

ساخت و کاربرد وزن‌های نمونه

زهرا رضایی قهرومدی

مرکز آمار ایران

چکیده. در این مقاله، مراحل مختلف وزن‌دهی و کاربرد آن در محاسبه براورد مشخصه‌های مورد نظر در آمارگیری‌های خانوار ارائه شده است. بهخصوص، تصحیح وزن‌های نمونه به منظور جبران عیوب و نقایص مختلف در نمونه انتخابی بیان شده است. در این مقاله توجه ما محدود به براوردهایی است که در پیش‌تر گزارش‌های آماری تهیه می‌شود که با ارائه مثال‌هایی از نمونه‌گیری‌های واقعی که در کشورهای در حال توسعه و کشورهای دیگری که موقعیت‌های آمارگیری مشابه دارند، به بررسی روش‌های مختلف وزن‌دهی می‌پردازیم.

۱- مقدمه

نمونه‌گیری‌های خانوار، عمدها بر اساس طرح‌های نمونه‌گیری پیچیده به منظور کنترل هزینه است. نتایج نمونه‌گیری دارای ویژگی‌هایی است که ممکن است منجر به اریبی و انحرافات دیگر بین نمونه و جامعه مرجع شود. این ویژگی‌ها شامل انتخاب واحدهای با احتمال‌های نابرابر، پوشش ندادن کل جامعه و بی‌پاسخی است. برای اصلاح این ویژگی‌ها از روش وزن‌دهی استفاده می‌کنیم تا نهایتاً به براوردهای مناسب در مورد مشخصه‌های مورد نظر دست یابیم. به عبارت دیگر، اهداف وزن‌دهی شامل موارد زیر است:

- تصحیح احتمال‌های نابرابر انتخاب؛
- تصحیح بی‌پاسخی واحد؛

- تصحیح توزیع نمونه موزون برای متغیرهای اصلی مورد نظر (مانند وزن، نژاد، جنسیت و...) بهمنظور تطابق با توزیع جامعه.

روش‌های مورد استفاده برای هر یک از حالت‌های فوق در بخش‌های بعدی با جزئیات بیش‌تر ارائه می‌شود. در بخش ۲ وزن‌های نمونه در غالب طرح نمونه‌گیری چند مرحله‌ای شامل تصحیح وزن‌های نمونه، به‌منظور محاسبه واحدهای تکراری و واحدهایی که واجد شرایط بودن آن‌ها نامشخص است، ارائه شده است. بخش ۳ وزن‌دهی را به‌منظور تصحیح احتمال‌های نابرابر انتخاب مورد بررسی قرار داده و در بخش ۴ و ۵ به ترتیب موضوع بی‌پاسخی و روش‌های جبران بی‌پاسخی، عدم پوشش و روش‌های جبران آن به‌همراه مثال‌های عددی بیان شده است. در بخش ۶ افزایش واریانس برآوردها که ناشی از وزن‌دهی است، به‌همراه مثال عددی ارائه شده، در بخش ۷ به مسئله پیرایش وزن‌ها پرداخته شده است و نهایتاً در بخش آخر نتیجه‌گیری از مطالب گفته شده، بیان می‌شود.

۲- ایجاد وزن‌های نمونه

وزن‌دهی اغلب به عنوان اولین مرحله تحلیل داده‌های نمونه‌گیری استفاده می‌شود، که این کار معمولاً با ساخت وزن پایه برای هر واحد نمونه‌گیری به‌منظور تصحیح احتمال‌های نابرابر انتخاب، آغاز می‌شود. در حالت کلی، وزن پایه هر واحد نمونه‌گیری، عکس احتمال انتخاب آن واحد در نمونه است. به این صورت که اگر احتمال وجود یک واحد در نمونه با p_i نمایش داده شود، وزن پایه w_i به صورت $w_i = 1/p_i$ است. به عنوان مثال یک واحد نمونه که با احتمال $1/50$ انتخاب شود، نماینده 50 واحد در جامعه‌ای است که نمونه از آن گرفته شده است.

