

برخی از کاربردهای رایج نمونه‌گیری سیستماتیک

لی-چون زنگ[†]

مرکز آمار نروژ

مترجم: شیرین گلچی

پژوهشکده‌ی آمار

چکیده. نمونه‌گیری سیستماتیک فنی است که در آمارگیری نمونه‌ای به صورت گسترده مورد استفاده قرار می‌گیرد و اجرای آن چه زمانی که واحدها با احتمال برابر انتخاب می‌شوند و چه در حالت انتخاب با احتمالات متناسب با اندازه‌های کمکی آسان است. این فن در صورتی که بتوان به اثرهای مطلوب طبقه‌بندی در فهرست کردن واحدها دست یافت، می‌تواند از کارایی زیادی برخوردار باشد. عیب اصلی آن این است که روش ناریبی برای برآورد واریانس نمونه‌گیری وجود ندارد و این که زمانی که مرتب کردن جامعه بر اساس اطلاعات نادرست صورت گرفته باشد، نمونه‌گیری سیستماتیک ممکن است ضعیف عمل کند. در این مقاله جنبه‌ای از نمونه‌گیری سیستماتیک را که در گذشته به آن توجه زیادی نشده است، مورد بررسی قرار می‌دهیم. نشان داده شده است که در برخی موقعیت‌های متعارف، زمانی که تحت یک مدل مفروض برای جامعه، نمونه‌گیری سیستماتیک به طور متوسط دارای کارایی یکسانی با دیگر روش‌های نمونه‌گیری تصادفی متناظر است، واریانس نمونه‌گیری در نمونه‌گیری سیستماتیک به مراتب بیش‌تر نوسان می‌کند. استفاده از نمونه‌گیری سیستماتیک با مخاطره‌ای در ارتباط است که به طور کلی با بزرگ شدن کسر نمونه‌گیری افزایش می‌یابد. این می‌تواند برای نمونه‌های بزرگ از جامعه‌های کوچک در حالت نمونه‌گیری یک‌مرحله‌ای، یا زیرنمونه‌های بزرگ از زیرجامعه‌های کوچک در حالت نمونه‌گیری چندمرحله‌ای، بسیار مخرب باشد.

[†]Zhang, L.C. (2008). On Some Common Practices of Systematic Sampling. *Journal of Official Statistics*, **24**, 557-569.

واژگان کلیدی: تصمیم آماری؛ مخاطره‌ی بیزی مرتبه دوم؛ طرح استوار؛ آمارگیری پانلی.

دریافت: ۱۳۸۸/۲/۱ پذیرش: ۱۳۸۸/۳/۴

۱- مقدمه

داده‌های گم‌شده مشکل عمومی در تحقیقات مبتنی بر آمارگیری است. نادیده گرفتن هر گونه داده‌های گم‌شده با استفاده از یک تحلیل موردی کامل می‌تواند نتایجی اریب را ایجاد کند. اریبی هنگامی اتفاق می‌افتد که شرکت‌کنندگان دارای داده‌های کامل، به‌طور سیستماتیک با شرکت‌کنندگان دارای داده‌های گم‌شده تفاوت داشته باشند. به‌خصوص مطالعات طولی مستعد چنین اریبی‌هایی هستند زیرا داده‌های گم‌شده به‌مرور زمان در نتیجه بی‌پاسخی و انصراف شرکت‌کنندگان انباشته می‌شود. یکی از روش‌های جبران داده‌های گم‌شده جانپی است. طی بیست سال گذشته انبوه نوشتگان درباره‌ی نظریه و روش‌شناسی جانپی رشد قابل ملاحظه‌ای داشته است و نرم‌افزارهایی نیز در راستای آن تکامل یافته است. ولی در زمینه‌ی جانپی داده‌های طولی کار نسبتاً کمی صورت گرفته است.

برای جانپی چند رهیافت نظری وجود دارد. راگوناتان [۱۵] این قبیل رهیافت‌ها را بازنگری و سه رده را شناسایی کرده است: معادلات برآورد موزون، جانپی چندگانه و فرمول‌های مبتنی بر درست‌نمایی. ابراهیم و دیگران [۸] رهیافت تماماً بیزی را به‌عنوان رده‌ی چهارم در نظر گرفتند.

معادلات برآوردگر موزون (WEE)، آمارهای ثبتي دارای داده‌های کامل را وزن دار می‌کنند تا موارد مشابه دارای داده‌های گم‌شده را جبران کنند. در این اواخر نوشتگان بر بهبود برآورد واریانس ([۱۹] و [۲۰]) تمرکز داشته‌اند زیرا WEE، هنگامی که تعدیل نشده باشد، واریانس واقعی داده‌ها را کم‌برآورد می‌کند. اجرای WEE در حال حاضر به جای شیوه‌های استاندارد در بسته‌های نرم‌افزار آماری اصلی، بر الگوریتم‌های مدل-ویژه و کاربر-تعریف‌شده تکیه دارد. جانپی چندگانه (MI) از شبیه‌سازی بیزی برای پرکردن داده‌های گم‌شده استفاده می‌کند که از تلفیق نتایج به دست آمده از چندین مجموعه‌ی جانپی‌شده‌ی مکرر به دست می‌آیند. برای مشاهده‌ی پوشش فراگیر جانپی چندگانه به رابین [۲۲] مراجعه کنید. مدل‌های تماماً بیزی (FB) روش‌های MI را با شبیه‌سازی توأم توزیع متغیرهای دارای داده‌های گم‌شده و نیز پارامترهای نامعلوم در یک معادله‌ی رگرسیون گسترش می‌دهند. در FB مدل‌های تحلیل و جانپی کاملاً و به‌طور هم‌زمان

مشخص هستند. همچنین فنون ماکسیمم درست‌نمایی نیز بر مدل‌های کاملاً مشخص تکیه دارند اما از این جهت که برآوردهای پارامترها با استفاده از تقریب‌های مبتنی بر درست‌نمایی به جای شبیه‌سازی بیزی ایجاد شده‌اند، با FB فرق دارند.

رهیافت‌های ماکسیمم درست‌نمایی برای جانهی اغلب در بسته‌های نرم‌افزاری اصلی، غیر قابل کنترل‌اند. اجرای روش‌ها بر فرض‌های محکم در مورد الگوهای گم‌شدگی که مکرراً در بررسی داده‌های پیچیده زیر پا گذاشته می‌شوند، تکیه دارد. در حالی که شیوه‌های MI در تعدادی بسته‌های نرم‌افزاری مانند SAS [۲۳]، Stata [۲۷]، S-Plus [۹] و R [۱۸] وجود دارند، این روش‌ها عموماً بر این فرض تکیه دارند که داده‌ها، نرمال چندمتغیره‌اند یا می‌توانند با توزیع نرمال چندمتغیره تقریب زده شوند [۲۴]. بیش‌تر کارهای اخیر درباره‌ی معادلات رگرسیونی زنجیره‌ای، منجر به افزوده شدن بسته‌هایی شده است که داده‌های رسته‌ای را نیز شامل می‌شوند: MICE در S-PLUS [۲۹]، Ice در Stata [۲۱] و IVEware برای SAS [۱۶]. ولی تألیف‌کنندگان هنوز مشکلاتی در اضافه کردن داده‌های طولی در روش‌های جانهی با استفاده از این برنامه‌ها دارند. فنون FB در جانهی طولی بسیار مناسب هستند، زیرا می‌توانند ساختار سلسله‌مراتبی را در فرایند مدل‌سازی ادغام کنند و مانند رگرسیون زنجیره‌ای، این قابلیت را دارند که به‌طور سیستماتیک به داده‌های رسته‌ای نیز پردازند. بسته‌های نرم‌افزار WinBUGS [۲۶] و MLwiN [۱۷] هر دو از چارچوب FB استفاده می‌کنند. کاولس [۵] و وودورت [۳۰] هر دو شرح مختصر و مفیدی برای WinBUGS تهیه کرده‌اند، درحالی‌که کارپنتر و کنوارد [۲] و کانگدون [۳] مثال‌های مقدماتی جانهی FB با داده‌های گم‌شده را مطرح کرده‌اند. پتیت [۱۳] و کیو و همکاران [۱۴] تحلیل‌های دقیقی درزمینه‌ی داده‌های رسته‌ای گم‌شده ارائه کرده‌اند.

