$ ، $C_2 = 9$ ، $C_3 = 16$. محاسبات را به صورت زیر انجام می‌دهیم

$$\begin{cases} N_1 S_1 / \sqrt{C_1} = \frac{(300)(10)}{2} = 1500 \\ N_2 S_2 / \sqrt{C_2} = \frac{(120)(4)}{3} = 160 \\ N_3 S_3 / \sqrt{C_3} = \frac{(180)(6)}{4} = 270 \end{cases}$$

$$\sum_{h=1}^3 N_h S_h / \sqrt{C_h} = 1930$$

بنابر فرمول بالا، داریم

$$n_1 = \frac{1500}{1930} n \simeq 0.777n$$

$$n_2 = \frac{160}{1930} n \simeq 0.082n$$

$$n_3 = \frac{270}{1930} n \simeq 0.139n$$

ب) بر اساس فرض مسأله $C - C_0 = 880$ می‌دانیم برای C ثابت، بنابر (۱۸.۴)

$$n = (C - C_0) \frac{\sum (N_h S_h / \sqrt{C_h})}{\sum N_h S_h \sqrt{C_h}}$$

داریم

$$\begin{cases} N_1 S_1 \sqrt{C_1} = 300(10)(2) = 6000 \\ N_2 S_2 \sqrt{C_2} = 120(4)(3) = 1440 \\ N_3 S_3 \sqrt{C_3} = 180(6)(2) = 2160 \end{cases}$$

و

$$\sum_{h=1}^3 N_h S_h \sqrt{C_h} = 11760$$

پس

$$n = \frac{(880)(1930)}{11760} \approx 145$$

و در نتیجه با توجه به نتایج قسمت الف، داریم

$$n_1 \approx 0.78n = 0.78(145) \approx 113$$

$$n_2 \approx 0.82n = 0.82(145) \approx 12$$

$$n_3 \approx 0.139n = 0.139(145) \approx 20$$

ج) از رابطه (۱۹.۴) برای واریانس از پیش تعیین شده V ، با توجه به رابطه $V_h = \frac{N_h}{N}$ ، داریم

$$n = \frac{(\sum N_h S_h \sqrt{C_h})(\sum N_h S_h / \sqrt{C_h})}{VN' + \sum N_h S_h^2}$$

از طرفی

$$\begin{cases} N_1 S_1^2 = 300 \times 10^2 = 30000 \\ N_2 S_2^2 = 120 \times 4^2 = 1920 \\ N_3 S_3^2 = 180 \times 6^2 = 6480 \end{cases}$$

و

$$\sum_{h=1}^r N_h S_h^2 = 38400$$

پس

$$n = \frac{(1930)(11760)}{0.4(600)^2 + 38400} \approx 125$$

و از آنجا

$$\begin{cases} n_1 = 125 \times 0.78 \approx 98 \\ n_2 = 125 \times 0.082 \approx 11 \\ n_3 = 125 \times 0.139 \approx 18 \end{cases}$$

مجموع n_i ها از n بیشتر می‌شود زیرا در محاسبه تقریبی آنها برای کم نشدن دقت نمونه‌گیری تقریب اضافی را منظور کرده‌ایم.
 (د) در نمونه‌گیری تصادفی از سه طبقه به ترتیب داریم

$$n_1 = 98, \quad \bar{Y}_1 = 30, \quad N_1 = 300, \quad s_1^2 = 80, \quad W_1 = \frac{300}{600} = \frac{1}{2}$$

$$n_2 = 11, \quad \bar{Y}_2 = 26, \quad N_2 = 120, \quad s_2^2 = 10, \quad W_2 = \frac{120}{600} = \frac{1}{5}$$

$$n_3 = 18, \quad \bar{Y}_3 = 24, \quad N_3 = 180, \quad s_3^2 = 30, \quad W_3 = \frac{180}{600} = \frac{3}{10}$$

لذا

$$\hat{Y}_N = \sum_{h=1}^r W_h \hat{Y}_h = \frac{1}{2}(30) + \frac{1}{5}(26) + \frac{3}{10}(24) = 27.4$$

برای تعیین برآورد واریانس متوسط مدت بستری بودن در هر بیمارستان، باید از فرمول برآورد واریانس میانگین نمونه‌گیری تصادفی استفاده کنیم. به شرح زیر

$$\hat{V}(\bar{Y}_1) = \left(\frac{1}{n_1} - \frac{1}{N_1} \right) s_1^2 = \left(\frac{1}{98} - \frac{1}{300} \right) 80 \approx 0.550$$

$$\hat{V}(\bar{Y}_2) = \left(\frac{1}{n_2} - \frac{1}{N_2} \right) s_2^2 = \left(\frac{1}{11} - \frac{1}{120} \right) 10 \approx 0.826$$

$$\hat{V}(\bar{Y}_3) = \left(\frac{1}{n_3} - \frac{1}{N_3} \right) s_3^2 = \left(\frac{1}{18} - \frac{1}{180} \right) 30 = 1.5$$

$$\begin{aligned}\hat{V}(\bar{Y}_{st}) &= \frac{1}{N^2} \sum_{h=1}^r N_h(N_h - n_h) \frac{S_h^2}{n_h} \\ &= \frac{1}{(600)^2} \left[300(300 - 98) \frac{10}{98} + 120(120 - 11) \frac{10}{11} + 180(180 - 18) \frac{30}{18} \right] \\ &\approx 0.305\end{aligned}$$

A

۷.۴ مقایسه دقت برآوردکننده‌ها در نمونه‌گیری تصادفی با طبقه‌بندی

و نمونه‌گیری تصادفی ساده

به صورت شهودی استنباط ما این است که نمونه‌گیری تصادفی با طبقه‌بندی همیشه برآوردکننده‌ای از میانگین را ارائه می‌دهد که واریانسش از واریانس برآوردکننده میانگین نمونه‌گیری تصادفی ساده کوچکتر است، اما این احساس شهودی همیشه درست نیست. مواردی وجود دارند که دقت نمونه‌گیری تصادفی ساده بیش از نمونه‌گیری تصادفی با طبقه‌بندی است. اگر مقادیر n_h خیلی از تخصیص اپتیمم دور باشند ممکن است واریانس برآوردکننده میانگین بزرگ باشد. در واقع حتی وقتی طبقه‌بندی برای حجم تثبیت‌شده نمونه با تخصیص اپتیمم انجام شود امکان دارد واریانس برآوردکننده بزرگ باشد. در این بخش مقایسه‌ای بین نمونه‌گیری تصادفی ساده و نمونه‌گیری تصادفی با طبقه‌بندی در حالت‌های تخصیص اپتیمم و متناسب انجام می‌دهیم. این مقایسه نشان خواهد داد در چه مواقعی انجام طبقه‌بندی مناسب است. قبل از آغاز مقایسه یادآور می‌شویم که واریانس میانگین نمونه تصادفی را با نماد V_{ran} و واریانس \bar{Y}_{st} در حالت تخصیص اپتیمم را با نماد V_{opt} و واریانس در حالت تخصیص متناسب را با نماد V_{prop} نشان می‌دهیم.

قضیه ۷.۴ اگر $\frac{1}{N_h}$ قابل اغماض باشد. آنگاه

$$V_{opt} \leq V_{prop} \leq V_{ran} \quad (22.4)$$

که در آن منظور از تخصیص اپتیمم تخصیص نین، برای n ثابت است، یعنی با $n_h \propto N_h S_h$ برهان. قبلاً دیدیم که به ترتیب مطابق با (۵.۲)، (۷.۴) و (۲۱.۴)

$$\begin{aligned}V_{ran} &= (1 - f) \frac{S^2}{n} \\ V_{prop} &= \frac{1 - f}{n} \sum W_h S_h^2 = \frac{\sum W_h S_h^2}{n} - \frac{\sum W_h S_h^2}{N} \\ V_{opt} &= \frac{(\sum W_h S_h)^2}{n} - \frac{\sum W_h S_h^2}{N}\end{aligned}$$

از طرفی تغییرات جامعه به صورت زیر است

$$\begin{aligned} S^r &= \frac{1}{N-1} \sum_{h=1}^L \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_N)^2 \\ &= \frac{1}{N-1} \sum_{h=1}^L \sum_{i=1}^{N_h} [(Y_{hi} - \bar{Y}_h) + (\bar{Y}_h - \bar{Y}_N)]^2 \end{aligned}$$

پس

$$\begin{aligned} (N-1)S^r &= \sum_{h=1}^L \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2 + \sum_{h=1}^L \sum_{i=1}^{N_h} (\bar{Y}_h - \bar{Y}_N)^2 \\ &\quad + 2 \sum_{h=1}^L \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)(\bar{Y}_h - \bar{Y}_N) \\ &= \sum_{h=1}^L \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2 + \sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 \end{aligned}$$

سومین جمله در عبارت وسط برابر با ۰ است (استدلال به عهده دانشجوست). پس

$$(N-1)S^r = \sum_{h=1}^L (N_h - 1)S_h^r + \sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 \quad (23.4)$$

زیرا برای طبقه h ام

$$S_h^r = \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (Y_{hi} - \bar{Y}_h)^2$$

اگر $\frac{1}{N_h}$ قابل اغماض باشد مسلماً $\frac{1}{N}$ نیز قابل اغماض است. پس از رابطه (۲۳.۴) داریم

$$\begin{aligned} S^r &= \sum_{h=1}^L \frac{N_h - 1}{N - 1} S_h^r + \sum_{h=1}^L \frac{N_h}{N - 1} (\bar{Y}_h - \bar{Y}_N)^2 \\ &= \sum_{h=1}^L \frac{\frac{N_h}{N} - \frac{1}{N}}{1 - \frac{1}{N}} S_h^r + \sum_{h=1}^L \frac{\frac{N_h}{N}}{1 - \frac{1}{N}} (\bar{Y}_h - \bar{Y}_N)^2 \\ &\approx \sum_{h=1}^L \frac{N_h}{N} S_h^r + \sum_{h=1}^L \frac{N_h}{N} (\bar{Y}_h - \bar{Y}_N)^2 \end{aligned}$$

پس

$$S^r \approx \sum_{h=1}^L W_h S_h^r + \sum_{h=1}^L W_h (\bar{Y}_h - \bar{Y}_N)^2 \quad (24.4)$$

اگر این مقدار را به جای S^2 در V_{ran} قرار دهیم، نتیجه می‌شود

$$V_{ran} = (1-f) \frac{S^2}{n} \simeq \frac{1-f}{n} \sum_{h=1}^L W_h S_h^2 + \frac{1-f}{n} \sum_{h=1}^L W_h (\bar{Y}_h - \bar{Y}_N)^2$$

که اولین جمله طرف دوم برابر V_{prop} است. پس

$$V_{ran} = V_{prop} + \frac{1-f}{n} \sum_{h=1}^L W_h (\bar{Y}_h - \bar{Y}_N)^2 \quad (25.4)$$

از این برابری نتیجه می‌شود که با توجه به نامنفی بودن جمله دوم طرف دوم

$$V_{prop} \leq V_{ran} \quad (26.4)$$

برابری وقتی برقرار می‌شود که \bar{Y}_h ها با \bar{Y}_N برابر باشند. از طرفی با توجه به روابط ابتدای برهان

$$V_{prop} - V_{opt} = \frac{1}{n} \left[\sum_{h=1}^L W_h S_h^2 - \left(\sum_{h=1}^L W_h S_h \right)^2 \right]$$

اگر قرار دهیم $\bar{S} = \sum_{h=1}^L W_h S_h$

$$\begin{aligned} V_{prop} - V_{opt} &= \frac{1}{n} \left[\sum_{h=1}^L W_h S_h^2 - \bar{S}^2 \right] \\ &= \frac{1}{n} \left[\sum_{h=1}^L W_h (S_h - \bar{S})^2 \right] \end{aligned} \quad (27.4)$$

زیرا

$$\begin{aligned} \sum_{h=1}^L W_h \bar{S}^2 - 2 \sum_{h=1}^L W_h S_h \bar{S} &= \bar{S}^2 \sum_{h=1}^L W_h - 2 \bar{S} \sum_{h=1}^L W_h S_h \\ &= \bar{S}^2 - 2 \bar{S} \cdot \bar{S} = -\bar{S}^2 \end{aligned}$$

پس $V_{prop} - V_{opt} \geq 0$ ، زیرا در (27.4) طرف آخر عبارتی نامنفی است. لذا

$$V_{opt} \leq V_{prop} \quad (28.4)$$

برابری وقتی برقرار است که S_h ها با \bar{S} برابر باشند. از ترکیب (26.4) و (28.4) نتیجه می‌شود که

$$V_{opt} \leq V_{prop} \leq V_{ran}$$

□

که برهان را تکمیل می‌کند.

نتیجه. از روابط (۲۵.۴) و (۲۷.۴) رابطه زیر نتیجه می‌شود

$$V_{ran} = V_{opt} + \frac{1}{n} \sum_{h=1}^L W_h (S_h - \bar{S})^2 + \frac{1-f}{n} \sum_{h=1}^L W_h (\bar{Y}_h - \bar{Y}_N)^2 \quad (29.4)$$

این رابطه نشان می‌دهد که اگر دو نمونه‌گیری تصادفی ساده و تصادفی با طبقه‌بندی در حالت تخصیص اپتیمم را در نظر بگیریم، دو مؤلفه در تقلیل واریانس دخالت دارند. مؤلفه اول (جمله آخر طرف دوم رابطه) وقتی کوچک است که میانگینهای طبقه‌ها خیلی از هم فاصله نداشته باشند، و مؤلفه دوم (جمله میانی طرف دوم) وقتی کوچک است که تفاوت بین انحراف معیارهای طبقات کم باشد. این مؤلفه، معرف تفاوت بین واریانس تخصیص اپتیمم و واریانس تخصیص متناسب است. اگر N_h قابل اغماض نباشد و اگر به جای S^2 از (۲۳.۴) استفاده کنیم، پس از محاسبات لازم، نتیجه می‌شود

$$V_{ran} = V_{prop} + \frac{1-f}{n(N-1)} \left[\sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 - \frac{1}{N} \sum_{h=1}^L (N - N_h) S_h^2 \right]$$

با توجه به عبارت داخل کروشه، اگر داشته باشیم

$$\sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 < \frac{1}{N} \sum_{h=1}^L (N - N_h) S_h^2 \quad (30.4)$$

آنگاه مقدار عبارت داخل کروشه منفی است و $V_{ran} \leq V_{prop}$. یعنی با شرط (۳۰.۴) نمونه‌گیری تصادفی ساده کاراتر از نمونه‌گیری با طبقه‌بندی و با تخصیص متناسب است. اگر همه S_h^2 ها، $h = 1, \dots, L$ ، باهم برابر باشند و مقدار مشترک آنها را با S_W^2 نشان دهیم به قسمی که تخصیص متناسب به مفهوم نیمن، اپتیمم باشد، در این صورت نابرابری (۳۰.۴) به صورت زیر درمی‌آید

$$\sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 < \frac{1}{N} \sum_{h=1}^L (N - N_h) S_W^2$$

و یا

$$\sum_{h=1}^L N_h (\bar{Y}_h - \bar{Y}_N)^2 < \frac{S_W^2}{N} (NL - N) = (L - 1) S_W^2$$

یا

$$\sum_{h=1}^L \frac{N_h (\bar{Y}_h - \bar{Y}_N)^2}{L - 1} < S_W^2 \quad (31.4)$$

اگر تحلیل واریانس را به خاطر داشته باشید از این رابطه نتیجه می‌گیرید که باید میانگین مربعات در

بین طبقات، کوچکتر از میانگین مربعات در درون طبقات باشد، یا به عبارت دیگر نسبت F باید کوچکتر از ۱ باشد. در چنین حالتی نمونه‌گیری تصادفی کاراتر از نمونه‌گیری تصادفی با طبقه‌بندی در حالت تخصیص متناسب است.