در طرح‌های چند مرحله‌ای، وزن‌های پایه معرف احتمال انتخاب در هر مرحله است. به عنوان مثال در طرح دو مرحله‌ای که نامین واحد نمونه‌گیری با احتمال p_{ij} در مرحله اول و نامین خانوار درون نامین واحد نمونه‌گیری با احتمال $p_{j(i)}$ در مرحله دوم انتخاب شود،

احتمال انتخاب هر خانوار در نمونه (p_{ij}) از رابطه زیر محاسبه می‌شود.

$$p_{ij} = p_i \times p_{j(i)}$$

بنا بر این عکس احتمال فوق، وزن پایه‌ای است که از تصحیح احتمال‌های نابرابر انتخاب حاصل شده است. به همین ترتیب، اگر $w_{ij,b}$ وزن پایه زامین خانوار، $w_{ij,nr}$ وزن جبران بی‌پاسخی و $w_{ij,nc}$ وزن جبران عدم پوشش در نظر گرفته شود، در این صورت وزن نهایی زامین خانوار از زامین واحد نمونه‌گیری شده از رابطه زیر به دست می‌آید.

$$w_{ij} = w_{ij,b} \times w_{ij,nr} \times w_{ij,nc}$$

شایان ذکر است که اندیس b , nr و nc در وزن‌ها به ترتیب برای نشان دادن وزن پایه، وزن مربوط به بی‌پاسخی و وزن مربوط به عدم پوشش است.

۱-۲- تصحیح وزن پایه برای واجدین شرایط نامعلوم

در مرحله جمع‌آوری داده در آمارگیری‌های خانوار، مثال‌ها و شواهدی وجود دارد که مسأله واجد شرایط بودن خانوارها مطرح می‌شود. به عنوان مثال، یک پرسشگر ممکن است در زمان جمع‌آوری داده هیچ فردی را در واحد مسکونی نمونه‌گیری شده، مشاهده نکند (حتی بعد از مراجعته مجدد). در چنین مواردی مشخص نیست که آیا واحد مسکونی اشغال شده است یا نه. اگر این واحد مسکونی واقعاً اشغال شده باشد، باید به عنوان یک واحد مسکونی بی‌پاسخ طبقه‌بندی شود. در غیر این صورت، این واحد مسکونی واجد شرایط نیست. بعضی مواقع پرسشگر فرض می‌کند که اگر هیچ فردی در واحد مسکونی مشاهده نشود، آن واحد مسکونی اشغال نشده است و در نتیجه واجد شرایط نیست، که این‌ها فرضیات ناصحیحی است که اغلب به اشتباه منجر به نرخ بالای پاسخ می‌شود.

بنا بر این زمانی که واجد شرایط بودن برخی از واحدهای مسکونی نمونه‌گیری شده نامعلوم است، وزنی‌دهی این مسأله را تعديل می‌کند. یک روش ساده، در نظر گرفتن نسبتی از واحدهای مسکونی در نمونه است که وضعیت واجد شرایط بودن یا نبودن آن‌ها مشخص است. از همین نسبت برای واحدهای مسکونی‌ای که واجد شرایط بودن آن‌ها

نامعلوم است، استفاده می‌شود. به عنوان مثال نمونه‌ای از ۳۰۰ واحد مسکونی که وضعیت پاسخ‌گویی آن‌ها در جدول زیر ثبت شده است را در نظر بگیرید.

جدول ۱ - وضعیت پاسخ‌گویی واحدهای مسکونی

تعداد واحدهای مسکونی	طبقه‌بندی جواب‌گویی
۲۱۵	واحدهای با مصاحبه کامل
۲۵	واحدهای واجد شرایط بی‌پاسخ
۱۰	واحدهای غیر واجد شرایط
۵۰	واحدهایی که واجد شرایط بودن آن‌ها نامعلوم است

در بین واحدهای مسکونی که وضعیت واجد شرایط بودن آن‌ها مشخص و معلوم است، نسبت واحدهای مسکونی واجد شرایط $\frac{215+25}{215+25+10} = 0.96$ است. بنا بر این می‌توان فرض کرد که همین نسبت (۰.۹۶) از واحدهای مسکونی که واجد شرایط بودن آن‌ها نامعلوم است، واجد شرایط هستند. به عبارت دیگر، ۹۶ درصد از ۵۰ واحد مسکونی که واجد شرایط بودن آن‌ها نامعلوم است، واجد شرایط هستند. اکنون ضریب زیر برای تعديل وزن واحدهای مسکونی واجد شرایط (واحدهای مسکونی با مصاحبه کامل و واحدهای مسکونی واجد شرایط که پاسخ نداده‌اند)، محاسبه می‌شود:

$$F_{ue} = \frac{\sum_c w_{ij,b} + \sum_{nr} w_{ij,b} + \varepsilon \times \sum_{ue} w_{ij,b}}{\sum_c w_{ij,b} + \sum_{nr} w_{ij,b}} = 1 + \varepsilon \frac{\sum_{ue} w_{ij,b}}{\sum_c w_{ij,b} + \sum_{nr} w_{ij,b}},$$