هدف این مقاله نشان دادن توانایی WinBUGS برای جبران داده‌های طولی گم‌شده با تمرکز خاص بر داده‌های گم‌شده‌ی متغیر کمکی است. ما این کار را با نگاه بر تحلیل طولی میزان وقوع دیابت در زنان استرالیایی انجام می‌دهیم. در بخش ۲ مثالی تشویق‌کننده از مطالعه‌ی طولی استرالیا درباره‌ی سلامت زنان را معرفی می‌کنیم. در بخش ۳ یک مدل کاملاً بیزی برای میزان وقوع دیابت بدون و با داده‌های گم‌شده‌ی متغیر کمکی را تعیین می‌کنیم. در بخش ۴ اجرای آن را در WinBUGS ارائه کرده و

نتایج را در بخش ۵ توصیف می‌کنیم. در بخش ۶ یک ماکروی کلی SAS (که WinBUGS نامیده می‌شود) برای تحلیل مدل‌های طولی با داده‌های گم‌شده‌ی متغیر کمکی ارائه می‌دهیم. با یک بحث و تعدادی توصیه در بخش ۷ نتیجه‌گیری می‌کنیم.

۲- جامعه‌های همگن

ابتدا نمونه‌گیری سیستماتیک با احتمالات برابر را از یک ترتیب ثابت جامعه که ممکن است با متغیر مورد نظر ناهمبسته تلقی شود، در نظر بگیرید. فرض کنید اندازه‌ی نمونه n و فاصله‌ی نمونه‌گیری k باشد. برای سادگی فرض می‌کنیم k به صورت طبیعی یک مقدار صحیح است که در $N = nk$ صدق می‌کند. m امین نمونه‌ی سیستماتیک را با s_m نمایش می‌دهیم، یعنی $s_m = \{m, m+k, m+2k, \dots, m+(n-1)k\}$. فرض کنید \bar{y}_m میانگین نمونه‌ی متناظر باشد که برآوردگر نارایی از میانگین جامعه است که با $\bar{Y} = \frac{1}{N} \sum_{i \in U} y_i$ نشان داده می‌شود. واریانس نمونه‌گیری \bar{y}_m به صورت $V_{sys} = k^{-1} \sum_{m=1}^k (\bar{y}_m - \bar{Y})^2$ به دست می‌آید که ممکن است از واریانس میانگین نمونه‌ی تصادفی ساده که با $V_{srs} = (n^{-1} - N^{-1}) \sigma^2$ نشان داده می‌شود و در آن $\sigma^2 = \frac{1}{N-1} \sum_{i \in U} (y_i - \bar{Y})^2$ بیش‌تر باشد یا نباشد.

همان‌طور که قبلاً اشاره شد دو نتیجه وجود دارد که نشان می‌دهد SRS و SYS دارای مخاطره‌های بیزی یکسان هستند، یعنی $E_{\theta}(V_{srs}) = E_{\theta}(V_{sys})$. ما مخاطره‌ی بیزی مرتبه‌ی دوم آن‌ها را تحت مدل همگن

$$(۳) \quad E(y_i) = \mu, \quad E((y_i - \mu)^r) = \mu_r < \infty$$

برای $i \in U$ ، که در آن برای سادگی، E را جایگزین E_{θ} کرده‌ایم و به ازای هر $i \neq j$ ، y_i از y_j مستقل است، بررسی خواهیم کرد. داریم $E(V_{srs}) = E(V_{sys}) = (\frac{1}{n} - \frac{1}{N}) \mu_r = (1 - \frac{1}{k}) \mu_{r,n}$ که در آن $\mu_{r,n}$ نمایانگر گشتاور مرکزی r ام \bar{y}_m است. این، حالت خاصی از قضیه‌ی کلی‌تر ۸-۵ کاکران [۲] است که در آن واریانس مدل y_i می‌تواند روی واحدها تغییر کند.

علاوه بر این، بنا بر نتیجه‌ی واریانس واریانس تجربی (به عنوان مثال [۱۲])، به دست می‌آوریم

$$V(V_{srs}) = \left(\frac{1}{n} - \frac{1}{N} \right)^r V(\sigma^r) = \left(1 - \frac{1}{k} \right)^r \mu_{r,n}^r \left(\frac{2}{N-1} + \frac{\gamma_r}{N} \right)$$

که در آن $\gamma_r = \frac{\mu_r}{\mu_r^r - r}$ ضریب کشیدگی y_i است. به همین ترتیب داریم

$$V(V_{sys}) = \left(\frac{k-1}{k} \right)^r V \left(\frac{1}{k-1} \sum_{m=1}^k (\bar{y}_m - \bar{Y})^r \right) = \left(1 - \frac{1}{k} \right)^r \mu_{r,n}^r \left(\frac{2}{k-1} + \frac{\gamma_{r,n}}{k} \right)$$

که در آن $\gamma_{r,n} = \frac{\gamma_r}{n}$ ضریب کشیدگی \bar{y}_m است. نتیجه می‌شود که ضریب تغییرات (CV) V_{sys} و V_{srs} به ترتیب عبارت‌اند از

$$CV(V_{srs}) = \left(1 - \frac{1}{k} \right) \left(\frac{2}{N-1} + \frac{\gamma_r}{N} \right)^{\frac{1}{r}} \quad (۴)$$

و

$$CV(V_{sys}) = \left(1 - \frac{1}{k} \right) \left(\frac{2}{k-1} + \frac{\gamma_r}{k} \right)^{\frac{1}{r}}$$

