

استفاده از روش چنگکزنی تعمیم یافته به منظور ایجاد سازگاری در جدول‌های حاصل از سرشماری

آرمان بیداربخت‌نیا

مرکز آمار ایران

چکیده. روش‌های تعديل وزنی، اغلب در آمارگیری‌های نمونه‌ای به منظور سازگار کردن توزیع‌های نمونه‌ای با توزیع‌های جامعه به کار می‌رود. این روش‌ها، علاوه بر آمارگیری‌های نمونه‌ای، ممکن است در سرشماری همراه با نمونه‌گیری نیز مورد استفاده قرار گیرند. در سرشماری همراه با نمونه‌گیری، مستله سازگاری جدول‌ها از اهمیت خاصی برخوردار بوده و مستلزم استفاده از روش‌های برآورد پیچیده‌ای است. این مقاله به معرفی دو روش پس‌طبقه‌بندی و چنگکزنی و چگونگی استفاده از این روش‌ها به منظور ایجاد سازگاری در جدول‌های حاصل از سرشماری عمومی نفوس و مسکن ۱۳۸۵ می‌پردازد.

۱- مقدمه

در آمارگیری‌های خانواری، اغلب برای جبران اختلافات ناشی از بی‌پاسخی و نقص چارچوب، از روش‌های تعديل وزنی برای برآورد پارامتر(های) جامعه استفاده می‌شود. در این روش‌ها وزن‌های پایه‌ای طرح با استفاده از اطلاعات کمکی موجود به گونه‌ای تعديل می‌شوند که توزیع‌های نمونه‌ای با توزیع‌های جامعه‌ای همگون شود. اطلاعات کمکی از منابع مختلفی فراهم می‌شود. داده‌های ثبتی، پیش‌بینی‌های جمعیتی، نتایج آمارگیری‌های معتبر مثل نیروی کار و هزینه و درامد خانوار و اطلاعات حاصل از مرحله اول در آمارگیری‌های چند واژگان کلیدی: تعديل وزنی؛ وزن پایه‌ای طرح؛ برآوردگر هورویتز- تامپسون؛ سرشماری همراه با نمونه‌گیری.

مرحله‌ای، از جمله منابعی است که می‌توان اطلاعات کمکی را از آن‌ها استخراج کرد. روش‌های متعددی برای انجام این تعدیل‌ها به کار می‌رود. پس طبقه‌بندی، چنگکازنی (Raking)، مدل‌بندی رگرسیون لوزتیکی و برآورد رگرسیونی تعمیم‌یافته، روش‌هایی هستند که با توجه به اهداف و نوع داده‌های کمکی می‌توان برای انجام تعدیل وزنی به کار برد. روش‌های نامبرده همگی عضو خانواده‌ای از برآوردهای بمنام برآوردهای کالیبره (Calibration Estimators) هستند.

تعديل‌های وزنی به واسطه کارکرد آن‌ها در ایجاد سازگاری بين توزيع‌های نمونه‌ای و توزيع‌های جامعه‌ای، علاوه بر آمارگیری‌های نمونه‌ای، ممکن است در سرشماری‌ها نيز مورد استفاده قرار گيرند.