۸.۴ حالت خاص تخصیص نین

وقتی n_h ها را برای مقداری تثبیت شده از n محاسبه می‌کنیم ممکن است یک یا چند n_h از N_h های مربوط بزرگتر شوند. اگر تخصیص مورد نظر، تخصیص نین باشد و مثلاً داشته باشیم $n_1 > N_1$ ، آنگاه خط‌مشی معمول این است: اگر جامعه بیش از دو طبقه داشته باشد

$$\hat{n}_1 = N_1, \quad \hat{n}_h = (n - N_1) \frac{W_h S_h}{\sum_{h=2}^L W_h S_h} \quad h \geq 2 \quad (۳۲.۴)$$

درواقع چون $n_1 > N_1$ نامفهوم است حداکثر $\hat{n}_1 = N_1$ است، سپس از رابطه $n_h = \frac{W_h S_h}{\sum_{h=1}^L W_h S_h} \cdot n$ مربوط به تخصیص نین استفاده می‌کنیم. بدین طریق که ابتدا N_1 را از حجم کل نمونه کم می‌کنیم و $n - N_1$ واحد باقیمانده نمونه را از بین $h - 1$ طبقه دیگر با استفاده از فرمول تخصیص نین انتخاب می‌کنیم که (۳۲.۴) را نتیجه می‌دهد. اگر $n_1 > N_1$ و $n_2 > N_2$ به دست آمده باشند، خط‌مشی معمول چنین است

$$\hat{n}_1 = N_1, \quad \hat{n}_2 = N_2, \quad \hat{n}_h = (n - N_1 - N_2) \frac{W_h S_h}{\sum_{h=3}^L W_h S_h} \quad h \geq 3 \quad (۳۳.۴)$$

استدلال مربوط به درستی رابطه (۳۳.۴) شبیه استدلال حالت قبل است. به طور کلی هر تعداد از n_h هایی را که بزرگتر از N_h مربوط باشند با N_h برآورد می‌کنیم و سپس مجموع این N_h ها را از کل n می‌کاهیم و $n - \sum N_h$ واحد باقیمانده را از طبقات باقیمانده با استفاده از فرمول تخصیص نین به دست می‌آوریم. نکته‌ای که در این مورد حائز اهمیت است این است که فرمول

$$V_{min}(\bar{Y}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L W_h S_h \right)^2 - \frac{1}{N} \sum_{h=1}^L W_h S_h^2$$

که در تخصیص نین برای محاسبه واریانس مینیم به کار می‌رود در این حالت برقرار نیست. اگر \sum' معرف مجموعیابی روی طبقه‌هایی باشد که برای آنها $\hat{n}_h < N_h$ است، فرمول واریانس مینیم به صورت زیر حاصل می‌شود

$$V_{min}(\bar{Y}_{st}) = \frac{(\sum' W_h S_h)^2}{n'} - \frac{\sum' W_h S_h^2}{N}$$

که در آن n' مجموع حجمهای نمونه‌های طبقه‌های با شرط $\hat{n}_h < N_h$ است.

۹.۴ برآورد حجم نمونه کل وقتی وزن نمونه‌ها معلوم است

برای تعیین مقدار n ، تحت تخصیص ایتیم فرمولهایی ارائه دادیم. در این بخش برای هر تخصیصی، فرمولهایی، معرفی می‌شود. فرض می‌کنیم که \bar{Y}_{st} دارای واریانس معلوم V بوده، s_h برآورد S_h و مقادیر $w_h = \frac{n_h}{n}$ معلوم باشند. می‌دانیم بنابر (۱۱.۴)

$$V \simeq \sum_{h=1}^L \frac{W_h^2 s_h^2}{n_h} - \frac{1}{N} \sum_{h=1}^L W_h s_h^2$$

اگر به جای n_h مقدار آن را در این رابطه قرار دهیم، نتیجه می‌شود

$$V \simeq \frac{1}{n} \sum_{h=1}^L \frac{W_h^2 s_h^2}{w_h} - \frac{1}{N} \sum_{h=1}^L W_h s_h^2$$

از این برابری، رابطه کلی زیر را برای \hat{n} به دست می‌آوریم

$$\hat{n} = \frac{\sum_{h=1}^L \frac{W_h^2 s_h^2}{w_h}}{V + \frac{1}{N} \sum_{h=1}^L W_h s_h^2} \quad (۳۴.۴)$$

اگر N بزرگ باشد، آنگاه

$$n_0 = \frac{1}{V} \sum_{h=1}^L \frac{W_h^2 s_h^2}{w_h} \quad (۳۵.۴)$$

و اگر N بزرگ نباشد، از ترکیب (۳۴.۴) و (۳۵.۴) داریم

$$\hat{n} = \frac{n_0}{1 + \frac{1}{NV} \sum_{h=1}^L W_h s_h^2} \quad (۳۶.۴)$$

در حالت‌های خاص، فرمولها به صورتهایی درمی‌آیند که برای محاسبه راحت‌ترند. چند مورد را در زیر می‌آوریم:

الف) تخصیص را نیمی می‌گیریم (n تثبیت شده است): وقتی $n_h \simeq n \frac{W_h s_h}{\sum W_h s_h}$ یعنی $w_h = \frac{W_h s_h}{\sum W_h s_h} = \frac{n_h}{n}$ آن‌گاه صورت کسر (۳۴.۴) به صورت زیر درمی‌آید

$$\sum_{h=1}^L \frac{W_h^2 s_h^2}{W_h s_h / \sum_{h=1}^L W_h s_h} = \left(\sum_{h=1}^L W_h s_h \right)^2$$

و در نتیجه (۳۴.۴) چنین نوشته می شود

$$\hat{n} = \frac{(\sum_{h=1}^L W_h s_h)^2}{V + \frac{1}{N} \sum_{h=1}^L W_h s_h^2} \quad (۳۷.۴)$$

ب) تخصیص متناسب. در تخصیص متناسب $W_h = \frac{N_h}{N}$ بنا بر این (۳۵.۴) را می توان به صورت زیر نوشت

$$n_o = \frac{1}{V} \sum_{h=1}^L W_h s_h^2 \quad (۳۸.۴)$$

و (۳۶.۴) به صورت زیر درمی آید

$$\hat{n} = \frac{n_o}{1 + \frac{n_o}{N}} \quad (۳۹.۴)$$

۱۰.۴ تعیین حجم نمونه وقتی واریانس برآوردکننده مجموع مقادیر واحدها از قبل مشخص شده است

اگر مایل باشیم n را به قسمی بیابیم که $V(\hat{T}_N)$ برابر مقدار معین V باشد، در فرمولهای قبل به جای V مقدار $\frac{V}{N^2}$ را قرار می دهیم تا مقدار n مورد نظر به دست آید، زیرا در فرمولهای بالا V برابر با واریانس \bar{Y}_{st} است در حالی که در این بخش V معرف $V(\hat{T}_N)$ یا $V(N\bar{Y}_{st})$ است که برابر با $N^2 V(\bar{Y}_{st})$ است. از قرار دادن $\frac{V}{N^2}$ به جای V در رابطه (۳۴.۴)، برای حالت کلی داریم

$$\hat{n} = \frac{\sum_{h=1}^L W_h^2 s_h^2 / w_h}{\frac{V}{N^2} + \frac{1}{N} \sum_{h=1}^L W_h s_h^2} = \frac{\sum_{h=1}^L \frac{N_h^2}{N^2} s_h^2 / w_h}{\frac{V}{N^2} + \frac{1}{N} \sum_{h=1}^L \frac{N_h}{N} s_h^2}$$

و یا به صورت خلاصه

$$\hat{n} = \frac{\sum_{h=1}^L N_h^2 s_h^2 / w_h}{V + \sum_{h=1}^L N_h s_h^2} \quad (۴۰.۴)$$

برای تخصیص نینم (n تثبیت شده) از (۳۷.۴) داریم

$$\hat{n} = \frac{(\sum_{h=1}^L W_h s_h)^2}{\frac{V}{N^2} + \frac{1}{N} \sum_{h=1}^L W_h s_h^2} = \frac{(\sum_{h=1}^L \frac{N_h}{N} s_h)^2}{\frac{V}{N^2} + \frac{1}{N} \sum_{h=1}^L \frac{N_h}{N} s_h^2}$$

و یا به صورت خلاصه

$$\hat{n} = \frac{(\sum_{h=1}^L N_h s_h)^2}{V + \sum_{h=1}^L N_h s_h^2} \quad (41.4)$$

و برای تخصیص متناسب از (۳۸.۴) و (۳۹.۴) داریم

$$n_o = \frac{1}{V/N^2} \sum_{h=1}^L W_h s_h^2 = \frac{N^2}{V} \sum_{h=1}^L \frac{N_h}{N} s_h^2$$

و یا

$$n_o = \frac{N}{V} \sum_{h=1}^L N_h s_h^2 \quad (42.4)$$

و

$$\hat{n} = \frac{n_o}{1 + \frac{n_o}{N}}$$

مثال ۳.۴* برای برآورد تعداد کسانی که از آنها در سال اول دانشکده‌های یک دانشگاه ثبت نام شده است، ۱۹۶ مسؤل ثبت نام دانشگاه را در نظر گرفته‌اند. این مسؤلین را ابتدا در ۷ طبقه منظور سپس یک طبقه را به علت کوچکی تعداد در ۶ طبقه دیگر ادغام کرده‌اند. مقادیر N_h ها در جدول زیر آمده‌اند. قرار است از هر طبقه تعدادی را به عنوان نمونه انتخاب کنند و از روی تعداد ثبت نام آنها، تعداد کل ثبت نامها را برآورد کنند. مقدار s_h را از روی ثبت نامهای دو سال قبل به صورت مقادیری از s_h که در جدول آمده‌اند برآورد کرده‌اند. می‌خواهند n حجم کل نمونه را به قسمی بیابند که ضریب تغییرات در کل ثبت نام ۵٪ باشد. کل تعداد ثبت نامهای دو سال پیش دانشگاه ۵۶۴۷۲ بوده است. با توجه به میزان ثبت نام دو سال پیش و تعریف ضریب تغییرات، خطای معیار مطلوب

$$\sigma = (0.05)(56472) = 2824$$

$$V = (2824)^2 = 7974976$$

ممکن است این انتقاد وارد باشد که ضریب تغییرات دو سال قبل برای ثبت نامهای فعلی صادق نبوده و ضریب تغییرات فعلی بیشتر است. اما فرض می‌کنیم که این ضریب ثابت مانده باشد. در جدول زیر مقادیر N_h ، s_h و $N_h s_h$ را قبل از تعیین n آورده‌ایم. فرمول مربوط برای تعیین n فرمول (۴۱.۴) است که برای تخصیص اپتیمم در برآورد مجموع به‌کار می‌رود. با ۱۹۶ واحد در این

* اقتباس از اثر کوکران: تکنیکهای نمونه‌گیری (۱۹۷۷) وایلی.

نمونه‌گیری با طبقه‌بندی برای برآورد نسبتها ۱۶۵

جامعه، نمی‌توان از fpc صرف‌نظر کرد. اما ابتدا fpc را نادیده می‌گیریم و n_0 را حساب می‌کنیم

$$n_0 = \frac{(\sum N_h s_h)^2}{V} = \frac{(26841)^2}{7974976} = 9034$$

و برای محاسبه n بدون نادیده گرفتن fpc ، داریم

$$\hat{n} = \frac{n_0}{1 + \frac{1}{V} \sum N_h s_h^2} = \frac{9034}{1 + \frac{4640387}{7974976}} \approx 57$$

با توجه به رابطه $\hat{n}_h = n \cdot \frac{N_h s_h}{\sum N_h s_h}$ می‌توان مقادیر حجم نمونه‌ای هر طبقه را مشخص کرد. این مقادیر هم در جدول زیر آمده‌اند

| طبقه | N_h | s_h | $N_h s_h$ | $N_h s_h^2$ | \hat{n}_h |
|-------|-------|-------|-----------|-------------|-------------|
| ۱ | ۱۳ | ۳۲۵ | ۴۲۲۵ | ۱۳۷۳۱۲۵ | ۹ |
| ۲ | ۱۸ | ۱۹۰ | ۳۴۲۰ | ۶۴۹۸۰۰ | ۷ |
| ۳ | ۲۶ | ۱۸۹ | ۴۹۱۴ | ۹۲۸۷۴۶ | ۱۱ |
| ۴ | ۴۲ | ۸۲ | ۳۴۴۴ | ۲۸۲۴۰۸ | ۷ |
| ۵ | ۷۳ | ۸۶ | ۶۲۷۸ | ۵۳۹۹۰۸ | ۱۳ |
| ۶ | ۲۴ | ۱۹۰ | ۴۵۶۰ | ۸۶۶۴۰۰ | ۱۰ |
| مجموع | ۱۹۶ | | ۲۶۸۴۱ | ۴۶۴۰۳۸۷ | ۵۷ |

▲

۱۱.۴ نمونه‌گیری با طبقه‌بندی برای برآورد نسبتها

فرض می‌کنیم هر واحد جامعه به یکی از دو رده مجزا از هم C و C' متعلق باشد. می‌خواهیم در این جامعه نسبت واحدهایی را بیابیم که در رده معین C می‌افتند. اگر جامعه در L طبقه سازمان یافته باشد، آنگاه فرض می‌کنیم A_h تعداد واحدهایی از طبقه h است که در رده C می‌افتند و a_h تعداد واحدهایی از نمونه تصادفی این طبقه باشد که در رده C هستند. پس

$$P_h = \frac{A_h}{N_h}, \quad p_h = \frac{a_h}{n_h} \quad h = 1, 2, \dots, L$$

که P_h ، نسبت واحدهایی از طبقه h جامعه و p_h نسبت واحدهایی از نمونه طبقه h است که در رده C می‌افتند. قبلاً در نمونه‌گیری تصادفی ساده دیدیم که p_h برآوردکننده ناریب P_h است. اگر p_{st} را به صورت

$$p_{st} = \sum_{h=1}^L W_h p_h \quad (43.4)$$

تعریف کنیم. واضح است که

$$\begin{aligned} E(p_{st}) &= \sum_{h=1}^L W_h E(p_h) = \sum_{h=1}^L W_h P_h = \sum_{h=1}^L W_h \frac{A_h}{N_h} \\ &= \sum_{h=1}^L \frac{N_h}{N} \cdot \frac{A_h}{N_h} = \frac{1}{N} \sum_{h=1}^L A_h \end{aligned}$$

اما $\sum_{h=1}^L A_h$ برابر مجموع واحدهایی از جامعه است که در رده C هستند. اگر قرار دهیم

$$P = \frac{\sum_{h=1}^L A_h}{N}$$

نتیجه می‌شود $E(p_{st}) = P$ یعنی p_{st} در (۴۳.۴) برآوردکننده نااریب نسبت واحدهایی از جامعه است که در رده C هستند.