که در آن ε برآورد نسبت واحدهای واجد شرایطی است که واجد شرایط بودن آن‌ها نامعلوم است. (در مثال فوق $\varepsilon = 0.96$). جمع‌بندی روی c ، nr و ue در فرمول فوق به ترتیب جمع وزن‌های پایه واحدهای مسکونی با مصاحبه کامل، بی‌پاسخ‌های واجد

شرایط و واحدهای مسکونی‌ای است که واجد شرایط بودن آن‌ها نامعلوم است. بنا بر این وزن پایه تعديل‌یافته واحدهای مسکونی واجد شرایط (واحدهای مسکونی با مصاحبه کامل و واحدهای مسکونی واجد شرایط که پاسخ نداده‌اند) با ضرب کردن وزن پایه در فاکتور F_{ue} به دست می‌آید.

۲-۲- تعديل وزن واحدهای تکراری

اگر مشخص شود که بعضی از واحدها در چارچوب تکرار شده‌اند، احتمال انتخاب بالاتر چنین واحدهایی می‌تواند به وسیله فاکتورهای تعديل جبران شود. اغلب واحدهای تکراری بعد از انتخاب نمونه، مشخص شده و احتمال انتخاب چنین واحدهایی به صورت زیر تعديل می‌شود:

فرض کنید n امین واحد نمونه‌گیری دارای احتمال انتخاب p_i است و k رکورد اضافی که در چارچوب نمونه‌گیری توسط این واحد تکراری، مشخص شده است هر یک دارای احتمال انتخاب معلوم $p_{ik}, \dots, p_{i2}, p_{i1}$ است. بنا بر این احتمال انتخاب تعديل‌یافته واحد نمونه‌گیری مورد نظر به صورت زیر تعریف می‌شود:

$$p_i = 1 - (1 - p_{i1})(1 - p_{i2}) \dots (1 - p_{ik})$$

وزن متناظر با این واحد نمونه‌گیری p_i است. اکنون روش‌های ساخت وزن‌های نمونه را به منظور تعديل سه ویژگی احتمال‌های انتخاب نابرابر، بی‌پاسخی و عدم پوشش به همراه چند مثال مورد بحث و بررسی قرار می‌دهیم.

۳- تعديل احتمال‌های انتخاب نابرابر

نمونه‌گیری دو مرحله‌ای شامل انتخاب مناطق به عنوان واحدهای انتخاب شده در مرحله اول و خانوارها به عنوان واحدهای انتخابی مرحله دوم می‌باشد. در مرحله اول، نمونه‌ای با حجم n از N واحد جامعه انتخاب می‌شود و سپس m خانوار از هر یک از M خانوار موجود

در مناطق نمونه‌گیری شده انتخاب می‌شوند. احتمال انتخاب هر خانوار بستگی به تعداد خانوارها در منطقه‌ای که در آن واقع شده است، دارد. فرض کنید، M معرف تعداد خانوارها در n امین منطقه است. بنا بر این احتمال انتخاب یک منطقه، n/N و احتمال شرطی انتخاب یک خانوار در n امین منطقه نمونه‌گیری شده، m/M است. احتمال کل انتخاب یک خانوار از رابطه زیر محاسبه می‌شود:

$$p_{ij} = p_i \times p_{j(i)} = \frac{n}{N} \times \frac{m}{M_i} = \frac{nm}{N} \times \frac{1}{M_i}$$