مرکز آمار ايران تصميم دارد تا به منظور بالا بردن كيفيت نتایج حاصل از سرشماري عمومي نفوس و مسكن، از طريق کاهش خطاهای غير نمونه‌گيری و همچنین ارائه اطلاعات بيشتر در نواحي کوچک جغرافيايی، برای زيرگروههای مختلف جمعيتي و در عين حال کاهش بار پاسخگو در سرشماري عمومي نفوس و مسكن ۱۳۸۵، اين سرشماري را به روش سرشماري هماه با نمونه‌گيری انجام دهد. به اين ترتيب برخى از اطلاعات (اقلام عمومي) برای تمام افراد، خانوارها و واحدهای مسکونی در جامعه جمع‌آوري می‌شوند و اطلاعات دیگری تحت عنوان اقلام تفصيلي فقط از بخشى از افراد، خانوارها و واحدهای مسکونی که به عنوان نمونه انتخاب شده‌اند، جمع‌آوري می‌گردد. به اين ترتيب در سرشماري عمومي نفوس و مسكن ۱۳۸۵ از سه نوع فرم برای جمع‌آوري اطلاعات مربوط به خانوارها استفاده می‌شود. فرم ۲ (پرسشنامه خانوار- اقلام عمومي) که فقط شامل اطلاعات عمومي است. فرم ۳ (پرسشنامه خانوار- اقلام عمومي و تفصيلي) که علاوه بر اقلام عمومي، اقلام تفصيلي را نيز شامل می‌شود و فرم ۴ (پرسشنامه خانوار مؤسسه‌اي).^۱ نحوه انجام سرشماري به اين ترتيب است که در هر يك از حوزه‌های سرشماري (محدوده کاريک مامور سرشماري)، برای هر کدام از خانوارهای مؤسسه‌اي، فرم ۴، برای هر کدام از خانوارهای معمولی غير ساكن، فرم ۳ و برای خانوارهای معمولی ساكن و خانوارهای گروهي، همزمان با تهيه فهرست خانوارهای حوزه، بر اساس يك فاصله نمونه‌گيری از پيش تعين شده و به روش سيستماتيك خطى، خانوارهایي به عنوان

خانوارهای نمونه انتخاب شده و برای این خانوارها فرم ۳ و برای بقیه خانوارهای حوزه فرم ۲ تکمیل می‌شود.

واضح است که در جدول‌های حاصل از نتایج سرشماری، جدول‌های مربوط به اقلام عمومی به طور مستقیم از فایل مشکل از مجموعه فرم‌های ۲ و ۳ و ۴ استخراج می‌شود، در حالی که جدول‌های مربوط به اقلام تفصیلی و همچنین تقاطع اقلام عمومی با تفصیلی (دو یا چند طرفه) فقط با استفاده از اطلاعات فرم ۳ و با به کارگیری روش‌های براورد، استخراج می‌شود.

آن‌چه در این جدول‌ها بیش از هر عامل دیگری اهمیت دارد، سازگاری آن‌ها با جدول‌هایی است که فقط از اقلام عمومی و به طور مستقیم حاصل شده‌اند. به این ترتیب در براوردهای انجام شده برای استخراج نتایج سرشماری عمومی نفوس و مسکن ۱۳۸۵، از تعديل‌های وزنی مناسب استفاده می‌شود که اطلاعات کمکی برای انجام این تعديل‌ها، شمارش‌های حاصل از فرم‌های ۲ و ۳ برای اقلام عمومی در هر یک از سطوح و یا ترکیب‌هایی از سطوح اقلام عمومی است.

دو روشی که در سرشماری عمومی نفوس و مسکن ۱۳۸۵ برای براورد اطلاعات مربوط به اقلام تفصیلی مورد استفاده قرار می‌گیرند، عبارتند از پس‌طبقه‌بندی و چنگکزی تعیین یافته که در این مقاله علاوه بر معرفی هر دو روش، به بررسی مفصل چنگکزی تعیین یافته در قالب یک مثال پرداخته می‌شود، ضمن این که مراحل محاسباتی این روش به‌گونه‌ای که برای برنامه‌نویسی نیز آسان باشد، شرح داده شده است.

۲- پس‌طبقه‌بندی

جامعه U با N عضو را در نظر بگیرید که نمونه S با n عضو از آن انتخاب شده است. پاسخ ز امین واحد جامعه را با $y_i, i=1, \dots, N$ و پاسخ اخذ شده از λ امین واحد نمونه را با $y_{\lambda}, \lambda=1, \dots, n$ نشان می‌دهیم. هدف، براورد مقدار کل جامعه برای متغیر مورد بررسی Y است. فرض کنید پس از انجام نمونه‌گیری جامعه را به G طبقه تقسیم‌بندی کرده و اندازه جمعیت در طبقه g ام را با N_g و اندازه نمونه متعلق به آن را با n_g نشان می‌دهیم که $G = g_1, \dots, g_m$. مثلاً اگر جمعیت فعلی کشور را بر اساس همه ترکیب‌های