قضیه ۸.۴ در نمونه‌گیری تصادفی با طبقه‌بندی

$$V(p_{st}) = \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2(N_h - n_h)}{N_h - 1} \cdot \frac{P_h Q_h}{n_h} \quad (44.4)$$

برهان. این، حالتی خاص از قضیه کلی واریانس برآورد میانگین در نمونه‌گیری با طبقه‌بندی است. قبلاً دیدیم که

$$V(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{S_h^2}{n_h}$$

اگر Y_{hi} را در حالت خاص، وقتی در رده C است برابر ۱ و در غیر این صورت برابر ۰ بگیریم، بنابراین آنچه در نمونه‌گیری تصادفی برای تعیین نسبتها دیدیم

$$S_h^2 = \frac{N_h}{N_h - 1} P_h Q_h \quad (45.4)$$

که در آن $Q_h = 1 - P_h$ اگر (۴۵.۴) را در رابطه قبلی قرار دهیم، نتیجه می‌شود

$$\begin{aligned} V(p_{st}) &= \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{N_h}{N_h - 1} \frac{P_h Q_h}{n_h} \\ &= \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2(N_h - n_h)}{N_h - 1} \cdot \frac{P_h Q_h}{n_h} \end{aligned}$$

نمونه‌گیری با طبقه‌بندی برای برآورد نسبتها ۱۶۷

تبصره. در کاربردها، حتی وقتی که fpc قابل اغماض نباشد $\frac{N_h}{N_h-1}$ را برابر با ۱ می‌گیرند و فرمول واریانس را به صورت تقریبی زیر به کار می‌برند

$$\begin{aligned} V(p_{st}) &\simeq \frac{1}{N^2} \sum_{h=1}^L N_h(N_h - n_h) \frac{P_h Q_h}{n_h} \\ &= \sum_{h=1}^L \frac{N_h^2}{N^2} \left(1 - \frac{n_h}{N_h}\right) \frac{P_h Q_h}{n_h} \\ &= \sum_{h=1}^L W_h^2 (1 - f_h) \frac{P_h Q_h}{n_h} \end{aligned} \quad (۴۶.۴)$$

فرع ۱. اگر بتوانیم fpc را نادیده بگیریم، خواهیم داشت

$$V(p_{st}) \simeq \sum_{h=1}^L W_h^2 \frac{P_h Q_h}{n_h} \quad (۴۷.۴)$$

فرع ۲. اگر تخصیص، متناسب باشد در (۴۲.۴) برابری $\frac{N_h^2}{N^2} = \frac{n_h}{n}$ را منظور می‌کنیم. چون $\frac{N_h^2}{N^2} = \frac{n_h}{n}$ پس

$$\begin{aligned} V(p_{st}) &= \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2}{N_h - 1} (N - n) \frac{n_h}{n} \frac{P_h Q_h}{n_h} \\ &= \frac{N - n}{N} \cdot \frac{1}{n \cdot N} \sum_{h=1}^L \frac{N_h^2}{N_h - 1} P_h Q_h \end{aligned} \quad (۴۸.۴)$$

$$\simeq \frac{1 - f}{n} \sum_{h=1}^L W_h P_h Q_h \quad (۴۹.۴)$$

در برابری آخر $1 - N_h$ را به تقریب برابر با N_h فرض کرده‌ایم.

فرع ۳. چون P_h در عمل مجهول است، لذا نمی‌توان $V(p_{st})$ را به دست آورد. برای تعیین $\hat{V}(p_{st})$ همان‌طور که در نمونه‌گیری تصادفی ساده برای نسبتها دیدیم از برآورد $\frac{P_h Q_h}{n_h}$ که تقریباً برابر با $\frac{p_h q_h}{n_h - 1}$ است، استفاده می‌کنیم، بنابراین

$$\hat{V}(p_{st}) \simeq \frac{1}{N^2} \sum_{h=1}^L \frac{N_h^2 (N_h - n_h)}{N_h - 1} \frac{p_h q_h}{n_h - 1} \quad (۵۰.۴)$$

فرع ۴. بهترین انتخاب مقادیر n_h ها برای مینیم کردن $V(p_{st})$ از قضیه کلی (۶.۴) نتیجه می‌شود:

الف) اگر حجم نمونه ثابت باشد برای اینکه واریانس؛ مینیمم شود باید $n_h \propto N_h S_h$ باشد
ولی از (۴۵.۴) داریم

$$S_h = \sqrt{\frac{N_h}{N_h - 1} P_h Q_h}$$

پس

$$n_h \propto N_h \sqrt{N_h / (N_h - 1)} \cdot \sqrt{P_h Q_h} \simeq N_h \sqrt{P_h Q_h}$$

که در آن $N_h - 1 \simeq N_h$ فرض شده است. پس از رابطه

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

داریم

$$n_h \simeq n \cdot \frac{N_h \sqrt{P_h Q_h}}{\sum_{h=1}^L N_h \sqrt{P_h Q_h}} \quad (۵۱.۴)$$

ب) اگر هزینه $C = C_0 + \sum_{h=1}^L C_h n_h$ ثابت باشد، وقتی واریانس p_{st} مینیمم است؛
داشته باشیم

$$n_h = n \frac{N_h S_h / \sqrt{C_h}}{\sum N_h S_h / \sqrt{C_h}}$$

که اگر به جای S_h مقدار تقریبی آن $\sqrt{P_h Q_h}$ را قرار دهیم

$$n_h \simeq n \frac{N_h \sqrt{P_h Q_h / C_h}}{\sum_{h=1}^L N_h \sqrt{P_h Q_h / C_h}} \quad (۵۲.۴)$$

مقدار n از رابطه (۱۸.۴) با قرار دادن $\sqrt{P_h Q_h}$ به جای S_h به دست می‌آید

$$n = \frac{(C - C_0) \sum (N_h \sqrt{P_h Q_h / C_h})}{\sum (N_h \sqrt{P_h Q_h C_h})} \quad (۵۳.۴)$$

مثال ۴.۴ جامعه زیر با دو طبقه داده شده است. نمونه‌هایی با حجم $n_1 = 2$ و $n_2 = 2$ از هر طبقه انتخاب کنید و نشان دهید که $E(p_{st}) = P$.

نمونه‌گیری با طبقه‌بندی برای برآورد نسبتها ۱۶۹

| طبقه I | طبقه II |
|--------|---------|
| ۱ | ۱ |
| ۰ | ۰ |
| ۱ | ۰ |

می‌توانیم $9 = \binom{3}{2} \binom{2}{2}$ نمونه ممکن به حجم $n = n_1 + n_2 = 4$ از این جامعه استخراج کنیم. این ۹ نمونه ممکن در زیر فهرست شده‌اند

| I | II | p_1 | p_2 | p_{st} |
|------|------|-------|-------|----------|
| | ۱, ۰ | ۰٫۵ | ۰٫۵ | ۰٫۵ |
| ۱, ۰ | ۱, ۰ | | ۰٫۵ | ۰٫۵ |
| | ۰, ۰ | | ۰ | ۰٫۲۵ |
| | ۱, ۰ | ۱ | ۰٫۵ | ۰٫۷۵ |
| ۱, ۱ | ۱, ۰ | | ۰٫۵ | ۰٫۷۵ |
| | ۰, ۰ | | ۰ | ۰٫۵ |
| | ۱, ۰ | ۰٫۵ | ۰٫۵ | ۰٫۵ |
| ۰, ۱ | ۱, ۰ | | ۰٫۵ | ۰٫۵ |
| | ۰, ۰ | | ۰ | ۰٫۲۵ |

تشخیص چگونگی تشکیل جدول به عهده خواننده است.

$$E(p_{st}) = \frac{1}{9} (0.5 + 0.5 + \dots + 0.25) = \frac{4.5}{9} = 0.5$$

از طرفی با توجه به جامعه اصلی

$$P = \frac{3}{6} = 0.5$$

پس

$$E(p_{st}) = P$$

Δ

مثال ۵.۴ یک بررسی مقدماتی در سه شهر کوچک انجام داده‌اند و نسبت خانواده‌هایی را که دو فرزند یا بیشتر دارند برآورد کرده‌اند. از داده‌های حاصل از بررسی مقدماتی استفاده کنید و تعیین کنید که برای برآورد p_{st} با تخصیص نینمن، وقتی حجم کل نمونه n است، حجم نمونه هر شهر چقدر باید باشد؟

| شهر | خانواده | p_h | $p_h q_h$ | $\sqrt{p_h q_h}$ | $N \cdot \sqrt{p_h q_h}$ |
|-------|---------|-------|-----------|------------------|--------------------------|
| A | ۲۰۰۰ | ۰٫۱۰ | ۰٫۰۹ | ۰٫۳ | ۶۰۰ |
| B | ۳۰۰۰ | ۰٫۱۵ | ۰٫۱۲۷۵ | ۰٫۳۵ | ۱۰۵۰ |
| C | ۵۰۰۰ | ۰٫۲۰ | ۰٫۱۶ | ۰٫۴ | ۲۰۰۰ |
| مجموع | ۱۰۰۰۰ | | | | ۳۶۵۰ |

از فرمول نینم داریم

$$\hat{n}_h \simeq n \cdot \frac{N_h \sqrt{p_h q_h}}{\sum N_h \sqrt{p_h q_h}}$$

با توجه به محاسبات جدول بالا که، سه ستون اول آن داده‌های حاصل از بررسی است، داریم

$$\hat{n}_1 = \frac{600}{3650} n = \frac{12}{73} n$$

$$\hat{n}_2 = \frac{1050}{3650} n = \frac{21}{73} n$$

$$\hat{n}_3 = \frac{2000}{3650} n = \frac{40}{73} n$$

۱۲.۴ اثر انحراف از تخصیص ایتیم

در این بخش کاهش دقت نتیجه نمونه‌گیری را به دلیل عدم توفیق در دستیابی به تخصیص ایتیم مورد بحث قرار می‌دهیم.

فرض کنید که قصد داریم برای مقداری مفروض از n ، از تخصیص نینم استفاده کنیم. حجم نمونه‌ای n'_h در طبقه h باید به صورت زیر باشد

$$n'_h = \frac{n(W_h S_h)}{\sum W_h S_h} \quad (54.4)$$

می‌دانیم که واریانس نینم حاصل از این تخصیص چنین است

$$V_{min}(\bar{Y}_{st}) = \frac{(\sum W_h S_h)^2}{n} - \frac{\sum W_h S_h^2}{N} \quad (55.4)$$

در عمل، چون مقادیر S_h ها معلوم نیستند فقط این تخصیص را با تقریب می‌توان انجام داد. اگر \hat{n}_h حجم نمونه‌ای مورد استفاده در طبقه h باشد مقدار واریانس که حاصل می‌شود به صورت

زیر است

$$V(\bar{Y}_{st}) = \sum \frac{W_h^2 S_h^2}{\hat{n}_h} - \frac{\sum W_h S_h^2}{N}$$

افزایش واریانس ناشی از تخصیص ناکامل عبارت است از

$$V(\bar{Y}_{st}) - V_{\min}(\bar{Y}_{st}) = \sum \frac{W_h^2 S_h^2}{\hat{n}_h} - \frac{1}{n} \left(\sum W_h S_h \right)^2 \quad (56.4)$$

از رابطه (56.4) داریم

$$W_h S_h = \frac{n'_h}{n} \sum W_h S_h$$

اگر این مقدار را در (56.4) قرار دهیم، نتیجه زیر حاصل می‌شود

$$\begin{aligned} V(\bar{Y}_{st}) - V_{\min}(\bar{Y}_{st}) &= \frac{(\sum W_h S_h)^2}{n^2} \left[\sum \frac{n_h'^2}{\hat{n}_h} - n \right] \\ &= \frac{(\sum W_h S_h)^2}{n^2} \sum \frac{(\hat{n}_h - n_h')^2}{\hat{n}_h} \end{aligned} \quad (57.4)$$

اگر در رابطه (55.4)، $N \rightarrow \infty$ ، آن‌گاه این رابطه به صورت زیر درمی‌آید

$$V_{\min}(\bar{Y}_{st}) = \frac{(\sum W_h S_h)^2}{n} \quad (58.4)$$

از تقسیم روابط (57.4) و (58.4) نتیجه می‌شود

$$\frac{V(\bar{Y}_{st}) - V_{\min}(\bar{Y}_{st})}{V_{\min}(\bar{Y}_{st})} = \frac{1}{n} \sum \frac{(\hat{n}_h - n_h')^2}{\hat{n}_h} \quad (59.4)$$

این مقدار، افزایش نسبی واریانس را، وقتی تخصیص به صورت ناکامل انجام می‌شود نشان می‌دهد. در این رابطه \hat{n}_h مقدار حجم نمونه در طبقه h ام است که به جای مقدار n_h' که مربوط به تخصیص اپتیمم است مورد استفاده واقع شده است. اگر fpc قابل چشمپوشی نباشد، در رابطه (59.4) نماد (=) به نماد « \geq » تبدیل می‌شود.