در صورتی که جامعه به اندازه‌ی کافی بزرگ باشد داریم $CV(V_{sys}) \approx \left(1 - \frac{1}{k} \right) \sqrt{\frac{r}{k-1}}$ که همراه با کسر نمونه‌گیری افزایش می‌یابد. برای مثال، کسر نمونه‌گیری کل در آمارگیری نیروی کار نیروژ در حدود $\frac{1}{14}$ است، به نحوی که $CV(V_{sys}) \approx 12\%$. در مقایسه، $CV(V_{srs}) = O\left(\frac{1}{\sqrt{N}}\right)$ و مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری تصادفی ساده ناچیز است. انتخاب نمونه‌های سیستماتیک از یک فهرست ظاهراً تصادفی، اما در حقیقت ثابت جامعه یک اقدام مبتنی بر شانس و بدون انتظار حصول بهبودی در مقایسه با نمونه‌گیری تصادفی ساده است. به سادگی، کنترل کمتری روی واریانس واقعی نمونه‌گیری وجود دارد که ممکن است به شکل قابل ملاحظه‌ای از امید ریاضی خود انحراف داشته باشد. این امر برای نمونه‌گیری سیستماتیک طبقه‌بندی‌شده در مقایسه با نمونه‌گیری تصادفی ساده‌ی طبقه‌بندی‌شده نیز صادق است. در نمونه‌گیری دومرحله‌ای که در آن نمونه‌گیری سیستماتیک برای زیرنمونه‌گیری واحدها داخل یک واحد نمونه‌گیری اولیه (PSU) مورد استفاده قرار می‌گیرد، آنچه روی

مخاطره‌ی بیزی مرتبه‌ی دوم تأثیرگذار است کسرهای نمونه‌گیری داخل PSUها است. بنابراین نمونه‌گیری سیستماتیک در حالت نمونه‌گیری چندمرحله‌ای می‌تواند مخاطره‌ی بیزی مرتبه‌ی دوم زیادی داشته باشد، حتی اگر کسر نمونه‌گیری کلی کوچک باشد.

۳- جامعه‌های رگرسیونی نسبتی

حال وضعیت نمونه‌گیری سیستماتیک πps را در نظر بگیرید. در این حالت قاعده‌ی «هر k امین» برای مجموع تجمعی یک متغیر کمکی به کار می‌رود که با x_i برای $i \in U$ نشان داده می‌شود. هر فهرست ثابت U می‌تواند مورد استفاده قرار بگیرد. برای سادگی فرض می‌کنیم x_i ها مقادیر صحیح‌اند. فرض کنید $X = \sum_{i \in U} x_i$ ، آن‌گاه طول فاصله توسط $k = \frac{X}{n}$ داده می‌شود که فرض می‌کنیم k نیز به صورت طبیعی یک عدد صحیح است. اگر از زاویه‌ی دیگری به قضیه نگاه کنیم، نمونه‌گیری سیستماتیک با احتمالات برابر همان نمونه‌گیری سیستماتیک πps است که در آن $x_i \equiv 1$. واحد i ممکن است در x_i نمونه‌ی سیستماتیک مختلف ظاهر شود. فرض می‌کنیم احتمال شمول به گونه‌ای است که به ازای همه‌ی $i \in U$ ، $\pi_i = \frac{n x_i}{X} < 1$. بر اساس هر نمونه‌ی سیستماتیک πps که با s_m برای $m = 1, \dots, k$ نشان داده می‌شود، برآوردگر Y عبارت است از

$$\hat{Y}_m = \sum_{i \in s_m} \frac{y_i}{\pi_i} = \frac{X}{n} \sum_{i \in s_m} b_i = X \bar{b}_m$$

که $b_i = \frac{y_i}{x_i}$ و $\bar{b}_m = \sum_{i \in s_m} \frac{b_i}{n}$

$$\text{داریم } E_{\text{sys}}(\bar{b}_m) = \frac{Y}{X} \text{ و } E_{\text{sys}}(\bar{y}_m) = Y$$

$$V_{\text{sys}}(\hat{Y}_m) = X^2 V_{\text{sys}}(\bar{b}_m) = X^2 \left\{ \frac{1}{k} \sum_{m=1}^k \left(\bar{b}_m - \frac{Y}{X} \right)^2 \right\}$$

حال \bar{b}_m بهترین برآوردگر خطی نارایب (BLUE) برای β تحت مدل زیر است

$$(5) \quad y_i = x_i \beta + x_i \varepsilon_i$$

که در آن $E(\varepsilon_i) = 0$ و $V(\varepsilon_i) = \mu_r$ و برای $i \neq j \in U$ ، $\text{COV}(\varepsilon_i, \varepsilon_j) = 0$ ، یعنی یک مدل رگرسیونی نسبتی با واریانس باقیمانده‌ی متناسب با x_i^r داریم

$$V_{\text{sys}}(\bar{b}_m) = \frac{1}{k} \sum_{m=1}^k (\bar{b}_m - \beta)^r - \left(\frac{Y}{X} - \beta \right)^r \doteq \frac{1}{k} \sum_{m=1}^k (\bar{b}_m - \beta)^r = \frac{1}{k} \sum_{m=1}^k \bar{\varepsilon}_m^r \stackrel{\text{def}}{=} Z$$

که در آن $\bar{\varepsilon}_m = \sum_{i \in s_m} \frac{\varepsilon_i}{n}$ ، به شرط آن که جامعه به اندازه‌ی کافی بزرگ باشد و لذا جمله‌ی $\left(\frac{Y}{X} - \beta \right)^r$ در مقایسه با Z در مرتبه‌ی پایین‌تری باشد. داریم $E(Z) = \mu_{r,n}$ به علاوه

$$Z^r = \frac{1}{k^r} \left(\sum_{m=1}^k \bar{\varepsilon}_m^r + \sum_{m,p:p \neq m} \bar{\varepsilon}_m^r \bar{\varepsilon}_p^r \right)$$

که در آن $\bar{\varepsilon}_m$ و $\bar{\varepsilon}_p$ لزوماً از هم مستقل نیستند زیرا ممکن است واحدهایی در هر دوی s_p و s_m ظاهر شوند. هر چند داریم

$$E(\bar{\varepsilon}_m^r \bar{\varepsilon}_p^r) = n^{-r} \sum_{\substack{(i_1, i_r) \in s_m \\ (j_1, j_r) \in s_p}} E(\varepsilon_{i_1} \varepsilon_{i_r} \varepsilon_{j_1} \varepsilon_{j_r})$$

که در آن $E(\varepsilon_{i_1} \varepsilon_{i_r} \varepsilon_{j_1} \varepsilon_{j_r})$ صفر نیست و در واقع مثبت است تنها اگر به صورت $E(\varepsilon_i^r \varepsilon_j^r)$ یا $E(\varepsilon_i^r \varepsilon_j^r)$ باشد. فرض کنید s_{mp} نمایانگر فصل مشترک s_m و s_p باشد. داریم

$$\sum_{\substack{(i_1, i_r) \in s_m \\ (j_1, j_r) \in s_p}} E(\varepsilon_{i_1} \varepsilon_{i_r} \varepsilon_{j_1} \varepsilon_{j_r}) = E \left(r \sum_{i \neq j \in s_{mp}} \varepsilon_i^r \varepsilon_j^r + \sum_{\substack{i \in s_m \\ j \in s_p}} \varepsilon_i^r \varepsilon_j^r \right)$$