و وزن خانوار نمونه‌گیری شده بر اساس این طرح نمونه‌گیری به صورت زیر است:

$$w_i = \frac{1}{p_{ij}} = \frac{N}{nm} \times M_i$$

۴- تعدیل وزن برای بی‌پاسخی

این یک امر نادر است که در آمارگیری‌ها، تمام اطلاعات مورد نظر از همه واحدهای نمونه‌گیری به دست آید. به عنوان مثال، برخی از خانوارها ممکن است پاسخگوی هیچ اطلاعاتی نباشند (بی‌پاسخی واحد) در حالی که خانوارهای دیگر ممکن است تنها به تعدادی از سوالات پاسخ دهند (بی‌پاسخی قلم). بنا بر این این مسأله بسیار حائز اهمیت است که موارد بی‌پاسخی را هر چند کم باشد، در نظر بگیریم تا با این کار امکان اریب شدن برآوردهای مورد نظر طرح را کاهش دهیم. به عنوان مثال فردی که در یک منطقه شهری زندگی می‌کند و درامد نسبتاً بالایی دارد، ممکن است با احتمال کمتری در نمونه‌گیری شرکت کند. بنا بر این عدم حضور بخش زیادی از این قشر جامعه می‌تواند در برآوردهای میانگین درامد خانوارها، پیشرفت تحصیلی و... تأثیرگذار باشد.

۱- کاهش اریبی بی‌پاسخی در آمارگیری‌های خانوار

اریبی ناشی از بی‌پاسخی برای میانگین نمونه، تابعی از دو فاکتور زیر است:

- نسبتی از جامعه که پاسخ نداده‌اند؛

- اختلاف میانگین مشخصه مورد نظر در جامعه بین گروه‌های پاسخگو و بی‌پاسخ.
- کاهش اریبی ناشی از بی‌پاسخی، مستلزم پایین بودن نرخ بی‌پاسخی و اختلاف کم بین افراد و خانوارهای پاسخگو و بی‌پاسخ است. از طریق حفظ رکوردهای ثبت شده مربوط به هر واحد نمونه‌گیری، امکان برآورد نرخ بی‌پاسخی برای کل نمونه و یا زیر حوزه‌های مورد نظر وجود دارد. بنا بر این مطالعات خاصی در این زمینه می‌تواند به ارزیابی اختلافات بین واحدهای بی‌پاسخ و پاسخگو بپردازد. (گروز و کاپر ۱۹۹۸)

۴-۲- جبران بی‌پاسخی

روش‌های متعددی برای افزایش نرخ پاسخ‌دهی و کاهش اریبی ناشی از بی‌پاسخی در آمارگیری‌های خانوار وجود دارد. در یک مصاحبه ممکن است مصاحبه‌کننده مجبور باشد چند مرتبه تماس بگیرد و یا به افراد مراجعه کند. در واقع مصاحبه‌کننده باید با تلاش‌های متعدد، مصاحبه خود را با خانوار نمونه‌گیری شده به اتمام برساند. نرخ پاسخ‌دهی بالا از طریق آموزش بهتر مصاحبه‌گران امکان‌پذیر است. با وجود تلاش‌های زیاد در این زمینه، بی‌پاسخی امری غیر قابل اجتناب است. بنا بر این طراحان نمونه‌گیری اغلب تعديل‌هایی را برای جبران بی‌پاسخی ارائه می‌دهند. سه روش زیر روش‌های تعديل بی‌پاسخی واحد است.

- تعديل حجم نمونه با در نظر گرفتن حجم نمونه بالاتر از آن‌چه مورد نیاز است، صورت گیرد تا از این طریق بتوانیم بی‌پاسخی‌های مورد انتظار را نیز به حساب آوریم؛
- جانشین کردن؛ فرایند جایگذاری خانوارهای بی‌پاسخ با خانوارهای دیگری که در نمونه نبوده به‌طوری که خانوارها نسبت به مشخصه خاص، شبیه خانوارهای بی‌پاسخ باشند؛
- تعديل وزن نمونه به‌منظور جبران بی‌پاسخی.

توصیه می‌شود بی‌پاسخی واحد در آمارگیری خانوار از طریق تعديل وزن‌های نمونه که در بخش ۴-۳ با یک مثال عددی ارائه خواهد شد، صورت گیرد. مشکلات و مسائل

متعددی در ارتباط با جانشین کردن وجود دارد که می‌توان در کالتون (۱۹۸۳) با این مشکلات آشنا شد.