اگر $g_h = |\hat{n}_h - n_h'| / \hat{n}_h$ ، به صورت زیر درمی‌آید

$$\frac{V - V_{\min}}{V_{\min}} = \sum_{h=1}^L \frac{\hat{n}_h}{n} g_h^2$$

که در واقع میانگین موزون g_h^2 است. بنابراین، حد بالایی محافظه‌کارانه $\frac{V - V_{\min}}{V_{\min}}$ برابر g^2 است که در آن g ، بزرگترین تفاوت نسبی در هر طبقه است.

مثلاً اگر $g = 0.2$ یا 20% باشد، نمونه‌گیری واریانس نمی‌تواند از $(0.2)^2$ تجاوز کند. اگر $g = 0.3$ یا 30% باشد، نمونه‌گیری واریانس حداکثر 9% است.

۱۳.۴ مسأله تخصیص، برای بیش از یک صفت

گاهی اوقات در یک نمونه‌گیری، چند صفت از واحدهای نمونه به‌طور همزمان مشاهده و اندازه‌گیری می‌شوند. چون بهترین تخصیص مقادیر n_h برای یک صفت الزاماً بهترین تخصیص برای صفت دیگر نیست (زیرا مقادیر s_h صفتها با هم متفاوت‌اند) بنابراین معلوم نیست بین مقادیر مختلف n_h که متناظر با صفت‌های مختلف برای هر طبقه به‌دست می‌آیند کدام n_h را باید برگزید، و نمونه‌گیری از آن طبقه را با آن حجم نمونه‌ای انجام داد. در این مورد اولین گامی که باید برداشت این است که بین صفت‌های مختلف، صفتی را که از نظر بررسی نمونه‌ای در درجه اول اهمیت‌اند انتخاب کنیم و بقیه را کنار بگذاریم. این کار موجب تقلیل صفات برای تعیین مقادیر n_h می‌شود. اگر داده‌های مقدماتی خوبی موجود باشند، آن‌گاه می‌توانیم تخصیص اپتیمم را برای هر یک از صفات بیابیم و ببینیم که در هر طبقه میزان ناهماهنگی حجم‌های حاصل از صفات چقدر است. اگر بین صفات، میزان همبستگی زیاد باشد اختلاف این حجمها نسبتاً کم خواهد بود. در هر حال بین این حجمها در هر طبقه باید با مصالحه یکی را انتخاب کرد. مثال زیر نحوه عمل را مشخص می‌کند.

مثال ۶.۴ داده‌های این مثال را جنس^۱ در یک بررسی ارائه داده است. ایالت آیووا از نظر جغرافیایی به ۵ منطقه تقسیم شد و سازمانهای عمده کشاورزی در هر منطقه مشخص شد. این ۵ منطقه به‌عنوان ۵ طبقه در یک بررسی از صنعت شیر به‌کار رفت. سه قلم از مهمترین اقلام تحت بررسی به‌صورت (۱) تعداد گاوهایی که در روز دوشیده شده‌اند، (۲) تعداد گالنه‌های شیر در روز، و (۳) کل درآمد روزانه از محصولات شیر مشخص شدند. از یک بررسی در چند سال قبل مقادیر s_h طبقات به‌صورتی بوده‌اند که در جدول ۱ نشان داده‌ایم. در جدول ۲، تخصیص‌های اپتیمم نیم بر پایه این s_h ها برای هر قلم، در نمونه‌ای به حجم ۱۰۰۰ سازمان داده شده است.

جدول ۱. وزنها و مقادیر s_h طبقات

| طبقه | $W_h = \frac{N_h}{N}$ | s_h | s_h | s_h |
|------|-----------------------|------------------|--------------------|-------------------------------|
| | | برای گاوهای شیره | برای گالنه‌های شیر | برای درآمد روزانه محصولات شیر |
| ۱ | ۰.۱۹۷ | ۴.۶ | ۱۱.۷ | ۳۳۲ |
| ۲ | ۰.۱۹۱ | ۳.۴ | ۹.۸ | ۳۵۷ |
| ۳ | ۰.۲۱۹ | ۳.۳ | ۷.۰ | ۲۴۶ |
| ۴ | ۰.۱۸۴ | ۲.۸ | ۶.۵ | ۱۷۳ |
| ۵ | ۰.۲۰۸ | ۳.۷ | ۹.۸ | ۲۷۹ |

مسأله تخصیص، برای بیش از یک صفت ۱۷۳

جدول ۲. مفادیر n_h ایشیم متناظر با سه صفت و مفادیر n_h تخصیص متناسب

| طبقه | تخصیص | | | متوسط n_h |
|------|--------|-------|--------|-------------|
| | متناسب | گاوها | گالنها | |
| ۱ | ۱۹۷ | ۲۵۲ | ۲۵۸ | ۲۳۶ |
| ۲ | ۱۹۱ | ۱۸۲ | ۲۰۹ | ۲۴۶ |
| ۳ | ۲۱۹ | ۲۰۳ | ۱۷۱ | ۱۹۲ |
| ۴ | ۱۸۲ | ۱۴۵ | ۱۳۲ | ۱۱۵ |
| ۵ | ۲۰۸ | ۲۱۶ | ۲۲۸ | ۲۰۹ |

به طوری که از جدول ۲ دیده می شود تخصیصهای ایشیم در ستونهای ۲ و ۳ و ۴ تفاوتی نسبتاً کمی دارند و تفاوت تمام آنها با تخصیص متناسب در یک جهت است. مثلاً در طبقه اول، با تخصیص متناسب، باید از ۱۹۷ سازمان نمونه گیری کرد در حالی که حجم نمونه برای طبقه اول با تخصیص ایشیم از ۲۳۶ تا ۲۵۸ است. در ستون آخر متوسط تقریبی سه حجم نمونه ای ایشیم را از طبقه اول برابر با ۲۵۰ نوشته ایم که مصالحه ای بین سه حجم نمونه ای حاصل از تخصیص ایشیم برای سه صفت است. در جدول ۳، واریانسهای نمونه ای مورد انتظار V_{exp} را برای تخصیصهای ایشیم سه گانه و برای حالت مصالحه و تخصیص متناسب آورده ایم. فرمولهایی که برای محاسبه به کار برده ایم به صورت زیرند

$$V_{exp} = \frac{(\sum W_h S_h)^2}{n}$$

$$*V_{com} = \sum \frac{(W_h S_h)^2}{n_h}$$

$$V_{prop} = \frac{\sum W_h S_h^2}{n}$$

اگر ممکن می بود که سه بار، هر بار برای یک صفت با تخصیص ایشیم، عمل نمونه گیری را انجام داد، دقت حاصل نزدیک دقت تخصیص حالت مصالحه می شد. آنچه جالب است این است که دقت تخصیص متناسب کمی کمتر از دقت تخصیص حالت مصالحه است.

جدول ۳. واریانسهای مورد انتظار برآورد می کنیم

| بوع تخصیص | گاوها | گالنها | درآمدها |
|-----------|-------|--------|---------|
| ایشیم | ۰٫۱۲۷ | ۰٫۸۱۰ | ۷۶٫۸ |
| مصالحه ای | ۰٫۱۲۸ | ۰٫۸۰۲ | ۷۷٫۶ |
| متناسب | ۰٫۱۳۱ | ۰٫۸۳۷ | ۸۰٫۹ |

۱۴.۴ ساختن طبقات

در این بخش، سؤالهایی را به صورت زیر مطرح می‌کنیم. بهترین مشخصه برای ساختن طبقات چیست؟ چگونه می‌توان کرانه‌های طبقات را تعیین کرد؟

واضح است که بهترین مشخصه متغیر Y ، توزیع فراوانی آن است. بعد از آن توزیع فراوانی متغیری کمکی است که همبستگی بسیار قوی با متغیر Y دارد. دالنیوس^۱، با فرض معلوم بودن تعداد طبقات، برای تعیین بهترین کرانه‌های طبقات تحت تخصیص متناسب و تخصیص نیم، معادلاتی را ارائه داده است. روشهای تقریبی سریعتر دیگری نیز به وسیله پژوهشگران دیگر عرضه شده‌اند. ما در اینجا تخصیص نیم را به دلیل کاربرد فراوان آن، و به دلیل اینکه عموماً بهتر از تخصیص متناسب عمل می‌کند در نظر می‌گیریم. بدو می‌پذیریم که طبقات با استفاده از مقادیر متغیر Y ساخته می‌شوند. ابتدا راهی نظری و سپس طریقی کاربردی را بررسی می‌کنیم.

فرض کنید Y_0 کوچکترین و Y_L بزرگترین مقدار Y در جامعه باشند. مسأله مورد نظر، یافتن کرانه‌های طبقاتی یعنی مقادیر Y_1, Y_2, \dots, Y_{L-1} هستند به قسمی که L طبقه به وجود آید و ضمناً واریانس برآوردکننده میانگین جامعه، یعنی

$$V(\bar{Y}_{st}) = \frac{1}{n} \left(\sum_{h=1}^L W_h S_h \right)^2 - \frac{1}{N} \sum_{h=1}^L W_h S_h^2$$

مینیمم شود. اگر $\frac{1}{N}$ قابل اغماض باشد می‌توان از جمله دوم طرف دوم صرف نظر کرد، لذا کافی است که Y_1, Y_2, \dots, Y_{L-1} را که مشخص‌کننده کران طبقه‌ها هستند به قسمی بیابیم که $\sum W_h S_h$ مینیمم شود. در این مجموع، چون Y_h در بازه‌های $[Y_{h-1}, Y_h]$ و $[Y_h, Y_{h+1}]$ ظاهر می‌شود لذا در جمله‌های $W_h S_h$ و $W_{h+1} S_{h+1}$ ظاهر خواهد شد. این دوبازه به ترتیب طبقه‌های h ام و $h+1$ ام را مشخص می‌کنند. با این توضیح داریم

$$\frac{\partial}{\partial Y_h} \left(\sum W_h S_h \right) = \frac{\partial}{\partial Y_h} (W_h S_h) + \frac{\partial}{\partial Y_h} (W_{h+1} S_{h+1}) \quad (۶۰.۴)$$

سعی می‌کنیم مقادیر دو جمله عبارت سمت راست (۶۰.۴) را حساب کنیم. برای این منظور فرض می‌کنیم $f_Y(y)$ تابع چگالی یا تابع فراوانی نسبی Y ها باشد. اگر t معرف Y هایی باشد که در طبقه h ام قرار دارند آن‌گاه وزن طبقه h ام، همان احتمال مربوط به این طبقه است، یعنی

$$W_h = \int_{Y_{h-1}}^{Y_h} f(t) dt = F(Y_h) - F(Y_{h-1})$$

که $F(t)$ تابع توزیع متناظر با $f(t)$ است، لذا

$$\frac{\partial W_h}{\partial Y_h} = f(Y_h) \quad (۶۱.۴)$$

میانگین Y های بازه (Y_{h-1}, Y_h) ، یعنی میانگین Y های واقع در طبقه h ام، با توجه به اینکه تابع چگالی در این بازه $f(t)/W_h$ است عبارت است از

$$\mu_h = \int_{Y_{h-1}}^{Y_h} t \frac{f(t)}{W_h} dt$$

چون t معرف Y های طبقه h ام، و $\frac{f(t)}{W_h}$ معرف تابع چگالی متناظر با این طبقه است، پس S_h^2 که می توان آن را تقریباً متناظر با واریانس t دانست به صورت زیر است

$$\begin{aligned} S_h^2 &= E(t - E(t))^2 = E(t^2) - [E(t)]^2 \\ &= \int_{Y_{h-1}}^{Y_h} t^2 \frac{f(t)}{W_h} dt - \left[\int_{Y_{h-1}}^{Y_h} t \frac{f(t)}{W_h} dt \right]^2 \end{aligned}$$

بنابراین

$$W_h S_h^2 = \int_{Y_{h-1}}^{Y_h} t^2 f(t) dt - \frac{\left[\int_{Y_{h-1}}^{Y_h} t f(t) dt \right]^2}{W_h}$$

با یادآوری دستور مشتقگیری از انتگرالها نسبت به پارامتری که در حدود انتگرال و عبارت زیر انتگرال وجود دارد، یعنی با یادآوری فرمول

$$\frac{\partial}{\partial \theta} \int_{\alpha(\theta)}^{\beta(\theta)} f(x, \theta) dx = f[\beta(\theta), \theta] \cdot \beta'(\theta) - f[\alpha(\theta), \theta] \alpha'(\theta) + \int_{\alpha(\theta)}^{\beta(\theta)} \frac{\partial}{\partial \theta} f(x, \theta) dx$$

از طرفین رابطه بالا به نسبت به Y_h مشتق می گیریم

$$\begin{aligned} S_h^2 \frac{\partial W_h}{\partial Y_h} + 2W_h S_h \cdot \frac{\partial S_h}{\partial Y_h} \\ = Y_h^2 f(Y_h) - \frac{2 \left[\int_{Y_{h-1}}^{Y_h} t f(t) dt \right] Y_h f(Y_h) \cdot W_h - \frac{\partial W_h}{\partial Y_h} \left[\int_{Y_{h-1}}^{Y_h} t f(t) dt \right]^2}{W_h^2} \\ = Y_h^2 f(Y_h) - 2Y_h f(Y_h) \frac{\int_{Y_{h-1}}^{Y_h} t f(t) dt}{W_h} + \left[\int_{Y_{h-1}}^{Y_h} t f(t) dt / W_h \right]^2 \cdot f(Y_h) \end{aligned}$$

و یا با توجه به مقدار μ_h

$$S_h^2 \frac{\partial W_h}{\partial Y_h} + 2W_h S_h \frac{\partial S_h}{\partial Y_h} = Y_h^2 f(Y_h) - 2Y_h f(Y_h) \mu_h + f(Y_h) \cdot \mu_h^2 \quad (۶۲.۴)$$

با توجه به (۶۱.۲)، به در طرف رابطه (۶۲.۲) مقدار ثابت $S_h^r \frac{\partial W_h}{\partial Y_h} = S_h^r f(Y_h)$ را نظیر به نظیر اضافه می‌کنیم. نتیجه می‌شود

$$2S_h^r \frac{\partial W}{\partial Y_h} + 2W_h S_h \frac{\partial S_h}{\partial Y_h} = Y_h^r f(Y_h) - 2Y_h \mu_h f(Y_h) + \mu_h^r f(Y_h) + S_h^r f(Y_h)$$