جمله‌ی اول با اختیار کردن $(i_1, i_r) = (j_1, j_r)$ و $(i_1, i_r) = (j_1, j_r)$ ، $i_1 \neq i_r$ ، جمله‌ی دوم با اختیار کردن $i_1 = i_r$ و $j_1 = j_r$ ، به دست می‌آید. فرض کنید s_m^p نمایانگر واحدهایی از s_m باشد که s_p شامل آن‌ها نمی‌شود و s_p^m نمایانگر واحدهایی از s_p باشد که s_m آن‌ها را شامل نمی‌شود. داریم

$$\sum_{i \in s_m, j \in s_p} \varepsilon_i^r \varepsilon_j^r = \sum_{i, j \in s_{mp}} \varepsilon_i^r \varepsilon_j^r + \sum_{i \in s_{mp}, j \in s_p^m} \varepsilon_i^r \varepsilon_j^r + \sum_{i \in s_{mp}, j \in s_m^p} \varepsilon_i^r \varepsilon_j^r + \sum_{i \in s_p^m, j \in s_m^p} \varepsilon_i^r \varepsilon_j^r$$

$$\sum_{i,j \in s_{mp}} \varepsilon_i^r \varepsilon_j^r = \sum_{i \in s_{mp}} \varepsilon_i^r + \sum_{i \neq j \in s_{mp}} \varepsilon_i^r \varepsilon_j^r$$

لذا به دست می‌آوریم

$$\begin{aligned} n^r E(\bar{\varepsilon}_m^r \bar{\varepsilon}_p^r) &= r \sum_{i \neq j \in s_{mp}} \mu_i^r + \left(\sum_{i \in s_{mp}} \mu_i^r - \sum_{i \in s_{mp}} \mu_i^r \right) \\ &+ \left(\sum_{i \in s_{mp}} \mu_i^r + \sum_{i \neq j \in s_{mp}} \mu_i^r + \sum_{\substack{i \in s_{mp} \\ j \in s_p^m}} \mu_i^r + \sum_{\substack{i \in s_{mp} \\ j \in s_m^p}} \mu_i^r + \sum_{\substack{i \in s_p^m \\ j \in s_m^p}} \mu_i^r \right) \\ &= r c_{mp} (c_{mp} - 1) \mu_i^r + c_{mp} (\gamma_r + r) \mu_i^r + n^r \mu_i^r = \mu_i^r (n^r + c_{mp} \gamma_r + r c_{mp}^r) \end{aligned}$$

که در آن c_{mp} تعداد واحدهای مشترک s_p و s_m است و دو جمله‌ی دربردارنده‌ی c_{mp} تنها در صورتی که s_{mp} تهی نباشد وجود خواهند داشت. چهارمین گشتاور مرکزی $\bar{\varepsilon}_m$ را با $\mu_{\epsilon,n}$ نمایش می‌دهیم. حال داریم

$$V(Z) = \frac{1}{k} (\mu_{\epsilon,n} - \mu_{r,n}^r) + \frac{\mu_{r,n}^r}{k^r n^r} \left(\gamma_r \sum_{m \neq p} c_{mp} + r \sum_{m \neq p} c_{mp}^r \right) = \mu_{r,n}^r \lambda_n(\gamma_r)$$

که در آن

$$\lambda_n(\gamma_r) = r \left(\frac{1}{k} + \frac{\sum_{m \neq p} c_{mp}^r}{X^r} \right) + \gamma_r \frac{\sum_{i=1}^N x_i^r}{X^r}$$

زیرا $\mu_{\epsilon,n} = \mu_{r,n}^r (\gamma_r + 3)$ ، $X = nk$ و $\sum_{m \neq p} c_{mp} = \sum_{i=1}^N x_i (x_i - 1)$ جمله‌ی $\sum_{m \neq p} c_{mp}^r$ به طور کلی حل‌نشده‌ی به نظر می‌رسد، اما با ترتیب معلوم واحدها می‌توان آن را

$$CV(V_{sys}(\hat{Y}_m)) \doteq CV(Z) = \sqrt{\lambda_n(\gamma_r)}$$

محاسبه کرد. نتیجه می‌شود که

در ضمن انواع گوناگونی از دیگر روش‌های نمونه‌گیری تصادفی πps وجود دارد. به

راحتی می‌توان نشان داد که در حالت نمونه‌گیری پواسون (PS)، ضریب تغییرات $V_{ps}(\hat{Y})$ از مرتبه‌ی $O\left(\frac{1}{\sqrt{N}}\right)$ تحت مدل (۵) است. نتایج به طور کلی برای هر طرح نمونه‌گیری با اندازه‌ی ثابت πps صادق است به شرط آن که واریانس نمونه‌گیری آن از طریق یک جمله‌ی تصحیح جامعه‌ی متناهی با واریانس نمونه‌گیری PS مرتبط باشد.

۴- یک مثال عددی

برای توضیح نتایج (۴) و (۶) از لحاظ عددی، فرض کنید نمونه‌ای را شامل ۱۰ واحد از جامعه‌ای به اندازه‌ی ۱۰۰ که با $U = \{1, 2, \dots, 100\}$ نشان داده می‌شود، در نظر گرفته‌ایم. متغیرهای کمکی به سادگی به صورت $x_i = i$ داده شده‌اند. متغیرهای آمارگیری y_i باید تحت مدل زیر شبیه‌سازی شوند

$$(7) \quad y_i = x_i + x_i^a \varepsilon_i$$

که در آن $\varepsilon_i \sim N^{iid}(0, \sigma^2)$ و $0 \leq a \leq 1$.

لذا واریانس شرطی y_i به شرط x_i برابر است با $x_i^2 \sigma^2$. در حالت $a=0$ ، $y_i - x_i$ از مدل همگن (۳) پیروی می‌کند. در حالت $a=1$ مدل (۵) را با $\beta=1$ داریم که استفاده از آن می‌تواند انگیزه‌ی نمونه‌گیری πps را فراهم کند.

ابتدا نمونه‌گیری πps را در نظر بگیرید. فرض کنید $a=0, 0/25, 0/5, 1$ و $\sigma=0/01, 0/1$ توجه کنید که در این‌جا تا زمانی که مقادیر منفی y بتوانند با احتمال‌های غیر قابل اغماض تولید شوند، که در این صورت پایه و اساس نمونه‌گیری πps زیر سؤال می‌رود، σ نمی‌تواند خیلی بزرگ شود. برای هر زوج (a, σ) یک جامعه‌ی $\theta = (y_1, \dots, y_{100})^T$ تولید می‌کنیم که برای آن سه واریانس نمونه‌گیری محاسبه می‌شود. اولین واریانس برای نمونه‌گیری سیستماتیک πps به دست می‌آید. دومین واریانس برای SPS است که یک روش نمونه‌گیری πps تصادفی تقریبی است و این واریانس با مونت کارلوی ساده محاسبه می‌شود. نهایتاً واریانس مجانبی نظری نمونه‌گیری را برای نمونه‌گیری سیستماتیک πps با جایگشت تصادفی θ پیش از انتخاب یک نمونه‌ی سیستماتیک محاسبه می‌کنیم که عبارت است از [۴]

$$V_{asy} = \sum_{i=1}^{\dots} \pi_i \left(1 - \frac{n-1}{n} \pi_i \right) \left(\frac{y_i}{\pi_i} - \frac{Y}{n} \right)^2$$

و می‌تواند برای محک زدن کارایی دو واریانس دیگر مورد استفاده قرار بگیرد.