ممکن حاصل از سطوح متغیرهای جنس (زن و مرد) وضع سواد (باسواد و بی‌سواد) و وضع فعالیت (شاغل و بیکار) طبقه‌بندی کنیم، خواهیم داشت:

اگر π احتمال انتخاب شدن عضو نام در نمونه باشد، وزن پایه برای عضو نام عبارت است از:

$$d_i = \frac{1}{\pi_i}$$

که این عدد در سرشماری همراه با نمونه‌گیری تقریباً برابر با فاصله نمونه‌گیری است. به این ترتیب براورد پس‌طبقه‌بندی برای مقدار کل جامعه، $t_y = \sum_{j=1}^N Y_j$ ، عبارت است از:

$$(1) \quad \hat{t}_{y(pos)} = \sum_{g=1}^G \frac{N_g}{\hat{N}_g} \hat{t}_{y(g)}$$

که در آن $\hat{t}_{y(g)}$ براوردگر هوروتیز-تامپسون برای مقدار کل متغیر Y در پس‌طبقه g ام و $\hat{N}_g = \sum_{i=1}^{n_g} d_{gi}$ براورد تعداد در پس‌طبقه g ام است و d_{gi} عبارت است از: وزن پایه عضو نام در پس‌طبقه g ام.

با توجه به این که با تمام اقلام سرشماری به شکل متغیرهای گستته برخورد می‌شود، لذا متغیر Y در روابط بالا همواره به صورت یک متغیر دو حالتی است که به شکل زیر تعریف می‌شود.

$$Y_i = \begin{cases} 1 & \text{نمونه نام دارای صفت مورد نظر باشد} \\ 0 & \text{نمونه نام دارای صفت مورد نظر نباشد} \end{cases}$$

به این ترتیب براورد جمعیت دارای صفت مورد نظر در پس‌طبقه g ام به صورت $\hat{t}_{y(g)} = \sum_{i=1}^{n_g} d_{gi}$ محاسبه می‌شود که n'_g عبارت است از تعداد عناصر نمونه متعلق به پس‌طبقه g ام که دارای صفت مورد نظر هستند و d_{gi} وزن پایه برای واحد نمونه نام در پس‌طبقه g ام است.

این روش را به دلیل این که برای تمام ترکیب‌های ممکن از سطوح متغیرهای پس‌طبقه‌بندی سازگاری ایجاد می‌کند، روش پس‌طبقه‌بندی کامل (Complete Poststratification) یا وزن‌دهی خانه‌ای (Cell Weighting) گویند. در مواردی اندازه نمونه در برخی از پس‌طبقه‌ها صفر است یک روش آن است که برخی پس‌طبقه‌ها در یکدیگر ادغام می‌شوند.

اما روشی که بیشتر مورد استفاده قرار می‌گیرد و کمتر با مشکل اندازه نمونه صفر مواجه است، روش «پس‌طبقه‌بندی ناقص» (Incomplete Poststratification) است که پس‌طبقه‌بندی را روی حاشیه‌ها انجام می‌دهد. این روش به طور کلی شبیه به روش پس‌طبقه‌بندی است، با این تفاوت که در این روش فقط برای حاشیه‌های یک متغیر سازگاری ایجاد می‌شود. مزیت این روش بر روش پس‌طبقه‌بندی کامل، طبق آن‌چه گفته شد، این است که در اینجا کمتر با مشکل اندازه نمونه صفر مواجه هستیم. برای نمونه، فرض کنید در مثال قبل فقط برای متغیر جنس (مرد و زن)، سازگاری مورد نیاز باشد. به این ترتیب $2 = G$ ، به این ترتیب برخلاف مثال قبل که برای تمام ترکیب‌های ممکن حاصل از سطوح سه متغیر جنس، وضع سواد و وضع فعالیت نیاز به ایجاد سازگاری بود، در این حالت که فقط سازگاری برای متغیر جنس لازم است، احتمال این که اندازه نمونه در سطوح مورد نظر برای سازگاری (هر یک از سطوح متغیر جنس) صفر باشد خیلی کمتر است.