اگر طرفین را بر $2S_h$ تقسیم کنیم پس از خلاصه کردن، داریم

$$S_h \frac{\partial W}{\partial Y_h} + W_h \frac{\partial S_h}{\partial Y_h} = \frac{1}{2S_h} f(Y_h) [Y_h^r - 2Y_h \mu_h + \mu_h^r + S_h^r]$$

یا

$$\frac{\partial(W_h S_h)}{\partial Y_h} = \frac{1}{2} f(Y_h) \frac{(Y_h - \mu_h)^r + S_h^r}{S_h} \quad (63.2)$$

به طریقی مشابه می‌توان نتیجه گرفت که

$$\frac{\partial(W_{h+1} S_{h+1})}{\partial Y_h} = -\frac{1}{2} f(Y_h) \frac{(Y_h - \mu_{h+1})^r + S_{h+1}^r}{S_{h+1}} \quad (64.2)$$

حال اگر بخواهیم که $V(\bar{Y}_{.t})$ مینیم شود باید با توجه به (۶۰.۴) داشته باشیم

$$\frac{\partial}{\partial Y_h} (W_h S_h) + \frac{\partial}{\partial Y_h} (W_{h+1} S_{h+1}) = 0$$

یا

$$\frac{\partial}{\partial Y_h} (W_h S_h) = -\frac{\partial}{\partial Y_h} (W_{h+1} S_{h+1})$$

بنابر رابط (۶۳.۲) و (۶۴.۲) نتیجه می‌شود که

$$\frac{(Y_h - \mu_h)^r + S_h^r}{S_h} = \frac{(Y_h - \mu_{h+1})^r + S_{h+1}^r}{S_{h+1}} \quad (h = 1, 2, \dots, L-1)$$

بدین طریق دستگامی با $L-1$ معادله به دست می‌آید که در صورت معلوم بودن برآوردهایی مقدماتی از $S_h, S_{h+1}, \dots, S_{L-1}$ و $\mu_h, \mu_{h+1}, \dots, \mu_{L-1}$ را می‌توان یافت. متأسفانه این معادلات به دلیل وابستگی S_h و μ_h به Y_h کاربرد عملی ندارند. دالنیوس و هاجز روشی تقریبی و سریع برای مینیم کردن $\sum W_h S_h$ به شرح زیر ارائه کرده‌اند. قرار می‌دهیم

$$Z_h = Z(Y_h) = \int_{Y_0}^{Y_h} \sqrt{f(t)} dt \quad (65.2)$$

اگر تعداد طبقات زیاد بوده و عرض آنها کم باشد $f(y)$ را می توان تقریباً در هر طبقه ثابت گرفت، یعنی توزیع Y را در هر طبقه توزیع یکنواخت فرض کرد. بنابراین

$$W_h = \int_{Y_{h-1}}^{Y_h} f(t) dt \simeq f_h(Y_h - Y_{h-1}) \quad (66.4)$$

از طرفی می دانیم که اگر متغیر Y روی فاصله (a, b) دارای توزیع یکنواخت باشد آن گاه

$$V(Y) = \frac{(b-a)^2}{12}$$

چون در اینجا فرض کرده ایم که Y روی فاصله (Y_{h-1}, Y_h) توزیع یکنواخت دارد. لذا به طور تقریبی

$$S_h^2 \simeq \frac{(Y_h - Y_{h-1})^2}{12}$$

ولذا

$$S_h \simeq \frac{(Y_h - Y_{h-1})}{\sqrt{12}} \quad (67.4)$$

با توجه به (65.4) می توان نوشت

$$Z_h = \int_{Y_0}^{Y_h} \sqrt{f(t)} dt, \quad Z_{h-1} = \int_{Y_0}^{Y_{h-1}} \sqrt{f(t)} dt$$

پس

$$Z_h - Z_{h-1} = \int_{Y_{h-1}}^{Y_h} \sqrt{f(t)} dt$$

چون توزیع در فاصله (Y_{h-1}, Y_h) یکنواخت است

$$Z_h - Z_{h-1} \simeq \sqrt{f_h}(Y_h - Y_{h-1}) \quad (68.4)$$

در این رابطه و در (66.4)، f_h مقدار ثابت $f(Y)$ در طبقه h ام است. از (68.4) داریم

$$\begin{aligned} \sum_{h=1}^L (Z_h - Z_{h-1})^2 &\simeq \sum_{h=1}^L f_h (Y_h - Y_{h-1})^2 \\ &= \sum_{h=1}^L f_h (Y_h - Y_{h-1})(Y_h - Y_{h-1}) \end{aligned}$$

با توجه به (۶۶.۴)

$$\sum_{h=1}^L (Z_h - Z_{h-1})^2 \simeq \sum_{h=1}^L W_h (Y_h - Y_{h-1})$$

اگر به جای $Y_h - Y_{h-1}$ مقدار آن را از (۶۷.۴) قرار دهیم نتیجه می شود که

$$\sum_{h=1}^L (Z_h - Z_{h-1})^2 \simeq \sqrt{12} \sum_{h=1}^L W_h S_h = \sum_{h=1}^L f_h (Y_h - Y_{h-1})^2 \quad (۶۹.۴)$$

برای اینکه $\sum W_h S_h$ مینیم شود باید مقدار

$$\sum_{h=1}^L (Z_h - Z_{h-1})^2$$

مینیم شود. مجموع h فاصله طبقاتی مقداری ثابت است، زیرا

$$\sum_{h=1}^L (Z_h - Z_{h-1}) = Z_L - Z_0$$

اینک h فاصله به صورت $Z_h - Z_{h-1}$ داریم که طول هر کدام مثبت و مجموع آنها مقداری ثابت است، بنابراین $\sum_{h=1}^L (Z_h - Z_{h-1})^2$ وقتی مینیم است که $Z_h - Z_{h-1}$ ، یعنی مقدار

$$\sqrt{f_h} (Y_h - Y_{h-1}), \quad h = 1, 2, \dots, L$$

به ازای هر مقدار h ، ثابت باشد. لذا اگر فاصله های طبقه ها را به صورتی اختیار کنیم که این مقادیر در هر طبقه تقریباً همانند باشند، کرانه های طبقه ها به گونه ای خواهند بود که واریانس برآوردکننده میانگین جامعه مینیم باشد.

با توضیحات بالا دالنیوس و هاجز قاعده عملی زیر را ارائه می دهند:

اگر $f(Y)$ معلوم باشد، مقادیر تجمعی $\sqrt{f(Y)}$ را محاسبه می کنیم. سپس Y_h ها را به قسمی می یابیم که بر دامنه مقادیر $\sqrt{f(Y)}$ فاصله های برابر ایجاد کنند. باید توجه داشت که برای تعیین کرانه ضروری است که f یعنی تابع توزیع فراوانی جامعه معلوم باشد. مثال زیر نحوه استفاده از این قاعده را نشان می دهد.

مثال ۷.۴ داده های جدول ۱، توزیع فراوانی درصد وامهای بانکی را در ۱۳۴۳۸ شعبه بانک کشور آمریکا نشان می دهد که به وامهای صنعتی تخصیص یافته اند. توزیع چاوله است و مد آن در دم چپ قرار دارد. در ستون مربوط به مقادیر تجمعی \sqrt{f} ، به عنوان مثال $\sqrt{۳۴۳۶} = ۵۸٫۹$ و $\sqrt{۳۴۶۴} + \sqrt{۲۵۱۶} = ۱۰۹٫۱$ و نظایر آن. فرض کنید مایلیم ۵ طبقه داشته باشیم. چون

مجموع مقادیر تجمعی \sqrt{f} برابر با ۳۹۵٫۵ است، نقاط تقسیم باید در ۷۷٫۹، ۱۵۵٫۸، ۲۳۳٫۷ و ۳۱۱٫۶ باشند. نزدیکترین نقاط موجود در جدول ۲ آمده‌اند.

جدول ۱. محاسبه کرانه‌های طبقات به وسیله قاعده \sqrt{f} تجمعی*

| مقدار تجمعی $\sqrt{f(t)}$ | $f(Y)$ | وامهای صنعتی کل وامها | مقدار تجمعی $\sqrt{f(Y)}$ | $f(Y)$ | وامهای صنعتی کل وامها |
|---------------------------|--------|--------------------------|---------------------------|--------|--------------------------|
| ۳۴۰٫۳ | ۱۲۶ | ۵۰-۵۵ | ۵۸٫۹ | ۳۴۶۴ | ۰-۵ |
| ۳۵۰٫۶ | ۱۰۷ | ۵۵-۶۰ | ۱۰۹٫۱ | ۲۵۱۶ | ۵-۱۰ |
| ۳۵۹٫۷ | ۸۲ | ۶۰-۶۵ | ۱۵۵٫۵ | ۲۱۵۷ | ۱۰-۱۵ |
| ۳۶۶٫۸ | ۵۰ | ۶۵-۷۰ | ۱۹۵٫۳ | ۱۵۸۱ | ۱۵-۲۰ |
| ۳۷۳٫۰ | ۳۹ | ۷۰-۷۵ | ۲۲۹٫۱ | ۱۱۴۲ | ۲۰-۲۵ |
| ۳۷۸٫۰ | ۲۵ | ۷۵-۸۰ | ۲۵۶٫۴ | ۷۴۹ | ۲۵-۳۰ |
| ۳۸۲٫۰ | ۱۶ | ۸۰-۸۵ | ۲۷۹٫۰ | ۵۱۲ | ۳۰-۳۵ |
| ۳۸۶٫۴ | ۱۹ | ۸۵-۹۰ | ۲۹۸٫۴ | ۳۷۶ | ۳۵-۴۰ |
| ۳۸۷٫۸ | ۲ | ۹۰-۹۵ | ۳۱۴٫۷ | ۲۶۵ | ۴۰-۴۵ |
| ۳۸۹٫۵ | ۳ | ۹۵-۱۰۰ | ۳۲۹٫۱ | ۲۰۷ | ۴۵-۵۰ |

* این مثال از مرجع [۳] اقتباس شده است.

جدول ۲. نزدیکترین نقاط موجود به نقاط تقسیم

| طبقه | کرنه‌ها | | | | |
|--------------------------------|---------|-------|--------|--------|---------|
| | ۱ | ۲ | ۳ | ۴ | ۵ |
| بازه مربوط به \sqrt{f} تجمعی | ۰-۵٪ | ۵-۱۵٪ | ۱۵-۲۵٪ | ۲۵-۴۵٪ | ۴۵-۱۰۰٪ |
| | ۵۸٫۹ | ۹۶٫۶ | ۷۳٫۶ | ۸۵٫۶ | ۷۴٫۸ |

توجه. رابطه (۶۹٫۴) یک پیامد فوری و جالب دارد: اگر $W_h S_h$ ثابت باشد، تخصیص نین، حجم نمونه‌ای ثابت $n_h = \frac{n}{L}$ را برای همه طبقات به دست می‌دهد. در مقایسه با روشهای دیگر، تخصیص $n_h = \frac{n}{L}$ در این حالت روشی رضایتبخش است. برای مطالعه بیشتر در زمینه ساختن طبقه‌ها، خواننده را به مرجع [۳] ارجاع می‌دهیم.

۱۵.۴ طبقه‌بندی بعد از انتخاب نمونه

گاهی اوقات وقتی مایلیم طبقه‌بندی را برحسب متغیری کلیدی انجام دهیم با این مشکل روبه‌رو هستیم که نمی‌توانیم واحدها را تا بعد از انتخاب در طبقه صحیح آنها قرار دهیم. مثلاً اگر بخواهیم نتیجه نظرخواهی خاصی را برحسب جنس نظردهنده طبقه‌بندی کنیم و نظرخواهی را با مصاحبه تلفنی انجام دهیم قبل از تماس با فرد نمی‌توانیم پاسخگو را در طبقه زنان یا مردان منظور کنیم.

همین‌طور اگر حسابرسی بخواهد حسابهایی را که با نمونه‌گیری برای حسابرسی انتخاب می‌کند در دو طبقه عمده فروش و خرده فروش قرار دهد تا بعد از انتخاب حساب و بررسی آن نمی‌تواند این طبقه‌بندی را انجام دهد.

فرض کنید نمونه تصادفی ساده‌ای به حجم n برای نظرخواهی انتخاب شود. پاسخگویان نمونه را پس از انتخاب می‌توان به دو قسمت n_1 مرد و n_2 زن تقسیم کرد. اگر بخواهیم \bar{Y}_N جامعه را به وسیله \bar{Y}_1 برآورد کنیم باید $\frac{N_1}{N}$ و $\frac{N_2}{N}$ را بدانیم. توجه کنید که در این حالت n_1 و n_2 دو متغیر تصادفی‌اند و از قبل نمی‌دانیم که در n تماس تلفنی با پاسخگویان چند پاسخگوی زن و چند پاسخگوی مرد خواهیم داشت. البته مجموع این دو متغیر تصادفی، مساوی مقدار ثابت n است. این نمونه‌گیری گرچه به ظاهر نمونه‌گیری با طبقه‌بندی است ولی در واقع به دلیل تصادفی بودن مقادیر n_1 و n_2 با تعریف نمونه‌گیری با طبقه‌بندی مطابقت نمی‌کند. می‌توان نشان داد که اگر $\frac{N_1}{N}$ معلوم و برای هر طبقه $20 \leq n_{1i}$ ، آنگاه طبقه‌بندی پس از انتخاب نمونه که گاهی آن را طبقه‌بندی پسین می‌گویند تقریباً همان دقت نمونه‌گیری تصادفی با طبقه‌بندی در تخصیص متناسب را دارد [۱۵]. طبقه‌بندی پس از انتخاب نمونه، غالباً وقتی نمونه تصادفی ساده دقیقاً برحسب گروه‌بندیهای جامعه‌ای حالت تعادلی ندارد، انتخابی مناسب است. فرض کنید مثلاً نمونه تصادفی ساده‌ای به حجم $n = 100$ از جامعه‌ای انتخاب کرده‌ایم که باید به‌طور برابر بین زن و مرد تقسیم شود. صفت مورد نظر، وزن پاسخگویان به کیلوگرم است و هدف نمونه‌گیری برآورد کردن متوسط وزن در جامعه است. نمونه منتخب اطلاعات زیر را فراهم کرده است:

| مرد | زن |
|--|--------------------------|
| $n_1 = 20$ | $n_2 = 80$ |
| $\bar{y}_1 = 90$ کیلوگرم | $\bar{y}_2 = 55$ کیلوگرم |
| $\bar{Y} = \frac{20}{100}(90) + \frac{80}{100}(55) = 62$ | |

با توجه به اینکه تعداد مردان پاسخگو نسبت به تعداد زنان پاسخگو کم است به نظر می‌رسد که برآورد $\bar{Y} = 62$ برآوردی کمتر از واقعیت است. می‌توانیم این برآورد را به صورت زیر تصحیح کنیم

$$\bar{Y}_{st} = \frac{N_1}{N} \bar{y}_1 + \frac{N_2}{N} \bar{y}_2 = \frac{1}{4}(90) + \frac{3}{4}(55) = 72.5$$

این برآورد به نظر بیشتر واقعی می‌آید زیرا در این محاسبه به زنان و مردان وزنه‌ای برابر نسبت داده شده است. توجه دارید که گرچه N_1 ، N_2 و N نامعلوم‌اند ولی از قبل فرض شده است که نسبت دو جنس برابرند. مسلماً این \bar{Y}_{st} واریانس دارد که با واریانس مربوط به \bar{Y}_{st} در نمونه‌گیری تصادفی با طبقه‌بندی تفاوت دارد، زیرا نمونه‌گیری با طبقه‌بندی، قبل از انتخاب نمونه طراحی نشده است. اما، واریانس تقریبی را می‌توان به صورت زیر به دست آورد.