جدول ۱- نتایج شبیه‌سازی برای نمونه‌گیری πps بر حسب درصد

CV واریانس نمونه‌گیری			کارایی نسبی		طرح
نظری	SPS	سیستماتیک	SPS	سیستماتیک	
۱۶	۱۶	۴۱	۱۰۰	۹۸	$\sigma = ۰/۰۱$ $a = ۱$
۱۶	۱۶	۴۱	۱۰۰	۱۰۱	$\sigma = ۰/۱$
۱۵	۱۵	۳۴	۱۰۰	۱۰۱	$\sigma = ۰/۰۱$ $a = ۰/۵$
۱۵	۱۵	۳۴	۹۹	۹۹	$\sigma = ۰/۱$
۱۸	۱۸	۲۸	۹۹	۱۰۰	$\sigma = ۰/۰۱$ $a = ۰/۲۵$
۱۸	۱۸	۲۹	۹۹	۹۹	$\sigma = ۰/۱$
۳۵	۳۶	۳۷	۹۹	۱۰۰	$\sigma = ۰/۰۱$ $a = ۰$
۳۷	۳۸	۴۰	۹۹	۱۰۱	$\sigma = ۰/۱$

کارایی نسبی: نسبت بین متوسط واریانس نمونه‌گیری πps سیستماتیک (یا نمونه‌گیری SPS) و متوسط نظری واریانس V_{asy} .

CV واریانس نمونه‌گیری: ضریب تغییرات واریانس‌های نمونه‌گیری.

شبیه‌سازی‌ها برای $\theta = ۱۰۰۰$ تولید شده به صورت مستقل تکرار می‌شوند. نتایج در جدول ۱ ارائه شده‌اند که در آن کارایی نسبی (RE) برابر است با نسبت $\frac{E(V_{\delta})}{E(V_{asy})}$. موارد زیر را ملاحظه می‌کنیم:

- ۱- مشاهده می‌شود که نمونه‌گیری سیستماتیک πps و SPS هر دو کارایی نسبی حول ۱۰۰ درصد به دست می‌دهند، به گونه‌ای که این دو روش نمونه‌گیری از لحاظ اصل مخاطره‌ی بیزی معادل‌اند.
- ۲- تحت مدل (۵) یعنی $a = ۱$ ، CV نمونه‌گیری سیستماتیک πps می‌تواند از (۶) مشتق شود. از نرمال بودن ϵ_i ‌ها داریم $\gamma_r = ۰$ و از U داده شده داریم

..... گزیده‌مطالب آماری، سال ۱۹، شماره ۲، پاییز و زمستان ۱۳۸۷، صص ۲۴۳-۲۲۵

ندارد. CV نظری توسط شبیه‌سازی‌ها تأیید می‌شود. مشاهده می‌شود که نمونه‌گیری سیستماتیک πps دارای مخاطره‌ی بیزی مرتبه‌ی دوم به مراتب بیش‌تری نسبت به نمونه‌گیری تصادفی πps است.

۳- برای a داده شده مخاطره‌ی بیزی مرتبه‌ی دوم روی σ تغییرات زیادی ندارد. برای $0 < a < 1$ مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری تصادفی πps تقریباً ثابت است و به مقدار قابل توجهی از مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری سیستماتیک πps کمتر است. مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری تصادفی πps با نزدیک شدن a به صفر، به سرعت افزایش می‌یابد، اما کوچک‌تر از مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری سیستماتیک πps باقی می‌ماند. به طور خلاصه نمونه‌گیری تصادفی πps با توجه به اصل تصمیم بیز استوار تحت مدل (۵) ترجیح داده می‌شود و این انتخاب نسبت به انحراف از فرض $a = 0$ ، که صورتی از تشخیص نادرست $\pi(\theta)$ است، استوار است.

حال نمونه‌گیری سیستماتیک با احتمالات برابر را در نظر بگیرید. یک نتیجه‌ی کلی وجود دارد که اظهار می‌دارد زمانی که واریانس داخل نمونه‌ای از واریانس جامعه بیش‌تر است، نمونه‌گیری سیستماتیک از SRS کارا تر است و این به خاطر تجزیه‌ی واریانس جامعه به صورت زیر است

$$\sum_{i \in U} (y_i - \bar{Y})^2 = \sum_{m=1}^k \sum_{i \in s_m} (y_i - \bar{y}_m)^2 + \sum_{m=1}^k n (\bar{y}_m - \bar{Y})^2.$$

یعنی تجزیه‌ی واریانس جامعه به تغییرات داخل k نمونه‌ی سیستماتیک و تغییرات بین این نمونه‌ها. از آن‌جا که V_{sys} با مؤلفه‌ی دوم متناسب است به ازای یک θ داده شده زمانی مینیمم می‌شود که مؤلفه‌ی اول ماکسیمم شود. بر اساس ترتیب متناظر واحدها، نمونه‌گیری سیستماتیک می‌تواند به صورت بالقوه به حصول کارایی نسبت به نمونه‌گیری تصادفی ساده منجر شود. برای مثال حالت غایی را تحت مدل (۷) با $\sigma = 0$ در نظر بگیرید یعنی $y_i = x_i$. ترتیب بهینه برای یک نمونه‌ی سیستماتیک شامل ۱۰ واحد، ایجاد

تناوب بين ترتيب كاهشي و افزايشي يكي از هر ۱۰ واحد از جامعه است ([۱۰] مثال ۳-۲-۴)، كه توسط $U_{opt} = (1, \dots, 10, 20, \dots, 11, 21, \dots, 30, 40, \dots, 31, \dots, 100, \dots, 91)$ نشان داده مي شود و در اين حالت $V_{sys}(\bar{y}_s) = 0$.

البته در عمل هرگز y_i ها دقيقاً معلوم نيستند. ولي، حال كه $\bar{x}_m = \frac{X}{N}$ ثابتي از نمونه گيري است، ترتيب U_{opt} تحت مدل (۷) با $a = 0$ بهينه باقي مي ماند. بروردگر مبتني بر يك نمونه ي سيستماتيک انتخاب شده از U_{opt} با احتمالات برابر عبارت است از

$$\hat{Y}_m = X + N\bar{\varepsilon}_m$$

كه در آن $X = \sum_{i=1}^N x_i$ و $\bar{\varepsilon}_m = \frac{1}{n} \sum_{i \in S_m} (y_i - x_i)$. اين بروردگر همان بروردگر تفاضل ([۱۰]، فصل ۳-۶) بر اساس يك نمونه ي تصادفي ساده است. به عبارت ديگر كارايي نمونه گيري سيستماتيک بر اساس U_{opt} مي تواند با استفاده ي توأم نمونه گيري تصادفي ساده و بروردگر تفاضل نيز حاصل شود. البته مخاطره ي بيزي مرتبه ي دوم براي استراتژي دوم تنها از مرتبه ي $O\left(\frac{1}{\sqrt{N}}\right)$ است. اين وضعيت در جدول (۲) نشان داده شده است كه در آن RE برابر $\frac{E(V_{sys})}{E(V_{srs})}$ است و CV_δ ، CV واريانس نمونه اي واقعي توليد شده توسط $\delta = SYS, SRS$ است. موارد زير را ملاحظه مي كنيم:

۱- همان طور كه انتظار مي رفت استفاده ي توأم از نمونه گيري تصادفي ساده و بروردگر تفاضل به اندازه ي نمونه گيري سيستماتيک بهينه تحت مدل (۷) با $a = 0$ كارآمد است. كارايي نمونه گيري سيستماتيک با انحراف از فرض $a = 0$ ، يعني حرکت a از ۰ به سمت ۱، به آهستگي كم مي شود.