۴-۴- تعدیل بی‌پاسخی وزن‌های نمونه

روش تعدیل وزن‌های نمونه به منظور جبران بی‌پاسخی که در آمارگیری‌های خانوار مورد استفاده قرار می‌گیرد، در چهار مرحله زیر خلاصه شده است.

مرحله اول: به کارگیری وزن‌های اولیه (برای احتمال‌های انتخاب نابرابر و تعدیل‌های دیگر که در بخش‌های ۲ و ۳ بحث شده است، در صورت لزوم).

مرحله دوم: تقسیم‌بندی نمونه به زیرگروه‌ها و محاسبه نرخ پاسخ موزون برای هر زیر گروه.

مرحله سوم: استفاده از عکس نرخ پاسخ برای زیرگروه‌ها به منظور تعدیل بی‌پاسخی.

مرحله چهارم: محاسبه وزن تعدیل شده بی‌پاسخی برای ناممیں واحد به صورت زیر است:

$$w_i = w_{i1} \times w_{i2}$$

که در آن w وزن اولیه، w_1 وزن تعدیل بی‌پاسخی و در نتیجه w_2 وزن نهایی برای واحد نام پس از تعدیل بی‌پاسخی است.

مثال ۱: با استفاده از روش نمونه‌گیری چند مرحله‌ای طبقه‌بندی شده، نمونه‌ای به حجم ۱۰۰۰ خانوار از دو منطقه (شمال و جنوب) یک کشور انتخاب شده است. خانوارها در شمال کشور با نرخ $1/100$ و در جنوب با نرخ $1/200$ نمونه‌گیری شده‌اند. نرخ پاسخ در مناطق شهری کمتر از نرخ پاسخ در مناطق روستایی است.

هدف برآورد نسبت و تعداد خانوارهایی است که دسترسی به مراقبت‌های اولیه بهداشت دارند. فرض کنید n تعداد خانوارهای نمونه‌گیری شده در طبقه h ام، t تعداد خانوارهای واحد شرایط پاسخگو و r تعداد خانوارهای واحد شرایطی است که به اطلاعات پرسشنامه پاسخ داده‌اند و به مراقبت‌های اولیه بهداشت نیز دسترسی دارند. بنا بر این وزن تعدیل شده بی‌پاسخی برای خانوارها در طبقه h ام به صورت

$$W_h = W_{vh} \times W_{rh}$$

است که در آن $n_h/r_h = n_h/r_h$ ، وزن مربوط به تعدیل بی‌پاسخی در طبقه h ام می‌باشد (برای سادگی از اندیس z صرف نظر شده است). اطلاعات در سطح طبقات در جدول زیر خلاصه شده است.

جدول ۲ - اطلاعات خانوارهای نمونه‌گیری شده به تفکیک طبقات

$w_h t_h$	$w_h r_h$	w_h	w_{rh}	w_{vh}	t_h	r_h	n_h	طبقات
۸۷۵۰	۱۰۰۰	۱۲۵	۱/۲۵	۱۰۰	۷۰	۸۰	۱۰۰	شمال-شهری
۲۴۰۰۰	۳۰۰۰	۲۵۰	۲/۵۰	۱۰۰	۱۰۰	۱۲۰	۳۰۰	شمال-روستایی
۲۵۴۰۰	۴۰۱۲۰	۲۳۶	۱/۱۸	۲۰۰	۱۵۰	۱۷۰	۲۰۰	جنوب-شهری
۳۹۹۶۰	۷۹۹۲۰	۲۲۲	۱/۱۱	۲۰۰	۱۸۰	۲۶۰	۴۰۰	جنوب-روستایی
۱۰۹۱۱۰						۵۰۰	۷۳۰	کل

بنا بر این نسبت خانوارهایی که دسترسی به مراقبت‌های اولیه بهداشت دارند به صورت زیر برآورد می‌شود:

$$\hat{p} = \frac{\sum w_h t_h}{\sum w_h r_h} = \frac{109110}{160040} = .682 / .682 / 2$$