در پس‌طبقه‌بندی برای حاشیه‌ها اگر پس‌طبقه‌بندی شامل بیش از یک متغیر باشدند مثلاً در مورد ذکر شده، اگر بخواهیم به طور همزمان برای حاشیه‌های سه متغیر سازگاری ایجاد شود، یعنی $6 = 2 + 2 + 2 = G$ ، دیگر روش پس‌طبقه‌بندی قابل استفاده نیست. برای ایجاد سازگاری با شرایط بالا، در سرشماری عمومی نفوس و مسکن ۱۳۸۵ از روش چنگکزی تعمیم‌یافته استفاده شده است که در ادامه معرفی می‌شود.

شایان ذکر است که در سرشماری عمومی نفوس و مسکن ۱۳۸۵، در مورد جدول‌های خانواری و واحد مسکونی، به دلیل نیاز به سازگاری در تمام ترکیب‌های ممکن از سطوح متغیرهای پس‌طبقه‌بندی، برای ایجاد سازگاری در جدول‌ها از روش پس‌طبقه‌بندی کامل استفاده شده است. البته مشکل خانه‌های با اندازه نمونه صفر در این

مورد نیز همچنان باقی است، که با ادغام برخی از سطوح می‌توان آن را برطرف نمود.

۳- چنگکزنی

در پس‌طبقه‌بندی برای حاشیه‌ها، اگر بیش از یک متغیر پس‌طبقه‌بندی داشته باشیم، روشی مشابه با پس‌طبقه‌بندی ناقص، ولی به صورت تکراری مورد استفاده قرار می‌گیرد. این روش را چنگکزنی یا برازش متناسب تکراری (Iterative Proportional Fitting) گوییم.

فرض کنید بخواهیم برای حاشیه‌های متغیرهای جنس و وضع سواد به‌طور همزمان سازگاری ایجاد شود. این کار باید به صورت تکراری و در چند مرحله انجام گیرد. در مرحله اول، برآوردهای پس‌طبقه‌بندی ناقص را روی حاشیه‌های سط्रی (یکی از متغیرهای جنس و وضع سواد) به دست می‌آوریم تا برای حاشیه‌های آن متغیر، سازگاری ایجاد شود. در این صورت مشاهده می‌شود که برای حاشیه‌های ستونی، سازگاری وجود ندارد. در مرحله دوم، برآوردهای پس‌طبقه‌بندی ناقص را برای حاشیه‌های ستونی به دست می‌آوریم تا برای آن‌ها نیز سازگاری ایجاد شود. با انجام مرحله دوم، سازگاری حاصل از مرحله نخست از بین می‌رود، لذا در مرحله سوم، مجدداً برای حاشیه‌های سطري، برآوردهای پس‌طبقه‌بندی ناقص را محاسبه می‌کنیم. این فرایند به‌طور یک در میان روی حاشیه‌های سطري و ستونی تا جایی ادامه می‌یابد که برای تمام حاشیه‌ها، سازگاری ایجاد شود.

در یک آمارگیری به وسعت سرشماری همراه با نمونه‌گیری، به‌طور همزمان به‌دلیل ایجاد سازگاری در تعداد زیادی متغیر با سطوح نسبتاً زیاد هستیم. برای رفع این مشکل در سرشماری عمومی نفوس و مسکن ۱۳۸۵ حالت تعمیم‌یافته‌ای از روش چنگکزنی مورد استفاده قرار می‌گیرد.

۴- چنگکزنی تعمیم‌یافته

در این روش به‌دلیل وزن‌هایی هستیم که بر اساس یک تابع فاصله مشخص، کمترین فاصله تا وزن‌های پایه‌ای را داشته و در عین حال برآوردهایی تولید کنند که با جمعیت کل

برای صفت مورد نظر در حاشیه‌ها برابر باشند.