۲- مخاطره ي بيزي مرتبه دوم نمونه گيري سيستماتيک به مراتب بيش تر از اين مخاطره براي نمونه گيري تصادفي ساده تحت فرض $a = 0$ است. براي هر دو روش تغييرات براي $0 \leq a \leq 1$ خيلي كم است. به طور خلاصه استفاده ي مركب نمونه گيري تصادفي ساده و بروردگر تفاضل بر اساس اصل تصميم بيز استوار تحت مدل (۷) با $a = 0$ بر نمونه گيري سيستماتيک بهينه ترجيح داده مي شود و اين انتخاب نسبت به انحراف از اين فرض استوار است.

جدول ۲- نتایج شبیه‌سازی برای نمونه‌گیری با احتمالات برابر بر حسب درصد: نمونه‌گیری سیستماتیک بر اساس U_{opt} در مقابل استفاده‌ی توأم از نمونه‌گیری تصادفی ساده و برآوردگر تفاضل

$a = 1$			$a = 0.5$			$a = 0$			
CV_{srs}	CV_{sys}	RE	CV_{srs}	CV_{sys}	RE	CV_{srs}	CV_{sys}	RE	
۱۹	۴۸	۱۱۲	۱۸	۴۸	۱۰۷	۱۶	۴۷	۱۰۴	$\sigma = 0.1$
۲۰	۴۷	۱۱۲	۱۷	۴۷	۱۰۷	۱۶	۴۹	۹۹	$\sigma = 0.1$

RE : نسبت بین متوسط واریانس نمونه‌گیری سیستماتیک و نمونه‌گیری تصادفی ساده.

CV : ضریب تغییرات واریانس‌های نمونه‌گیری.

۵- نمونه‌گیری سیستماتیک برای چندین موقعیت

یک نمونه‌ی سیستماتیک، پس از این که انتخاب شد، ممکن است برای چندین موقعیت متعاقب مورد استفاده قرار بگیرد. چنین نمونه‌ای ممکن است تشکیل یک گروه در یک طرح پانلی چرخشی بدهد، مانند طرح LFS (آمارگیری نیروی کار) در اغلب کشورها. همچنین می‌تواند هسته‌ی اصلی یک آمارگیری پانلی باشد، با واحدهای تکمیلی که گاه‌گاه برای به حساب آوردن بازسازی طبیعی جامعه به نمونه اضافه می‌شوند. برای ساده کردن بحث در این جا فرض می‌کنیم که یک نمونه‌ی سیستماتیک در موقعیت اول انتخاب شده و پیش از کنار گذاشته شدن برای تمامی موقعیت‌های بعدی مورد استفاده قرار گرفته است، و جامعه‌ی U در طول این دوره یکسان باقی مانده است.

در این صورت نتایج (۴) و (۶) مستقیماً برای کل دوره‌ی فعال پانل صدق می‌کنند. به صورت واضح‌تر، فرض کنید $y_i = (y_{i1}, \dots, y_{it}, \dots, y_{i4})^T$ متغیر مورد بررسی متناظر با $i \in U$ باشد. نتایج (۴) و (۶) مستقیماً برای هر تابعی از y_i صدق می‌کنند. برای مثال، فرض کنید y_i از ۴ معیار وضعیت اشتغال در هر یک از ۴ فصل سال تقویمی در آمارگیری نیروی کار تشکیل شده باشد. اشتغال متوسط سالیانه از میانگین y_{i1} تا y_{i4} حاصل می‌شود. با انتخاب یک نمونه‌ی سیستماتیک در موقعیت اول، خطر نوسان واریانس برآوردگر نرخ اشتغال متوسط سالیانه و همچنین خطر نوسان واریانس در هر فصل را می‌پذیریم. علاوه بر این یکی از کاربردهای مهم داده‌های پانلی، برآورد تغییرات در جامعه

است. فرض کنید $z_i = (z_{i2}, \dots, z_{iT})^T$ ، که در آن برای $t = 2, \dots, T$ ، $z_{it} = y_{it} - y_{i,t-1}$ تغییرات دوره به دوره است. باز هم نتایج (۴) و (۶) مستقیماً برای هر z_{it} صدق می‌کنند، به نحوی که برآورد تغییرات ممکن است به خاطر نمونه‌گیری سیستماتیک مخاطره‌ی بیزی مرتبه‌ی دوم زیادی داشته باشد.

ملاحظات بالا خودهمبستگی قوی محتمل داخل y_i را که اغلب در جامعه‌های طبیعی دیده می‌شود، در نظر نمی‌گیرند. به علاوه؟، یک آزمون شرطی نیز مورد نیاز است. ساده‌ترین حالت را که در آن $T = 2$ و y_{it} یک متغیر رشته‌ای مانند وضعیت اشتغال است، در نظر بگیرید. به عنوان یک مدل ساده از وابستگی بین y_{i2} و y_{i1} ، احتمالات تغییر وضعیت مارکوف p_{ab} را برای $y_{i2} = b$ به شرط $y_{i1} = a$ ، مستقل برای $i \in U$ در نظر می‌گیریم. به شرط $y_{i1} = a$ این به یک مدل جامعه‌ی همگن (۳) منتهی می‌شود. میانگین نمونه‌ی سیستماتیک $z_i = y_{i2} - y_{i1}$ به صورت زیر داده می‌شود

$$\bar{z}_m = \sum_{a: y_{i1}=a} \frac{n_a}{n} \bar{z}_{m,a}$$

که در آن n_a تعداد واحدهای با $y_{i1} = a$ است و $\bar{z}_{m,a}$ میانگین تغییرات میان آن‌ها است. صورت بسته‌ی واریانس شرطی $(\bar{z}_m | \{y_{i1}; i \in U\})$ به طور کلی پیچیده است. در عوض هر ترتیبی را که در آن واحدها با توجه به مقدار y_{i1} تقسیم‌بندی شده‌اند در نظر بگیرید. فرض کنید $\frac{N_a}{k}$ به طور طبیعی و به ازای کلیه‌ی a ها یک مقدار صحیح است که در آن N_a تعداد واحدهای با $y_{i1} = a$ در جامعه است. در این صورت n_a و $\bar{y}_{m,t=1}$ هر دو ثابت‌های نمونه‌گیری هستند، به نحوی که واریانس \bar{z}_m به سادگی، واریانس $\bar{y}_{m,t=2}$ است. حال نتیجه‌ی (۴) می‌تواند برای $\bar{y}_{m,a,t=2}$ صدق کند، یعنی داخل هر بخش از y_{i1} تحت مدل تغییر وضعیت مارکوف به نحوی که مخاطره‌های بیزی مرتبه‌ی دوم $\bar{y}_{m,a,t=2}$ به شرط $\{y_{i1}; i \in U\}$ مستقیماً به سمت \bar{z}_m می‌رود. در نظر گرفتن این حالت خاص نشان می‌دهد که مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری سیستماتیک می‌تواند برای برآوردگرهای تغییرات در جامعه‌های خودهمبسته و همچنین زمانی که واریانس به صورت شرطی ارزیابی می‌شود زیاد باشد. این موضوع را به وسیله‌ی یک مطالعه‌ی شبیه‌سازی در بخش بعد آزمون می‌کنیم.