مثل سابق $W_h = \frac{N_h}{N}$ ، $h = 1, \dots, L$ ، ولی n_h ها تصادفی اند و داریم

$$\begin{cases} E(n_h) = nW_h & h = 1, 2, \dots, L \\ V(n_h) = nW_h(1 - W_h) \end{cases}$$

اگر n_h ثابت می‌بود

$$\begin{aligned} \hat{V}(\bar{Y}_{st}) &= \sum_{h=1}^L W_h^2 \frac{N_h - n_h}{N_h} \frac{s_h^2}{n_h} \\ &= \sum_{h=1}^L W_h^2 \frac{s_h^2}{n_h} - \frac{1}{N} \sum_{h=1}^L W_h s_h^2 \end{aligned} \quad (70.4)$$

اما n_h تصادفی است. در این وضعیت یک تقریب کلی از $\hat{V}(\bar{Y}_{st})$ بدین طریق به دست می‌آید که به جای $\frac{1}{n_h}$ مقدار امید آن را قرار دهیم. متأسفانه تعیین امید عکس یک متغیر تصادفی در حالت کلی ساده نیست، اما می‌توان نشان داد که تقریبی خوب به صورت زیر است

$$E\left(\frac{1}{n_h}\right) \approx \frac{1}{nW_h} + \frac{1 - W_h}{n^2 W_h^2} \quad (71.4)$$

اگر (71.4) را در (70.4) قرار دهیم، نتیجه می‌شود

$$\begin{aligned} \hat{V}_P(\bar{Y}_{st}) &= \frac{1}{n} \sum_{h=1}^L W_h s_h^2 + \frac{1}{n^2} (1 - W_h) s_h^2 - \frac{1}{N} \sum_{h=1}^L W_h s_h^2 \\ &= \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{h=1}^L W_h s_h^2 + \frac{1}{n^2} \sum_{h=1}^L (1 - W_h) s_h^2 \end{aligned}$$

و سرانجام

$$\hat{V}_P(\bar{Y}_{st}) = \frac{N - n}{Nn} \sum_{h=1}^L W_h s_h^2 + \frac{1}{n^2} \sum_{h=1}^L (1 - W_h) s_h^2 \quad (72.4)$$

که در آن، زیرنویس P نمادی برای طبقه‌بندی پسین* است. در (72.4) جمله اول، واریانس میانگین نمونه طبقه‌بندی تحت تخصیص متناسب است. جمله دوم همیشه نامنفی است و مقدار افزایش واریانس را که می‌توان به جای طبقه‌بندی پیشین از طبقه‌بندی پسین انتظار داشت نشان می‌دهد. توجه کنید که در مخرج کسر جمله n^2 قرار دارد و معنای آن است که این جمله معمولاً خیلی کوچک است. به طور خلاصه، تقریب (71.4) تنها وقتی خوب است که n بزرگ و n_h مثبت باشد. همین طور، مقدار افزایش واریانس، به شرط بزرگ بودن n ، کوچک خواهد بود. پس، طبقه‌بندی پسین تنها وقتی

* حرف P نمادی برای Poststratification، به معنای طبقه‌بندی پسین، است.

نتایج خوبی می‌دهد که n بزرگ و تمام n_h ها نیز نسبتاً بزرگ باشند. یک پیامد کاربردی از این نتیجه آن است که طبقه‌بندی پسین را نمی‌توان در طبقات زیاد انجام داد.

مثال ۸.۴ مدیر شرکتی می‌داند که ۴۰٪ حسابهای قابل وصولش از عمده‌فروشیها و ۶۰٪ بقیه از خرده‌فروشیها به دست می‌آیند. اما، مشخص کردن حسابهای فردی بدون بررسی هر پرونده و مطالعه آن، مشکل است. حسابرسی می‌خواهد که نمونه‌ای به حجم $n = ۱۰۰$ از بین حسابها، برای برآورد متوسط مقدار حسابهای قابل وصول شرکت تهیه کند. نمونه‌ای تصادفی استخراج می‌کند که بعد از بررسی، ۷۰٪ واحدهای نمونه مربوط به عمده‌فروشیها و ۳۰٪ مربوط به خرده‌فروشیهاست. پس از انجام طبقه‌بندی روی این نمونه، نتایج نمونه‌ای زیر حاصل شده‌اند

| خرده‌فروشی | عمده‌فروشی |
|-------------------|-------------------|
| $n_2 = 30$ | $n_1 = 70$ |
| $\bar{y}_2 = 280$ | $\bar{y}_1 = 520$ |
| $s_2 = 90$ | $s_1 = 210$ |

پارامتر μ ، متوسط مقدار حسابهای قابل وصول شرکت را برآورد کنید. چون نسبت مشاهده شده حسابهای عمده‌فروشی یعنی ۷۰٪، از نسبت واقعی این حسابها در جامعه، یعنی ۴۰٪ خیلی فاصله دارد، به نظر می‌رسد که طبقه‌بندی بعد از نمونه‌گیری تصادفی طبقه‌بندی نامناسب است. می‌دانیم

$$\bar{Y}_{st} = \left(\frac{N_1}{N}\right)\bar{y}_1 + \left(\frac{N_2}{N}\right)\bar{y}_2 = (0.4)(520) + (0.6)(280) = 376$$

با استفاده از (۷۲.۴) و با اغماض از ضریب تصحیح جامعه متناهی، داریم

$$\begin{aligned}\hat{V}_P(\bar{Y}_{st}) &= \frac{1}{n} \sum_{h=1}^2 W_h s_h^2 + \frac{1}{n^2} \sum_{h=1}^2 (1 - W_h) s_h^2 \\ &= \frac{1}{100} [0.4(210)^2 + 0.6(90)^2] + \frac{1}{(100)^2} [0.6(210)^2 + 0.4(90)^2] \\ &= 225 + 297 = 227.97\end{aligned}$$

در عبارت $\hat{V}_P(\bar{Y}_{st})$ ، اولین جمله، مقداری است که اگر نمونه را از قبل طبقه‌بندی می‌کردیم (و این نتایج نمونه‌ای حاصل می‌شد) به دست می‌آمد. جمله دوم بهایی است که برای عدم طبقه‌بندی از قبل می‌پردازیم. ▲

۱۶.۴ نمونه‌گیری مضاعف برای طبقه‌بندی

تا اینجا فرض بر این بود که مقادیر $W_h = \frac{N_h}{N}$ ، $h = 1, 2, \dots, L$ ، قبل از شروع نمونه‌گیری، مقادیر ثابت معلومی هستند. ولی همیشه این طور نیست. مثلاً ممکن است بخواهیم جامعه رأی‌دهندگان

را برحسب جنس، سطح درآمد، یا سطح تحصیلات طبقه‌بندی کنیم ولی از روی فهرست ثبت اسامی رأی‌دهندگان، اطلاعات لازم برای طبقه‌بندی این صفات موجود نیست. اندیشه اصلی نمونه‌گیری مضاعف (نمونه‌گیری دو مرحله‌ای) برای طبقه‌بندی، نسبتاً ساده است، ولی برآورد واریانس میانگین در این حالت پیچیده است. فرض کنید اطلاعات مقدماتی نظیر جنس رأی‌دهندگان برای طبقه‌بندی به آسانی به دست آیند اما اطلاعات مفصلتر، مانند نظر سیاسی آنها یا علت مخالفت آنها با موضوع رأی‌گیری به آسانی در اختیار نبوده و نیاز به مصاحبه داشته باشد. بدیهی است انجام مصاحبه با همه رأی‌دهندگان کاری وقتگیر، پرهزینه و اصولاً نشدنی است. در این صورت نمونه‌ای بزرگ برای تشکیل طبقات و نمونه‌ای خیلی کوچکتر برای گردآوری داده‌ها اختیار می‌کنند. مثلاً برای مشخص کردن جنس رأی‌دهندگان نمونه‌ای بزرگ انتخاب می‌کنیم (مرحله اول) و برای تکمیل اطلاعات به تصادف با عده‌ای مصاحبه می‌نماییم (مرحله دوم). فرض کنید نمونه مرحله اول به حجم n' برای تعیین اینکه واحدها در کدام طبقات می‌افتند مورد استفاده قرار گیرد. گیریم

$$w'_h = \frac{n'_h}{n'} \quad h = 1, \dots, L$$

نسبت اولین بخشی از نمونه باشد که در طبقه h ام می‌افتد. واضح است که w'_h برآوردکننده نااریب W_h است (چرا؟). بدیهی است نمونه اول به تصادف انتخاب شده است. در مرحله دوم، از n'_h واحد متعلق به طبقه h ام، به صورتی تصادفی، n_h واحد انتخاب می‌کنیم. اندازه‌های تحت بررسی (اطلاعات مفصلتر) را از این n_h واحد به دست می‌آوریم و میانگین اندازه‌ها را \bar{Y}_h و میزان تغییر آنها را s_h می‌گیریم. این عمل را برای همه L طبقه انجام می‌دهیم. اگر \bar{Y}_N میانگین جامعه‌ای این صفت باشد، داریم

$$\hat{\bar{Y}}_N = \bar{Y}'_{st} \simeq \sum_{h=1}^L w'_h \bar{Y}_h$$

اگر کسرهای نمونه‌گیری مرحله دوم در هر طبقه، یعنی $\frac{n_h}{N_h}$ کوچک باشند و N بزرگ فرض شود، واریانس تقریبی \bar{Y}'_{st} به صورت زیر برآورد می‌شود

$$\hat{V}(\bar{Y}'_{st}) = \frac{n'}{n' - 1} \sum_{h=1}^L \left[\left(w'^2_h - \frac{w'_h}{n'} \right) \frac{s^2_h}{n_h} + \frac{w'_h (\bar{Y}_h - \bar{Y}'_{st})^2}{n'} \right]$$

(برای اثبات به [۳] رجوع شود)

اگر n' آن قدر بزرگ باشد که w'/n'_h قابل اغماض باشد این برآورد واریانس به صورت زیر خلاصه می‌شود

$$\hat{V}(\bar{Y}'_{st}) = \sum_{h=1}^L \left[\frac{w'^2_h s^2_h}{n_h} + \frac{w'_h (\bar{Y}_h - \bar{Y}'_{st})^2}{n'} \right] \quad (۷۳.۴)$$

برآورد واریانس \bar{Y}_{st} را در حالت کلی به خاطر بیاورید. این برآورد واریانس را به صورت زیر می‌توان نوشت

$$\hat{V}(\bar{Y}_{st}) = \frac{1}{N^2} \sum_{h=1}^L N_h^2 \left(\frac{N_h - n_h}{N_h} \right) \left(\frac{s_h^2}{n_h} \right) \quad (74.4)$$

اگر (۷۴.۴) را با (۷۳.۴) مقایسه کنیم می‌بینیم که صرف نظر از تصحیح جامعه متناهی، قسمت اول سمت راست برابری (۷۳.۴) با سمت راست (۷۴.۴) یکی است، که در آن به جای $\frac{N_h}{N}$ یعنی W_h مقدار w'_h گذاشته شده است و قبلاً هم گفتیم که w'_h برآوردکننده ناریب W_h است. جمله دوم (۷۳.۴) مؤلفه اضافی واریانس است که به دلیل عدم اطلاع از وزنه‌های دقیق طبقات به وجود آمده است. مطلب را با مثالی تشریح می‌کنیم.