یک ارائه $X_{n \times L}$ را برای متغیرهای کمکی در نظر بگیرید که در آن n تعداد اعضای نمونه و L تعداد سطوحی است که می‌خواهیم برای آن‌ها سازگاری ایجاد شود. برای عضو i نمونه، بردار متغیرهای کمکی به صورت $(x_{iL}, x_{iL-1}, \dots, x_{i1}) = X_i$ تعریف می‌شود. وزن پایه برای عضو i نمونه برابر با d_i و وزن نهایی برابر با $w_i = g_i d_i$ است که w_i را فاکتور تعديل وزنی گوییم. برای به دست آوردن برآوردهای تعديل شده، به دنبال w_i هایی هستیم که علاوه بر مینیمم کردن تابع فاصله

$$(2) \quad \Delta(w, d) = \sum_{i=1}^n \left(w_i \log\left(\frac{w_i}{d_i}\right) - w_i + d_i \right)$$

در معادلات سازگاری $\sum_{i=1}^n w_i X_i = 1$ نیز صدق کنند. λ_i عبارت است از یک بردار

ستونی L بعدی که هر یک از عناصر آن مقدار کل معلوم جامعه برای یکی از سطوح متغیرهای کمکی است. مقادیر w_i مورد نظر پس از حل معادلات لاغرانژ و با استفاده از روش‌های مبتنی بر تکرار، حاصل می‌شود. این تکرارها تا جایی ادامه می‌یابد که سازگاری کامل برقرار شده یا بر اساس معیارهای از پیش تعیین شده به همگرایی مورد نظر بررسیم. روش محاسبه w_i با استفاده از حل معادلات لاغرانژ به گونه‌ای که برای برنامه‌نویسی نیز مناسب باشد، در پیوست آمده است. بر پایه این روش، به کمک رابطه‌های (۵) تا (۷) مقادیر w_i مشخص می‌شوند. به این ترتیب عامل تعديل g_i با استفاده از رابطه زیر به دست می‌آید:

$$(3) \quad g_i = \exp(X_i \lambda_i)$$

که λ_i بردار ضرایب لاغرانژ است. وزن نهایی برای عضو i نمونه به صورت $w_i = g_i d_i$ حاصل می‌شود.

۵- مثال

یک جامعه ۱۰۰۰۰۰ نفری را در نظر بگیرید. فرض کنید ویژگی‌های جنس، سن و وضع

ساد بهترتب در ۲ سطح، ۵ سطح از تمام افراد جامعه مورد پرسش قرار می‌گیرد. فرض کنید برای ۳۰۰۰ نفر از این جامعه، به طور تصادفی علاوه بر ویژگی‌های فوق، توانایی صحبت کردن به زبان فارسی^۱ نیز در ۳ سطح پرسش می‌شود. جدول‌هایی که از اطلاعات مربوط به تمام افراد به طور مستقیم حاصل می‌شود، جدول جنس-سن-جنس- وضع سواد است و جدول‌های حاصل از نمونه که باید برآورد شود، سن-جنس- توانایی صحبت کردن به زبان فارسی و جنس-وضع سواد- توانایی صحبت کردن به زبان فارسی است. شایان ذکر است که ویژگی‌های سن و جنس برای تمام افراد و وضع سواد فقط برای افراد بالای ۶ سال پرسیده شده است. به این ترتیب علاقه‌مند هستیم تا جدول‌های حاصل از برآورد با جدول‌های مستقیم از جامعه سازگاری داشته باشد. به این ترتیب سطوحی که می‌خواهیم برای آن‌ها سازگاری ایجاد شود، عبارتند از $10 = 2 \times 5$ سطح برای جدول سن-جنس و $4 = 2 \times 2$ سطح از جدول جنس-وضع سواد، یعنی $L = 10 + 4 = 14$.

نتایج مستقیم از جامعه در جدول ۱ و جدول ۲ نشان داده شده است.

با توجه به وزن پایه که عبارت است از $d_i = \frac{100}{3} (i=1, \dots, 3000)$ ، برآورده‌گر هوروتیز- تامپسون برای هر یک از جدول‌ها نمونه‌ای در جدول ۳ و جدول ۴ قابل مشاهده است. با مقایسه این دو جدول با جدول ۱ و جدول ۲ مشخص می‌شود، که سرجمع‌های حاصل، با یکدیگر سازگاری ندارند. اما پس از تعديل وزن‌های پایه برای واحدهای نمونه با استفاده از روش چنگکزنی تعمیم‌یافته که نتایج آن در جدول ۵ و جدول ۶ مشاهده می‌شود، خواهیم دید که سازگاری مورد نظر با جدول ۱ و جدول ۲ حاصل شده است. همه محاسبات با برنامه‌نویسی به وسیله نرم‌افزار SAS/IML انجام شده است. برنامه مورد نظر پس از ۳ تکرار به همگرایی رسیده است.