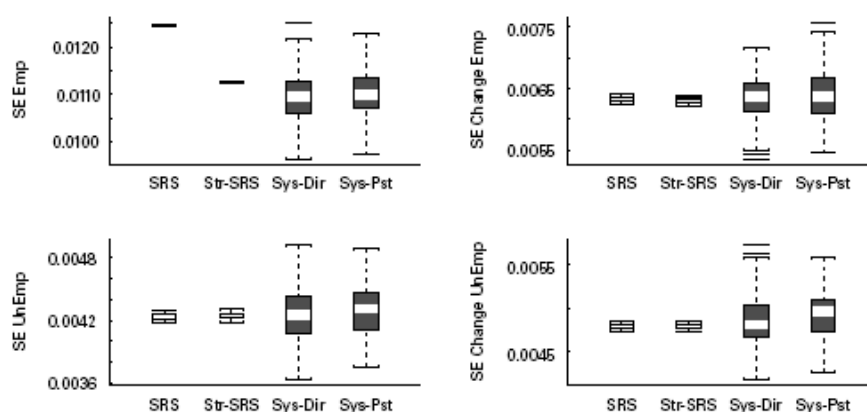
۶- شبیه‌سازی: پویایی بازار کار

در این بخش پویایی بازار کار را با استفاده از سرشماری سال ۲۰۰۱ نیروژ و آمارگیری نیروی کار نیروژ شبیه‌سازی می‌کنیم. از سرشماری سال ۲۰۰۱ وضعیت اشتغال را به دست می‌آوریم که به سه دسته‌ی « شاغل»، « بیکار» و « غیر فعال» طبقه‌بندی شده است و متغیر مورد نظر جامعه در $t=1$ در نظر گرفته می‌شود. سپس از آمارگیری نیروی کار در فصل آخر سال ۲۰۰۴ و فصل اول سال ۲۰۰۵ یک ماتریس تغییر وضعیت 3×3 برای وضعیت اشتغال بین دو فصل به دست می‌آوریم. با استفاده از این احتمالات تغییر وضعیت مارکوف قادریم یک وضعیت اشتغال در جامعه را در $t=2$ شبیه‌سازی کنیم. جامعه در داخل هر یک از ۱۹ شهرستان نیروژ به تفکیک شهر، سن، جنس و شماره‌ی شناسایی فردی (PIN) طبقه‌بندی شده است که PIN می‌تواند با وضعیت اشتغال مورد نظر ناهمبسته در نظر گرفته شود.

چهار راهبرد متفاوت را در نظر می‌گیریم: ۱- نمونه‌گیری سیستماتیک با احتمالات برابر در $t=1$ و برآورد بر اساس وزن‌دهی مستقیم که با $Sys-Dir$ نشان داده می‌شود، ۲- نمونه‌گیری تصادفی ساده در $t=1$ و برآورد بر اساس وزن‌دهی مستقیم که با SRS نشان داده می‌شود، ۳- نمونه‌گیری تصادفی طبقه‌بندی‌شده با تخصیص متناسب نسبت به جنس و سن (کلاً ۲۲ گروه) که به برآورد طبقه‌بندی‌شده منتهی و با $Str-SRS$ نمایش داده می‌شود، و ۴- نمونه‌گیری سیستماتیک با احتمالات برابر و برآورد پس طبقه‌بندی با ۲۲ گروه سنی-جنسیتی به‌عنوان پس طبقه، که با $Sys-Pst$ نمایش داده می‌شود.

شبیه‌سازی‌ها برای هر یک از ۱۹ شهرستان نیروژ به صورت جداگانه انجام می‌شود تا بازتابی از طرح طبقه‌بندی شده‌ی آمارگیری نیروی کار نیروژ باشد. یک نمونه‌ی انتخاب شده در $t=1$ ، در $t=2$ نیز مورد استفاده قرار می‌گیرد و اندازه‌ی نمونه‌های داخل شهرستان از آمارگیری نیروی کار نیروژ گرفته شده‌اند. نتایج برای همه‌ی شهرستان‌ها خیلی شبیه است. در این جا ما وضعیت را تنها برای Østfold در شکل ۱ نشان می‌دهیم. در این حالت نمونه‌گیری سیستماتیک می‌تواند به عنوان یک طبقه‌بندی تلویحی نسبت به شهر، سن و جنسیت در نظر گرفته شود. اثرات طبقه‌بندی فقط برای نرخ اشتغال در $t=2$ قابل توجه‌اند (Emp)، که در حدود ۲۰ درصد کاهش واریانس را در مقایسه با SRS

نشان می‌دهد. هرچند، بیش‌تر اثر طبقه‌بندی می‌تواند با طبقه‌بندی تنها بر اساس جنسیت و سن به دست آید. توجه کنید که طبقه‌بندی بر اساس شهر علاوه بر جنسیت و سن، به خاطر تعداد زیاد طبقات غیرکاربردی است. در همه‌ی حالت‌های دیگر نمی‌توان انتظار بهبود کارایی را با استفاده از نمونه‌گیری سیستماتیک داشت.



شکل ۱- نمودار جعبه‌ای انحراف معیار (SE) نرخ اشتغال در $t = 2$ (Emp)، تغییر نرخ اشتغال (Change Emp)، نرخ بیکاری در $t = 2$ (UnEmp) و تغییر نرخ بیکاری (Change UnEmp) برای شهرستان Østfold؛ نمونه‌گیری تصادفی ساده با وزن‌دهی مستقیم (SRS)، نمونه‌گیری تصادفی طبقه‌بندی‌شده با تخصیص متناسب (Str-SRS)، نمونه‌گیری سیستماتیک (Sys-Dir) و نمونه‌گیری سیستماتیک با برآورد پس طبقه‌بندی (Sys-Pst)

مشاهده می‌شود که در حالی که مخاطره‌های بیزی مرتبه‌ی دوم SRS و $Str-SRS$ برای جامعه‌ای به این اندازه (در حدود صد و هفتاد و نه هزار نفر) ناچیز است، این مخاطره‌ها تحت نمونه‌گیری سیستماتیک همچنین زمانی که واریانس‌ها مانند این‌جا به صورت شرطی محاسبه می‌شوند، قابل ملاحظه‌اند. ضریب تغییرات V_{sys} برای Emp ۱۱ درصد، برای تغییر Emp ۱۵/۸ درصد، برای $UnEmp$ ۱۶/۴ درصد و برای تغییر $UnEmp$ ۱۵/۴ درصد است. این مقادیر با مقدار غیر شرطی نظری که توسط (۴) داده می‌شوند و برای Østfold تقریباً برابر است با $\sqrt{2f} = \sqrt{\frac{2}{134/6}} = 12/2\%$ ، قابل مقایسه‌اند. لذا در موقعیت‌های خاص، نوسان واریانس ممکن است اثر مورد انتظار