مثال ۹.۴ می‌خواهند از روی فهرست ثبت نامها و حجم دانشکده‌های کشوری، متوسط تعداد ثبت نام را در سال معینی برآورد کنند. دانشکده‌های خصوصی، ثبت نامهایی کمتر از دانشکده‌های دولتی دارند. لذا دو طبقه خصوصی و دولتی در نظر گرفته شده‌اند. تعداد دانشکده‌های خصوصی و دولتی را نمی‌دانند. ولی تشخیص آن با انتخاب هر واحد نمونه سریعاً مشخص می‌شود. از نمونه‌ای به حجم $n' = 141$ دانشکده نتیجه شده است که

$$n'_1 = 84 \text{ خصوصی} \quad n'_2 = 57 \text{ دولتی}$$

زیرنمونه‌هایی به حجم ۱۱ دانشکده خصوصی و ۱۲ دانشکده دولتی به تصادف از نمونه بالا انتخاب کرده، تعداد ثبت نامهای آنها را در سال مورد نظر تعیین نموده‌اند. داده‌ها در جدول زیر آمده‌اند.

| خصوصی $n_1 = 11$ | | دولتی $n_2 = 12$ | |
|---------------------|---------|---------------------|---------|
| تعداد ثبت نام | دانشکده | تعداد ثبت نام | دانشکده |
| ۱۶۱۸ | ۱۲۲ | ۷۳۳۲ | ۴۵۲ |
| ۱۱۴۰ | ۸۸ | ۲۳۵۶ | ۱۳۱ |
| ۱۰۰۰ | ۶۵ | ۲۱۸۷۹ | ۹۹۶ |
| ۱۲۲۵ | ۵۵ | ۹۳۶ | ۵۰ |
| ۷۹۱ | ۷۹ | ۱۲۹۳ | ۱۰۶ |
| ۱۶۰۰ | ۷۹ | ۵۸۹۴ | ۳۲۶ |
| ۷۴۶ | ۴۰ | ۸۵۰۰ | ۵۰۶ |
| ۱۷۰۱ | ۷۵ | ۶۴۹۱ | ۳۷۱ |
| ۷۰۱ | ۳۲ | ۷۸۱ | ۱۰۷ |
| ۶۹۱۸ | ۴۲۸ | ۷۲۵۵ | ۲۹۸ |
| ۱۰۵۰ | ۱۱۰ | ۲۱۳۶ | ۱۲۸ |
| | | ۵۳۸۰ | ۲۸۰ |

بروز میانگین تعداد ثبت نامها و برآورد واریانس برآوردکننده میانگین را محاسبه کنید.
سایر آنچه در بالا گفتیم

$$\bar{Y}'_{st} = w'_1 \bar{Y}_1 + w'_2 \bar{Y}_2$$

در این مثال که داده‌ها واقعی هستند W'_1 و W'_2 را در کل جامعه نداریم. w'_1 و w'_2 را از روی نمونه
را بدست می‌آوریم

$$w'_1 = \frac{84}{141} \quad . \quad w'_2 = \frac{57}{141}$$

این دو مقدار برآورد نااریب W'_1 و W'_2 هستند. با توجه به داده‌های جدول با محاسبه \bar{Y}_1, \bar{Y}_2 و
 s_1, s_2 به ترتیب داریم

$$\begin{aligned} \bar{Y}'_{st} &= \left(\frac{84}{141}\right)(1681) + \left(\frac{57}{141}\right)(5853) \\ &= (0.60)(1681) + (0.40)(5853) = 3349.8 \end{aligned}$$

$$\hat{V}(\bar{Y}'_{st}) = \frac{1}{n_1}(w'_1 s_1)^2 + \frac{1}{n_2}(w'_2 s_2)^2 + \frac{1}{n'}[w'_1(\bar{Y}_1 - Y'_{st})^2 + w'_2(\bar{Y}_2 - Y'_{st})^2]$$

برای نوشتن این رابطه از (۷۳.۴) استفاده کرده‌ایم. پس

$$\begin{aligned} \hat{V}(\bar{Y}'_{st}) &= \frac{1}{11}[(0.60)(1773)]^2 + \frac{1}{12}[(0.40)(5763)]^2 \\ &\quad + \frac{1}{141}[(0.60)(1681 - 3349.8)^2 + (0.40)(5853 - 3349.8)^2] \\ &= 545708.05 + 29626.52 = 575334.57 \end{aligned}$$

در این محاسبه با توجه به داده‌های جدول از $s_1 = 1773$ و $s_2 = 5763$ استفاده شده است. جمله
دوم در محاسبه برآورد واریانس، یعنی 29626.52 نسبتاً بزرگ به نظر می‌رسد. توجه دارید که این
مقدار، افزایش واریانس ناشی از به‌کار بردن w'_1 و w'_2 به جای W_1 و W_2 است که نامعلوم‌اند. البته با
وجود ظاهر بزرگ این افزایش، با کمی دقت ملاحظه می‌شود که این مقدار فقط ۵٪ کل واریانس است. ▲

تمرینها

۱. میزان محصول سیب درختان باغی برحسب ده کیلوگرم به شرح جدول زیر است. محصولها
برحسب ردیف درختان در جدول آمده‌اند. سه ردیف آخر، درختان جوانتر هستند. جامعه را به دو

طبقه درختان جوان و درختان دیگر تقسیم می‌کنیم. از طبقه اول ۲ درخت با محصولهای ۵ و ۴ و از طبقه دوم ۶ درخت با محصولهای ۶، ۴، ۷، ۶، ۷، ۹ به تصادف انتخاب می‌کنیم. برآورد میانگین محصول درختان باغ را به دست آورید. برآورد واریانس این برآوردکننده را معین کنید. اگر از کل جدول استفاده کنیم واریانس دقیق این میانگین چقدر است؟ در سطح معنادار بودن ۵ درصد، از روی برآورد میانگین محصول درختان باغ، بازه اطمینانی برای میانگین محصول درختان باغ به دست آورید. اگر دو نمونه بالا یک نمونه تصادفی (ساده) از جامعه درختان به حساب آید، برآورد و میانگین محصول درختان باغ و برآورد واریانس این برآوردکننده چقدر است؟

| | | | | | | | |
|----|----|----|----|----|----|----|----|
| ۷ | ۸ | ۵ | ۶ | ۶ | ۱۰ | ۷ | ۶ |
| ۵ | ۴ | ۴ | ۷ | ۶ | ۶ | ۳ | ۸ |
| ۴ | ۷ | ۸ | ۱۰ | ۸ | ۴ | ۶ | ۶ |
| ۶ | ۴ | ۶ | ۴ | ۸ | ۷ | ۹ | ۸ |
| ۶ | ۳ | ۹ | ۹ | ۷ | ۸ | ۱۱ | ۹ |
| ۵ | ۳ | ۴ | ۳ | ۵ | ۲ | ۴ | ۳ |
| ۵ | ۳ | ۳ | ۴ | ۵ | ۴ | ۳ | ۴ |
| ۴ | ۵ | ۴ | ۳ | ۳ | ۴ | ۳ | ۵ |
| ۴۲ | ۳۷ | ۴۳ | ۴۶ | ۴۸ | ۴۵ | ۴۶ | ۴۹ |

۲. جمعیت سه شهر کوچک به ترتیب $N_1 = 40000$ ، $N_2 = 20000$ و $N_3 = 30000$ نفر است. می‌خواهیم برای بررسی مشخصه‌ای در این ۳ شهر نمونه‌ای تصادفی با طبقه‌بندی شامل ۴۰۰ واحد انتخاب کنیم. از روی سرشماریهای گذشته داریم $S_1 = 20$ ، $S_2 = 12$ ، $S_3 = 14$. می‌دانیم هزینه انتخاب هر واحد ۳۶ است. حجم نمونه‌ای که باید از هر شهر انتخاب کنیم در دو حالت زیر به دست آورید

الف) وقتی تخصیص، متناسب با حجم است. ب) وقتی تخصیص ایتیم است.
 ۳. در یک بررسی جمعیت شناختی روستاهای یک استان از روش نمونه‌گیری تصادفی با طبقه‌بندی استفاده می‌شود. هر شهرستان یک طبقه منظور می‌شود. اگر هزینه جمع‌آوری اطلاعات برای هر واحد ۱۰ باشد و هزینه‌های اداری و غیره روی هم ۱۰۰۰۰۰۰ باشند، اندازه ایتیم n حجم نمونه را طوری بیابید که واریانس میانگین نمونه دارای کوچکترین مقدار شود. هزینه کل پیش‌بینی شده برای این بررسی ۸۰۰۰۰۰۰ است. اطلاعات موجود در جدول زیر آمده است. با این اطلاعات حجم نمونه هر طبقه را تعیین کنید.

| طبقه | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ |
|-------------|-----|-----|-----|------|------|-----|-----|-----|-----|
| تعداد روستا | ۱۱۰ | ۸۲ | ۶۶ | ۵۲ | ۲۳ | ۶۸ | ۱۱۰ | ۹۰ | ۱۷۰ |
| جمعیت متوسط | ۴۰۰ | ۸۰۰ | ۹۰۰ | ۱۱۰۰ | ۱۹۰۰ | ۶۰۰ | ۴۵۰ | ۳۸۰ | ۳۰۰ |
| S_h | ۵۰۰ | ۹۰۰ | ۹۲۰ | ۱۱۶۰ | ۱۹۵۰ | ۵۰۰ | ۷۸۰ | ۵۰۰ | ۵۰۰ |

۴. یک مؤسسه برای تعیین متوسط درآمد روزانه روستایان سه روستا قصد دارد نمونه‌گیری

انجام دهد. تعداد روستاییان ۳ روستا به ترتیب ۱۵۰، ۶۰، و ۹۰ نفرند. از روی نمونه‌گیری مقدماتی واریانسهای درآمدهای ۳ روستا به ترتیب $\frac{۱۴۹}{۶}$ ، $\frac{۵۹}{۱۵}$ ، و $\frac{۸۹}{۱۱}$ به دست آمده است. هزینه تهیه یک واحد

نمونه در سه روستا به ترتیب ۴، ۹، ۱۶ است.
الف) مؤسسه ابتدا تصمیم می‌گیرد n فرد، جمعاً از ۳ روستا، انتخاب کند. با اطلاعات بالا، از هر روستا چه سهمی از n را باید برای نمونه‌گیری در نظر گرفت تا واریانس برآوردکننده میانگین مینیمم باشد؟
ب) اگر بودجه‌ای برابر ۴۴۰ که شامل هزینه‌های اداری نیست برای انجام تحقیق تخصیص

یابد، با این بودجه دقیقاً حجم نمونه‌ای که باید از هر روستا انتخاب شود چقدر است؟
ج) قبل از اجرای نمونه‌گیری، ممکن است تصمیم بگیرند که برآورد متوسط درآمد با دقتی خاص محاسبه شود. بدین معنا که واریانس برآوردکننده برابر $\frac{۸۵}{۱۱}$ باشد. در این صورت از هر روستا چه تعداد واحد نمونه باید برگزیند؟

د) اگر از ۳ روستا به ترتیب ۱۶، ۶، ۹ روستایی به تصادف انتخاب شوند و میانگین درآمد روزانه این ۳ نمونه به ترتیب ۱۰۰، ۱۲۰، ۱۱۰ باشد، برآوردهای ناریبی برای متوسط درآمد روزانه جامعه هر روستا بیابید. اگر واریانس این ۳ نمونه به ترتیب ۱۵، ۵، و $\frac{۱۱}{۳}$ باشد برآورد واریانس برآوردکننده متوسط درآمد هر یک از ۳ روستا و برآورد واریانس برآوردکننده متوسط درآمد روزانه جامعه ۳ روستا را تعیین کنید.

۵. برای تعیین برآورد متوسط مدت زمانی که مبتلایان به یک نوع بیماری اعصاب در بیمارستانهای شهری بستری می‌شوند، تعداد بیمارانی را که در یک سال در ۳ بیمارستان موجود در شهر بستری بوده‌اند در نظر می‌گیرند. تعداد بیماران به ترتیب ۶۰۰، ۲۲۰، و ۳۶۰ نفر بوده‌اند. از روی یک بررسی مقدماتی واریانس تعداد روزهای بستری بودن بیماران تقریباً برابر با $\frac{۵۱۱}{۱۱}$ ، $\frac{۲۲۹}{۱۱}$ ، و $\frac{۲۵۹}{۱۱}$ به دست آمده است. هزینه کسب اطلاع درباره یک بیمار در ۳ بیمارستان به ترتیب ۲، ۲، و ۹ است.

الف) اگر بودجه‌ای معادل ۱۱۸۲ که هزینه‌های اداری را شامل نیست برای تحقیق تخصیص دهند حجم نمونه هر بیمارستان چقدر باید باشد تا واریانس برآوردکننده مورد نظر مینیمم شود؟
ب) اگر از ۳ بیمارستان به ترتیب ۱۵۰، ۲۲، و ۵۲ بیمار به عنوان نمونه انتخاب کنند و متوسط مدت بستری بودن آنها به ترتیب ۲۰، ۲۵، و ۲۲ روز باشد برآوردی ناریب برای متوسط مدت بستری بودن در هر ۳ بیمارستان به صورت جداگانه و همچنین برآوردی ناریب برای متوسط مدت بستری بودن جامعه بیماران این نوع بیماری به دست آورید. اگر تغییرات این سه نمونه به ترتیب ۹، ۱۵، و ۳۵ باشد برآورد واریانس برآوردکننده متوسط مدت بستری بودن در هر بیمارستان و برآورد واریانس برآوردکننده متوسط مدت بستری بودن در بیمارستانها را محاسبه کنید.

ج) از سوابق موجود در بیمارستانها دریافته‌اند که از نمونه‌های بالا به ترتیب ۳، ۵، و ۱۰ نفر از نظر ارثی گرفتار این بیماری شده‌اند. برآوردی ناریب برای نسبت افراد بیماری را که در کل جامعه بیماران مذکور به دلیل وراثت بیمار شده‌اند بیابید و تعداد این افراد را برای کل جامعه در یک سال برآورد کنید. برآورد واریانس این برآوردکننده را به دست آورید.

د) اگر بخواهند قبل از نمونه‌گیری قسمت الف، متوسط مدت بستری بودن را با دقتی خاص

بیاوند به قسمی که واریانس برآوردکننده این متوسط برابر r^2 باشد در این صورت حجم نمونه‌ای که از هر بیمارستان باید بگیرند چقدر است؟

۶. در یک نمونه‌گیری تصادفی با طبقه‌بندی، تعداد طبقات ۴ و وزن طبقات به ترتیب عبارت‌اند از $\frac{1}{4}, \frac{1}{6}, \frac{1}{4}$ و $\frac{1}{4}$ و $N = 240$. اگر تعداد واحدهای نمونه در ۴ طبقه به ترتیب $20, 15, 10$ و 15 و تغییرات نمونه‌ای در طبقات به ترتیب $12, 8, 10, 12$ و میانگینهای نمونه‌ای در طبقات به ترتیب $14, 8, 12, 15$ باشند برآورد میانگین جامعه و برآورد واریانس این برآوردکننده را بیابید. در صورتی که توزیع \bar{Y}_{st} را نرمال بگیریم بازه اطمینانی برای \bar{Y}_N میانگین جامعه، با ضریب اطمینان 95% بیابید. اگر این توزیع نرمال نباشد با استفاده از روش ساترتوایت بازه اطمینانی برای \bar{Y}_N با همان ضریب اطمینان بیابید.

۷. از دفاتر ثبت سابقه ۴ بیمارستان تعداد کل افرادی که گروه خونی آنها معین شده است به ترتیب عبارت‌اند از $1200, 1540, 760$ و 1300 . از این ۴ دفتر به تصادف $50, 60, 30$ و 60 نفر را انتخاب می‌کنند. از همان دفاتر دریافته‌اند که $20, 22, 10, 18$ نفر دارای گروه خونی A هستند. الف) برآوردی نااریب برای نسبت افرادی که در کل جامعه مراجعه کنندگان دارای گروه خونی A هستند بیابید و تعداد این افراد را برآورد کنید.

ب) برآورد واریانس برآوردکننده این نسبت را تعیین کنید.