و تعداد خانوارهای برآورد شده که دسترسی به مراقبت‌های اولیه بهداشت دارند به صورت

$$\hat{t} = \sum w_h t_h = 109110$$

است. شایان ذکر است که برآورد غیر وزنی نسبت خانوارهایی که دسترسی به مراقبت‌های اولیه بهداشت دارند، تنها با استفاده از اطلاعات خانوارهایی که به پرسشنامه پاسخ داده‌اند، به صورت زیر برآورد می‌شود:

$$\hat{p}_{uw} = \frac{\sum t_h}{\sum r_h} = \frac{۵۰۰}{۷۳۰} = .۶۸۵ \quad , \quad \% ۶۸.۵$$

و برآورد نسبت خانوارهایی که دسترسی به مراقبت‌های اولیه بهداشت دارند با استفاده از وزن‌های اولیه و بدون در نظر گرفتن تعديل بی‌پاسخی به صورت

$$\hat{p}_1 = \frac{\sum w_{ih} t_h}{\sum w_{ih} r_r} = \frac{۸۳۰۰۰}{۱۲۶۰۰۰} = .۶۵۹ \quad , \quad \% ۶۵.۹$$

است. نتایج فوق نشان می‌دهد که اختلاف قابل قبولی بین نسبت برآورد شده با در نظر گرفتن وزن‌های اولیه در مقایسه با نسبت برآورد شده با در نظر گرفتن وزن‌های تعديل شده بی‌پاسخی وجود دارد. در صورتی که این اختلاف بین نسبت غیر وزنی و نسبت تعديل شده بی‌پاسخی جزئی است.

بعد از تصحیح وزن‌های بی‌پاسخی، در بخش بعدی، تعديل وزن‌ها به منظور جبران عدم پوشش بیان می‌شود.

۵- تعديل وزن برای عدم پوشش

عدم پوشش به عیب چارچوب نمونه‌گیری برای پوشش جامعه هدف و برخی از واحدهای نمونه که احتمال انتخاب در نمونه انتخابی را ندارند، اشاره می‌کند. این مشکل تنها یکی از عیوب متعدد چارچوب‌های نمونه‌گیری است که در انتخاب نمونه در آمارگیری‌های کشورهای در حال توسعه با آن مواجه هستیم. برای دستیابی به اطلاعات بیشتر در زمینه مشکلات موجود در چارچوب‌های نمونه‌گیری و برخی از راه حل‌های ممکن به یانسانه (۲۰۰۳) مراجعه شود. عدم پوشش یکی از نگرانی‌های عمده در آمارگیری‌های خانوار بالاخص در کشورهای در حال توسعه است. بنا بر این شناسایی، ارزیابی و کنترل عدم پوشش در آمارگیری‌های خانوار باید به عنوان یکی از مسائل اصلی در ادارات ملی آمار در کشورهای در حال توسعه مورد بررسی قرار گیرد.

در این بخش به شناسایی برخی از منابع عدم پوشش در بررسی‌های خانوار اشاره می‌شود و از یک روش جبران عدم پوشش، که به تعديل آماری وزن‌ها از طریق

پساطبقه‌بندی معروف است، اشاره می‌شود.

۱-۵- منابع عدم پوشش در آمارگیری‌های خانوار

اغلب آمارگیری‌های خانوار در کشورهای در حال توسعه بر اساس نمونه‌گیری چند مرحله‌ای طبقه‌بندی صورت می‌گیرد. واحدها در مرحله اول عموماً مناطق جغرافیایی هستند. در مرحله دوم فهرستی از خانوارها یا واحدهای مسکونی از واحدهای انتخاب شده، تهیه می‌شود و در آخرین مرحله نمونه‌ای از اعضای خانه یا ساکنین خانوار از نمونه انتخابی تهیه می‌شود. بنا بر این عدم پوشش ممکن است در هر سه سطح انتخاب واحدهای مرحله اول، خانوار و افراد وجود داشته باشد.