جدول ۱- جمعیت افراد به تفکیک جنس و وضع سواد

کل	وضع سواد		جنس
	با سواد	بی سواد	
۳۶۶۷۰	۲۹۳۰۳	۷۳۶۷	زن
۳۶۸۱۰	۲۹۲۸۶	۷۴۲۴	مرد
۷۳۴۸۰	۵۸۶۸۹	۱۴۷۹۱	کل

جدول ۲- جمعیت افراد به تفکیک جنس و سن

کل	جنس		گروه سنی
	مرد	زن	
۱۲۴۶۳	۶۲۱۴	۶۲۴۹	۱
۴۸۹۰۰	۲۲۴۷۹	۲۶۴۲۱	۲
۳۱۷۴۶	۱۵۸۶۳	۱۵۸۸۳	۳
۶۳۲۴	۳۲۲۸	۳۰۸۶	۴
۵۶۷	۲۸۶	۲۸۱	۵
۱۰۰۰۰	۵۰۰۸۰	۴۹۹۲۰	کل

جدول ۳- برآورد هورویتز- تامپسون تعداد افراد بر حسب سن و توانایی صحبت کردن به زبان فارسی به تفکیک جنس

کل	زن					مرد				
	گروه سنی			توانایی صحبت کردن به زبان فارسی		گروه سنی			توانایی صحبت کردن به زبان فارسی	
	۳	۲	۱	سنی	۱	۲	۳	۱	۲	۳
۵۸۰۰	۶۷	۱۴۲۳	۴۳۰۰	۱	۶۷۰۰	۱۶۷	۱۴۰۰	۵۱۳۳	۱	
۲۴۵۶۷	۳۶۶	۵۶۶۷	۱۸۵۲۳	۲	۲۳۴۵۷	۲۲۲	۵۷۰۰	۱۷۴۲۳	۲	
۱۵۵۳۳	۱۳۳	۲۷۰۰	۱۱۷۰۰	۳	۱۶۹۳۳	۵۳۳	۴۰۲۲	۱۱۸۵۷	۳	
۲۹۲۳	۶۷	۵۳۲	۲۳۲۳	۴	۳۶۰۰	۱۰۰	۹۶۷	۲۵۲۳	۴	
۱۶۷	-	۳۳	۱۲۳	۵	۳۰۰	-	۶۷	۲۲۳	۵	
۴۹۰۰۰	۶۳۳	۱۱۳۶۷	۳۷۰۰۰	کل	۵۱۰۰۰	۱۱۲	۱۲۶۶۷	۳۷۲۰۰	کل	

جدول ۴- برآورد هورویتز- تامپسون تعداد افراد بر حسب وضع سواد و توانایی صحبت کردن به زبان فارسی به تفکیک جنس

زن						مرد					
توانایی صحبت کردن به زبان						توانایی صحبت کردن به زبان					
کل	فارسی			وضع سواد		کل	فارسی			وضع سواد	
	۳	۲	۱				۳	۲	۱		
۷۶۶۷	۳۳	۱۲۲۲	۶۱۰۰	بی‌سواد		۷۶۶۷	۱۶۷	۱۹۲۳	۵۸۰۰	با سواد	
۲۸۶۲۳	۴۳۲	۶۷۲۳	۲۱۴۶۷	باسواد		۲۸۶۲۳	۷۲۳	۷۸۲۳	۲۱۲۳	کل	
۲۶۱۰۰	۴۶۶	۸۰۶۷	۲۷۵۶۷	کل		۳۶۱۰۰	۹۰۰	۹۷۶۷	۲۷۰۳۳	کل	

جدول ۵- برآورد تعديل شده تعداد افراد بر حسب سن و توانایی صحبت کردن به زبان فارسی به تفکیک جنس