طبقه‌بندی روی برآورد Emp را از بین ببرد. توجه کنید که مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری سیستماتیک نمی‌تواند به وسیله‌ی پس طبقه‌بندی کاهش داده شود. به خصوص برای برآورد تغییرات که در این جا نگرانی اصلی ما است ضریب تغییرات واریانس نمونه‌گیری سیستماتیک تقریباً برابر حالت برآورد سطح است. لذا استفاده از نمونه‌گیری سیستماتیک می‌تواند باعث شود که واریانس نمونه‌گیری واقعی یک برآوردگر تغییرات، به شدت در طول زمان تغییر کند. برای مثال، اگر واریانس واقعی ۱۵ درصد بالاتر از امید ریاضی بین فصل اول و دوم و ۱۵ درصد پایین‌تر از امید ریاضی بین فصل دوم و سوم باشد، آن‌گاه دو برآورد تغییرات دارای ۳۰ درصد اختلاف در واریانس‌های نمونه‌گیری هستند که از استفاده‌ی تنها از نمونه‌گیری سیستماتیک ناشی شده است. حال که در این جا ضریب تغییرات واریانس هر دوی برآوردگرهای تغییرات در حدود ۱۵ درصد است، این یک سناریوی غیر معمول نیست. در حالت افراطی‌تری که دو انحراف معیار بالا یا پایین‌تر از واریانس نمونه‌گیری ورد انتظارند، واریانس واقعی یک برآوردگر تغییرات تقریباً دو برابر (یعنی $\frac{\sqrt{x}}{\sqrt{y}}$) دیگری است. قطعاً نگاه‌داشتن این مقدار به عنوان یک ویژگی طرح نمونه‌گیری نامطلوب است.

۷- خلاصه

در قسمت‌های قبل مفهوم مخاطره‌ی بیزی مرتبه‌ی دوم و اصل تصمیم استوار بیز را معرفی کردیم. تعدادی از وضعیت‌هایی را در نظر گرفتیم که در آن‌ها نمونه‌گیری سیستماتیک به طور معمول جایگزین روش‌های نمونه‌گیری تصادفی دیگری می‌شود که به همان اندازه کارا هستند. نشان داده شده است که این کار می‌تواند مخاطره‌های بیزی مرتبه‌ی دوم بزرگ، یعنی نوسانات واریانس نمونه‌گیری واقعی در هر دوی آمارگیری نمونه‌ای مقطعی و طولی ایجاد کند. این می‌تواند برای نمونه‌های بزرگ گرفته شده از جامعه‌های کوچک، یا زیرنمونه‌های بزرگ از زیرجامعه‌های کوچک، بسیار زیان‌آور باشد. استفاده از نمونه‌گیری سیستماتیک برای راحتی در چنین موقعیت‌هایی بر پایه‌ی شانس و بدون هیچ‌گونه توقعی برای بهبود کارایی است.

نمونه‌گیری سیستماتیک به فراوانی در موقعیت‌هایی که ما مورد بررسی قرار نداده‌ایم نیز به کار می‌رود. کاکران [۲] مثال‌های متعددی، شامل استفاده‌ی معمول از نمونه‌گیری سیستماتیک یک یا دوبعدی در آمارگیری‌های جنگل‌داری و اراضی ذکر کرده است. چنین موقعیت‌هایی می‌توانند مشابه آنچه در این جا انجام شده است مورد مطالعه قرار بگیرند، اما نیازمند مدل‌های جامعه‌ی نسبتاً خاص شامل همبستگی‌هایی در طول زمان و مکان‌اند که خارج از چارچوب این مقاله‌اند. به علاوه، همان‌طور که قبلاً اشاره شد، جمع‌آوری یک نمونه‌ی سیستماتیک ممکن است در چنین موقعیت‌هایی هنوز هم آسان‌تر باشد. سرانجام، تعادل در نمونه‌گیری سیستماتیک می‌تواند در آمارگیری‌های اراضی و بوم‌شناسی نقش برجسته‌ای ایفا کند. برای مثال گستردگی یکنواخت نواحی نمونه‌گیری برای هموارسازی فضایی داده‌ها ممکن است مهم‌تر از تصادفی بودن نمونه تلقی شود.

ما نمونه‌گیری سیستماتیک را از نقطه نظر یک تصمیم آماری مورد مطالعه قرار داده‌ایم که در آن تابع زیان به صورت واریانس نمونه‌گیری برآوردگر آمارگیری تعریف می‌شود. توجه کنید که برای استنباط مدل-مبنا که در آن واریانس یک برآوردگر تنها تحت مدل جامعه ارزیابی می‌شود مخاطره‌ی بیزی مرتبه‌ی دوم نمونه‌گیری سیستماتیک با مخاطره‌ی بیزی مرتبه‌ی دوم یک روش نمونه‌گیری تصادفی دیگر تفاوتی ندارد به شرط آن که نمونه‌گیری در هر دو حالت فاقد اطلاعات مفید باشد. در واقع نمونه‌گیری سیستماتیک گاهی به‌عنوان گام اول در ساختن نمونه‌های متعادل مختلف مفید در نظر گرفته می‌شود [۱۱].

مرجع‌ها

- [1] Berger, J.O. (1985). *Statistical Decision Theory and Bayesian Analysis*, (2nd ed.). Springer. NY: New York.
- [2] Cochran, W.G. (1977). *Sampling Techniques*, (3rd ed.). New York: John Wiley and Sons.
- [3] French, S. (1986). *Decision Theory: An Introduction to the Mathematics of Rationality*. Chichester: Ellis Horwood Ltd.
- [4] Hartley, H.O. and Rao, J.N.K. (1962). Sampling with Unequal Probabilities and Without Replacement. *Annals of Mathematical Statistics*, **33**, 350-374.

- [5] Madow, W.G. (1949). On the Theory of Systematic Sampling, II. *Annals of Mathematical Statistics*, **20**, 333-354.
- [6] Madow, W.G. (1953). On the Theory of Systematic Sampling, III. *Annals of Mathematical Statistics*, **24**, 101-106.
- [7] Madow, W.G. and Madow, L.H. (1944). On the Theory of Systematic Sampling. *Annals of Mathematical Statistics*, **15**, 1-24.
- [8] Ohlsson, E. (1998). Sequential Poisson Sampling. *Journal of Official Statistics*, **14**, 149-162.
- [9] Raj, D. (1965). Variance Estimation in Randomized Systematic Sampling with Probability Proportionate to Size. *Journal of the American Statistical Association*, **60**, 278-284.
- [10] Särndal, C.E., Swensson, B., and Wretman, J. (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag.
- [11] Valliant, R., Dorfman, A.H., and Royall, R.M. (2000). *Finite Population Sampling and Inference*. New York: Wiley.
- [12] Wetherill, B.G. (1981). *Intermediate Statistical Methods*: Chapman and Hall: London.
- [13] Zigmond, A.S.; Snaith, R.P. (1983). The Hospital Anxiety and Depression Scale. *Acta Psychiatrica Scandinavica*, **67**, 361-370.

شیرین گلچی

فوق لیسانس آمار

تهران، خیابان دکتر فاطمی، خیابان باباطاهر، خیابان شهید فکوری، شماره ۱۴۵، پژوهشکده‌ی آمار.

پیام‌نگار: golchishirin@yahoo.com