ج) اگر در این مسأله تعداد افراد نمونه معلوم نباشد و قصد داشته باشیم جمعاً 70 نفر انتخاب کنیم و از سوابق قبلی این بیمارستانها بدانیم نسبت افرادی که دارای گروه خونی A هستند به ترتیب $4\%, 34\%, 32\%$ و 35% است. معین کنید که برای نمونه‌گیری مورد نظر از هر بیمارستان، با تخصیص نیم، تقریباً چند نفر باید انتخاب کنیم؟

۸. در یک نمونه‌گیری از دو طبقه مایلیم به جای n_1 و n_2 در انتساب نیمن داشته باشیم $n_1 = n_2$. اگر $V(\bar{Y}_{st})$ معرف واریانس در حالت $n_1 = n_2$ و $V_{opt}(\bar{Y}_{st})$ واریانس مربوط به تخصیص نیمن باشد نشان دهید وقتی N بزرگ باشد

$$\frac{V(\bar{Y}_{st}) - V_{opt}(\bar{Y}_{st})}{V_{opt}(\bar{Y}_{st})} = \left(\frac{r - 1}{r + 1} \right)^2$$

که در آن $r = \frac{n_1}{n_2}$.

۹. فرض کنید دو گروه از مصرف‌کنندگان یک قلم کالا را در نظر گرفته و نظر آنها را درباره کیفیت کالا خواسته‌ایم. اگر ۱ معرف نظر موافق و ۰ معرف نظر ناموافق باشد، داریم

۱ گروه: ۱, ۱, ۰, ۱, ۰, ۱, ۱, ۱

۲ گروه: ۱, ۰, ۰, ۰, ۱, ۰, ۱, ۰

الف) P و P_h ها را حساب کنید.

ب) دو نمونه به حجم $n_h = 3$ انتخاب کنید و p_1 و p_2 را به دست آورید.

ج) با استفاده از نتایج قسمت (ب) مقدار P را برآورد کنید.
 ۱۰. برای جامعه‌ای فرضی از دانشجویان ۴ دانشکده، اطلاعاتی به صورت جدول زیر داریم. تعداد دانشجویان هر دانشکده در ستون دوم جدول، حجم نمونه هر دانشکده در ستون سوم و نسبت دانشجویانی که از هر دانشکده حداقل یک بار در سال به پزشک مراجعه کرده‌اند در ستون آخر ثبت شده‌اند. نسبت دانشجویانی را برآورد کنید که در طول سال گذشته حداقل یک بار به پزشک مراجعه کرده‌اند. همچنین $\hat{V}(p_{st})$ را حساب کنید

| دانشکده | N_h | n_h | p_h |
|---------|-------|-------|-------|
| ۱ | ۲۰۰۰ | ۱۰۰ | ۰٫۲ |
| ۲ | ۱۶۰۰ | ۸۰ | ۰٫۳ |
| ۳ | ۱۲۰۰ | ۶۰ | ۰٫۴ |
| ۴ | ۱۲۰۰ | ۶۰ | ۰٫۳ |

۱۱. به مثال ۱۰٫۴ رجوع کنید و آنچه را که درباره جمعیت ۶۴ شهر در ۱۹۳۰ خواسته بودیم درباره جمعیت ۶۴ شهر در ۱۹۲۰ نیز بیابید.

۱۲. پژوهشگری می‌خواهد برآوردی از متوسط فروش سالیانه، ۵۶ شرکت را به دست آورد. برای این کار نمونه‌ای به حجم ۱۵ شرکت انتخاب می‌کند. داده‌های فراوانی مربوط به این شرکتها به صورت طبقه‌بندی با فاصله‌های ۵۰۰۰۰ دلار موجودند که در جدول زیر ارائه شده‌اند. برای $L = 3$ بهترین تخصیص کرانه‌های این طبقات را تعیین کنید.

| فراوانی | درآمد (برحسب ۱۰۰۰ دلار) | فراوانی | درآمد (برحسب ۱۰۰۰ دلار) |
|---------|----------------------------|---------|----------------------------|
| ۵ | ۳۰۰-۳۵۰ | ۱۱ | ۱۰۰-۱۵۰ |
| ۸ | ۳۵۰-۴۰۰ | ۱۴ | ۱۵۰-۲۰۰ |
| ۳ | ۴۰۰-۴۵۰ | ۹ | ۲۰۰-۲۵۰ |
| ۲ | ۴۵۰-۵۰۰ | ۴ | ۲۵۰-۳۰۰ |

۱۳. جامعه‌ای به ۳ طبقه افراز شده است. حجم طبقات به ترتیب ۲۰۰، ۳۰۰ و ۴۰۰ است. مقادیر S_h طبقه‌ها به ترتیب ۲ و ۴ و ۵ هستند. بودجه کل نمونه‌گیری، صرف نظر از هزینه‌های مشترک اداری و غیره ۳۶۰۰ و هزینه نمونه‌گیری هر واحد در هر طبقه برابر ۵۰ است.

الف) اگر تخصیص نیمین مورد نظر باشد، حجم نمونه در هر طبقه چقدر است؟
 ب) اگر از ۳ طبقه به ترتیب ۸، ۲۴ و ۴۰ واحد به تصادف بدون جایگذاری انتخاب کنیم، به شرط آنکه میانگین این نمونه‌ها به ترتیب ۸، ۱۰ و ۱۲ باشند، برآورد میانگین جامعه چقدر است؟
 برآورد واریانس این برآوردکننده چقدر است؟

ج) با ضریب اطمینان ۹۵ درصد، به شرط آنکه توزیع تقریبی \bar{Y}_{st} نرمال باشد، بازه اطمینانی برای میانگین جامعه بیابید.

۱۴. نمونه‌گیری تصادفی با طبقه‌بندی همیشه برآوردکننده‌ای با واریانس کوچکتر از واریانس برآوردکننده متناظرش در نمونه‌گیری تصادفی ساده تولید نمی‌کند. مسأله زیر این مطلب را نشان می‌دهد. یک توزیع‌کننده مواد غذایی در شهری بزرگ می‌خواهد بداند که آیا تقاضاها آن قدر هستند که فرآورده‌ای جدید را به بخش توزیعش اضافه کند یا نه. برای یاری به اخذ تصمیم، طراحی می‌کند که این فرآورده جدید را به نمونه‌ای تصادفی از فروشگاههایی که مشتری او هستند بدهد و میانگین فروش ماهیانه این فرآورده را برآورد کند. چهار زنجیره از فروشگاهها مشتری او هستند. هر زنجیره را یک طبقه می‌گیرد. در زنجیره ۱، جمعاً ۲۴ فروشگاه، در زنجیره ۲، به تعداد ۳۶ فروشگاه، در زنجیره ۳، به تعداد ۳۰ فروشگاه، و در زنجیره ۴، به تعداد ۳۰ فروشگاه وجود دارند. فروشگاه توان مادی و اداری تهیه داده‌های ماهیانه ۲۰ فروشگاه را دارد. توزیع‌کننده از تخصیص متناسب برای انبام یک نمونه‌گیری با طبقه‌بندی استفاده می‌کند. از چهار طبقه نمونه‌هایی تصادفی اختیار می‌کند. نتایج فروش ماهیانه در فروشگاههای منتخب به شرح جدول زیر است

| طبقه ۱ | طبقه ۲ | طبقه ۳ | طبقه ۴ |
|--------|--------|--------|--------|
| ۹۴ | ۹۱ | ۱۰۸ | ۹۲ |
| ۹۰ | ۹۹ | ۹۶ | ۱۱۰ |
| ۱۰۲ | ۹۳ | ۱۰۰ | ۹۴ |
| ۱۱۰ | ۱۰۵ | ۹۳ | ۹۱ |
| | ۱۱۱ | ۹۳ | ۱۱۳ |
| | ۱۰۱ | | |

میانگین فروش ماهیانه این فرآورده را در کل ۴ زنجیره برآورد کرده و حدود خطای آن را مشخص کنید. اگر نمونه‌ها را یک نمونه تصادفی ساده از جامعه فروشگاهها بگیریم واریانس برآوردکننده حاصل از این نمونه را با واریانس \bar{Y}_{st} قسمت بالا مقایسه کنید. چرا کارایی نمونه‌گیری تصادفی ساده بیش از نمونه‌گیری تصادفی با طبقه‌بندی است؟

۱۵. متغیر تصادفی Y دارای توزیع یکنواخت روی بازه $(a, a + c)$ است. دامنه تغییرات متغیر را به L قسمت برابر تقسیم می‌کنیم تا L طبقه ایجاد شود. از هر طبقه نمونه‌ای تصادفی به حجم $\frac{n}{L}$ می‌گیریم. بار دیگر از کل جامعه نمونه‌ای تصادفی به حجم n تهیه می‌کنیم. اگر واریانس برآوردکننده میانگین جامعه را در نمونه‌گیری تصادفی ساده با V_1 و در نمونه‌گیری با طبقه‌بندی با V_2 نمایش دهیم ثابت کنید که $V_2 = \frac{V_1}{L}$.

تمرینهای چهارگزینه‌ای

۱. متغیر تصادفی Y دارای توزیع یکنواخت روی بازه $(۲, ۵)$ است. دامنه تغییرات متغیر را به ۳ قسمت برابر تقسیم می‌کنیم تا سه طبقه به وجود آیند. از هر طبقه نمونه‌ای تصادفی به حجم ۳۰ می‌گیریم. بار دیگر از کل جامعه نمونه‌ای تصادفی به حجم ۹۰ می‌گیریم. اگر واریانس برآوردکننده

میانگین جامعه را در نمونه‌گیری از کل جامعه با V_1 و در نمونه‌گیری از طبقات با V_2 نشان دهیم، آنگاه
 الف) $V_1 = 3V_2$ (ب) $V_1 = \frac{1}{3}V_2$ (ج) $V_1 = \frac{1}{4}V_2$ (د) $V_1 = V_2$

۲. در یک نمونه‌گیری تصادفی با طبقه‌بندی از تخصیص نیمی استفاده می‌شود. جامعه از ۴ طبقه تشکیل شده است. در این طبقات مقادیر W_i, S_i با اعداد ۲، ۲، ۳، ۲ متناسب‌اند. اگر حجم نمونه طبقه سوم ۶۰ باشد حجم کل نمونه برابر است با

الف) ۱۶۵ (ب) ۱۵۰ (ج) ۱۳۵ (د) ۱۲۰

۳. در یک نمونه‌گیری با طبقه‌بندی تابع هزینه به صورت $C = \sum_{i=1}^2 C_i n_i$ است. n_i حجم نمونه در طبقه i و C_i هزینه تهیه یک واحد نمونه در طبقه i است. اگر W_i و S_i به ترتیب وزن و انحراف معیار تقریبی طبقه i باشد، داریم

| طبقه i | W_i | S_i | C_i |
|----------|-------|-------|-------|
| ۱ | ۰٫۴ | ۱۰ | ۴۰۰ |
| ۲ | ۰٫۶ | ۲۰ | ۹۰۰ |

اگر n حجم کل نمونه در طبقات فرض شود، برای $V(\bar{Y}_{st})$ تثبیت شده، وقتی C مینیمم است که $\frac{n}{11}$ برابر باشد با

الف) $\frac{1}{3}$ (ب) $\frac{2}{5}$ (ج) $\frac{1}{3}$ (د) $\frac{2}{7}$

۴. جامعه‌ای از دو طبقه تشکیل شده است. حجم طبقات به ترتیب ۱۰ و ۱۵ است. می‌خواهیم نمونه‌ای به حجم ۱۰، با روش متناسب با حجم، از دو طبقه اختیار کنیم. تعداد نمونه‌های ممکن برابر است با

الف) $\binom{25}{10}$ (ب) ۲۱۰ (ج) $\binom{15}{6}$ (د) $210 \binom{15}{6}$

۵. وقتی از نمونه‌گیری با طبقه‌بندی برای برآورد میانگین جامعه استفاده می‌کنیم که
 الف) پراکندگی در درون هر طبقه زیاد باشد (ب) پراکندگی در درون هر طبقه کم باشد
 ج) چارچوب نمونه‌گیری را نداشته باشیم (د) حجم نمونه بزرگ باشد

۶. در یک نمونه‌گیری تصادفی با طبقه‌بندی که تنها دارای دو طبقه به وزنهای $\frac{1}{3}$ و $\frac{2}{3}$ است، S^2 ی نمونه‌های دو طبقه به ترتیب ۴ و ۷ است. اگر حجم جامعه ۹۹ و حجم کل نمونه ۹ باشد و نمونه‌گیری با تخصیص متناسب صورت گیرد، برآورد واریانس برآوردکننده میانگین جامعه برابر است با تقریباً
 الف) ۰٫۶ (ب) ۰٫۵۴ (ج) ۰٫۶ (د) ۰٫۵

۷. در یک نمونه‌گیری با طبقه‌بندی و با L طبقه، واریانس برآوردکننده مجموع واحدهای جامعه برابر مجموع N_i, S_i^2 هاست. در این صورت کسر نمونه‌گیری در کل جامعه برابر است با

الف) $\frac{1}{L}$ (ب) $\frac{1}{L}$ (ج) $\frac{1}{L}$ (د) $\frac{1}{L}$

۸. در یک نمونه‌گیری با دو طبقه، کسرهای نمونه‌گیری طبقات کوچک و قابل اغماض‌اند. اگر حجم نمونه طبقه اول ۱۶ و وزن این طبقه $\frac{1}{3}$ باشد و اگر مقادیر S_i^2 ی طبقات برابر با ۴ باشند و بخواهیم واریانس \bar{Y}_{st} برابر با $\frac{1}{18}$ شود حجم طبقه دوم را باید برابر با مقدار زیر اختیار کنیم

الف) ۳۲ (ب) ۶۴ (ج) ۸۰ (د) ۹۶

۱۹۲ نمونه‌گیری تصادفی با طبقه‌بندی

۹. در یک نمونه‌گیری با طبقه‌بندی، اطلاعات زیر در دست است

| طبقه نام | W_i | S_i^2 | C_h |
|----------|---------------|---------|-------|
| ۱ | $\frac{1}{3}$ | ۴ | ۴ |
| ۲ | | ۹ | ۹ |

اگر بودجه نمونه‌گیری صرف‌نظر از مخارج اداری برابر 220 باشد برای اینکه $V(\bar{Y}_{st})$ مینیمم شود حجم کل نمونه را باید برابر مقدار زیر انتخاب کرد

الف) 60 ب) 484 ج) 30 د) 21

۱۰. در یک نمونه‌گیری با طبقه‌بندی، ۴ طبقه داریم. در این طبقه‌ها $N_h S_h / \sqrt{C_h}$ متناسب با اعداد ۲، ۳، ۳ و ۴ هستند. اگر تخصیص اپتیمم بوده و حجم کل نمونه 180 باشد حجم نمونه طبقه ۴ برابر است با

الف) 84 ب) 30 ج) 45 د) 60