زن						مرد					
گروه توانایی صحبت کردن به زبان فارسی						گروه توانایی صحبت کردن به زبان فارسی					
کل	سنی	۳	۲	۱		کل	سنی	۳	۲	۱	
۶۲۴۹	۷۲	۱۵۴۴	۴۶۲۳	۱		۶۲۱۴	۱۵۵	۱۲۹۸	۴۷۶۱	۱	
۲۴۴۲۱	۲۶۵	۵۶۳۶	۱۸۴۲۰	۲		۲۴۴۷۹	۳۴۵	۵۹۳۸	۱۸۱۹۷	۲	
۱۵۸۸۳	۱۳۷	۳۷۸۸	۱۱۹۵۸	۳		۱۵۸۶۳	۵۰۴	۴۲۵۶	۱۱۱۰۳	۳	
۳۰۸۶	۷۱	۵۶۳	۲۴۵۲	۴		۳۲۲۸	۸۸	۸۷	۲۲۸۰	۴	
۲۸۱	۰	۵۶	۲۲۵	۵		۲۸۶	۰	۶۳	۲۲۲	۵	
۴۹۹۲۰	۶۴۵	۱۱۵۸۷	۳۷۶۸۸	کل		۵۰۰۸۰	۱۰۹۲	۱۲۴۲۵	۲۶۵۶۳	کل	

جدول ۶- برآورد تعديل شده تعداد افراد بر حسب وضع سواد و توانایی صحبت کردن به زبان فارسی به تفکیک جنس

زن						مرد					
توانایی صحبت کردن به زبان						توانایی صحبت کردن به زبان					
کل	فارسی			وضع سواد		کل	فارسی			وضع سواد	
	۳	۲	۱				۳	۲	۱		
۷۳۶۷	۳۲	۱۳۱۵	۶۰۲۰	بی‌سواد		۷۲۲۴	۱۵۷	۱۸۲۷	۵۴۴۰	با سواد	
۲۹۳۰۳	۴۴۲	۶۸۸۸	۲۱۹۷۳	باسواد		۲۹۳۸۶	۷۰۹	۷۶۹۶	۲۰۹۸۱	کل	
۳۶۶۷۰	۴۷۴	۸۲۰۳	۲۷۹۹۳	کل		۳۶۸۱۰	۸۶۶	۹۵۲۳	۲۴۶۲۱	کل	

۶- سپاسگزاری

با سپاس فراوان از اعضای محترم کمیته سرشماری توأم با نمونه‌گیری^۳ که کاربردی کردن این روش و استفاده از آن در سرشماری عمومی نفوس و مسکن ۱۳۸۵ بدون راهنمایی‌های ارزشمند و پیشنهادهای سازنده ایشان، میسر نبود.

توضیحات

^۱ برای اطلاع از جزئیات مربوط به فرم‌های مذبور، به راهنمای‌های مربوط در سرشماری عمومی نفوس و مسکن ۱۳۸۵ مراجعه شود.

^۲ در اجرای آزمایشی سال ۱۳۸۴ سرشماری نفوس و مسکن ۱۳۸۵ «توانایی صحبت کردن به زبان فارسی»، یکی از پرسش‌هایی است که در آزمایش پرسیده شده اما در اجرای اصلی این سرشماری منظور نشده است.

^۳ «کمیته سرشماری توأم با نمونه‌گیری» یکی از کمیته‌های گروه تهیه طرح سرشماری عمومی نفوس و مسکن ۱۳۸۵ است.

مرجع‌ها

- [۱] Devill, J. C., and Samdal, C-E. and Sautory. O. (1993), “Generalized Raking Procedures in Survey Sampling,” journal of the American Statistical Association, 88, 1013-1020.
- [۲] Groves, R.M, Dillman, D.A., Eltinge, J.L., Little, R.J.(2002), Survey Nonresponse, New York : Willey.
- [۳] Kalton, G., Flores-Cervantes, I. (2003), “Weighting Methods,” Journal of Official Statistics, Vol. 19, No. 2, 81-97.
- [۴] Singh, A.C., and Mohl, C.A. (1996), “Understanding calibration estimators in survey sampling,” Survey Methodology. 22, 107-115.

پیوست: نحوه محاسبه w_i ‌ها در روش چنگکزنی تعمیم یافته

برای محاسبه w_i ‌ها به شیوه زیر عمل می‌کنیم.

فرض کنید I متغیر کمکی داریم. می‌خواهیم وزن‌هایی تولید کنیم که برآوردهای نمونه‌ای حاصل از آن وزن‌ها برای L سطح حاصل از حاشیه‌های متغیرهای کمکی و یا ترکیب‌هایی از سطوح آن‌ها، با مقدارهای معلوم جامعه سازگار باشد.

ارائه $X_{n \times L}$ را طوری تشکیل می‌دهیم که به‌ازای هر واحد نمونه یک سطر و به‌ازای هر سطح، یک ستون داشته باشد. در هر یک از ستون‌ها، عنصر i ام نمونه ($i = 1, \dots, n$) بسته به این‌که عنصر مزبور به سطح مربوط تعلق داشته یا نداشته باشد، به ترتیب مقدار یک یا صفر می‌گیرد.

d : عبارت است از یک بردار ستونی با n سطر که درایه‌های آن وزن پایه برای اعضای نمونه هستند.

t_x : یک بردار ستونی با L سطر که درایه‌های آن عبارت است از جمعیت حاصل از سرشماری، برای سطح m ام، ($m = 1, \dots, L$).

w : یک بردار ستونی با n سطر که درایه i ام آن برابر است با وزن تولید شده برای نمونه i ام در تکرار τ ام.

x : یک بردار ستونی با L سطر که درایه‌های آن عبارت است از برآورد موزون تعداد افراد برای سطح m ام در تکرار τ ام.

$$(4) \quad \hat{t}_x^{(v)} = X^T w^{(v)}$$

برای محاسبه وزن‌های w باید ابتدا فاکتور f را با استفاده از رابطه زیر به دست آوریم.

$$(5) \quad f^{(v)} = X (X^T X)^{-1} (\hat{t}_x^{(v)} - \hat{t}_x^{(v-1)})$$

در رابطه‌ی بالا داریم:

$\Gamma_{(v-1)}$: یک ارائه قطری $n \times n$ که عنصر i ام روی قطر آن عبارت است از $w_i^{(v-1)}$ که وزن تولید شده برای نمونه i ام در تکرار $(1-v)$ ام است. برای حالتی که $v=1$ باشد باید موارد زیر را در نظر داشت.

$$w^{(1)} = d$$

به این ترتیب داریم:

$$\hat{t}_x^{(1)} = X d$$

Γ_v : ارائه قطری است که عنصر i ام روی قطر آن عبارت است از وزن پایه برای نمونه i ام.

پس از محاسبه $f^{(v)}$ در تکرار v ام، بردار $g^{(v)}$ به صورت زیر محاسبه می‌شود.

$$(6) \quad g^{(v)} = g^{(v-1)} \# \exp(f^{(v)})$$

که نماد $\#$ نشان‌دهنده این است که درایه‌های متناظر دو بردار در یکدیگر ضرب معمولی شود و $\exp(f^{(v)})$ برداری است که هر عنصر آن، تابع \exp عنصر متناظر در $f^{(v)}$ است.

سپس بردار وزن $w^{(v)}$ در تکرار v ام با استفاده از رابطه زیر محاسبه می‌شود.

$$(7) \quad w^{(v)} = d \# g^{(v)}$$

پس از انجام هر تکرار یک شاخص محاسبه می‌شود.

$$d = norm_{(v)} - norm_{(v-1)}$$

$$norm_{(v)} = \| t_x - \hat{t}_x^{(v)} \|$$

که نماد $\| \cdot \|$ نشان‌گر نرم تفاضل هندسی دو بردار است. در صورتی که $d \leq \epsilon$ باشد تکرارها متوقف شده و بردار $w^{(v)}$ در آخرین تکرار، بردار وزن‌های نهایی است که با استفاده از آن می‌توان تمام جدول‌های مربوط را به دست آورد. مقدار ϵ عدد مثبت و کوچکی است که بر اساس عواملی از قبیل زمان اجرای برنامه و میزان سازگاری مورد نظر تعیین می‌شود.



پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتمال جامع علوم انسانی