چون واحدهای مرحله اول نمونه‌گیری عمده‌ای بر اساس مناطق شمارش شده و اطلاعات سرشماری نفوس و مسکن قبلی است، انتظار می‌رود که کل حوزه جغرافیایی جامعه هدف پوشش داده شود. همچنین عدم پوشش در اولین مرحله انتخاب نمونه عموماً کم است. عدم پوشش در این مرحله در کشورهای در حال توسعه، به سختی عدم پوشش در مراحل بعدی نیست. اما به هر حال عدم پوشش در اولین مرحله انتخاب واحد نمونه‌گیری در بیشتر آمارگیری‌ها وجود دارد. به عنوان مثال، یک آمارگیری ممکن است برای براورده جمعیت وارد شده در یک کشور یا یک منطقه از کشور طراحی شده باشد، که بعضی از مناطق مربوط به اولین مرحله انتخاب نمونه ممکن است به جهت مسائلی چون جنگ و ناآرامی‌های درون کشوری و بلایای طبیعی یا دلایل دیگر غیر قابل دسترس باشند. همچنین مناطق پرت با خانوارهای کم جمعیت اغلب از چارچوب نمونه‌گیری در آمارگیری‌های خانوار حذف می‌شوند، زیرا برای پوشش آن‌ها باید هزینه زیادی پرداخت شود و این واحدها نماینده کوچکی از جامعه می‌باشند و تأثیر کمی در نتایج کلی جامعه دارند. بنا بر این در نتایج ارائه شده از چنین آمارگیری‌هایی، باید حذف از چنین مناطقی به وضوح توضیح داده شود. زمانی که بخشی از جامعه را پوشش نداده‌ایم، نباید چنین تلقی شود که نتایج آمارگیری برای کل کشور یا منطقه می‌باشد، بنا بر این مسئله عدم پوشش باید بهطور کامل در گزارش‌های آماری ثبت شده باشد.

عوامل متعدد دیگری در ایجاد عدم پوشش مؤثر است، که به تعدادی از آن‌ها اشاره

می‌کنیم. یکی از این عوامل تعریف ناصحیح از مفاهیم مربوط به هر آمارگیری است. عوامل دیگری شامل از قلم افتادگی غیر عمده واحدهای مسکونی از فهرست‌های تهیه شده در زمان عملیات میدانی، از قلم افتادگی بهجهت خطاهای اندازه‌گیری، عدم حضور اعضاي غایب خانوار و از قلم افتادگی بهجهت عدم درک صحیح از مفاهیم آمارگیری وجود دارد. جزئیات بیشتر در مورد منابع عدم پوشش در لپکوسکی (۲۰۰۳) و منابع موجود در آن بیان شده است. روش‌های زیر برای کنترل عدم پوشش در آمارگیری‌های خانوار وجود دارد:

- بهبود دستورالعمل‌های میدانی از طریق استفاده از چارچوب‌های چندگانه و بهبود دستورالعمل‌های فهرست‌برداری؛
- جبران عدم پوشش از طریق تعديل وزن‌ها.

در مثال زیر از روش تعديل وزن‌ها (روش دوم) استفاده شده است. روش تعديل وزن معروف به پساطبقة‌بندی است. در این روش، جبران عدم پوشش از طریق تعديل توزیع نمونه‌گیری وزنی برای متغیرهای خاص، انجام می‌گیرد بهطوری که این توزیع با توزیع جامعه مورد نظر مطابقت داشته باشد.

مثال ۲: در مثال قبل فرض کنید تعداد خانوارها در شمال کشور ۴۵۰۲۵ و در جنوب کشور ۱۱۵۸۰۰ است. کل نمونه وزنی بهترتبیب برای مناطق شمالی و جنوبی با اعمال وزن بی‌پاسخی برابر با ۴۰۰۰۰ و ۱۲۰۰۴۰ است. بنابراین برای محاسبه وزن‌ها به روش زیر عمل می‌کنیم:

مرحله اول: محاسبه فاکتورهای پساطبقة‌بندی

$$w_{\tau_h} = \frac{N_h}{\hat{N}_h} = \frac{45025}{40000} = 1/126 \quad \text{برای منطقه شمال}$$

$$w_{\tau_h} = \frac{N_h}{\hat{N}_h} = \frac{115800}{120040} = .0965 \quad \text{برای منطقه جنوب}$$

که در آن \hat{N}_h کل نمونه وزنی است.

مرحله دوم: محاسبه وزن‌های نهایی تعديل شده

$$w_f = w_h \times w_{rh}$$

که نتایج در جدول زیر خلاصه شده است:

جدول ۳- اعمال وزن پساطبقه‌بندی با استفاده از اطلاعات جدول قبلی

$w_f t_h$	$w_f r_h$	w_f	w_h	t_h	r_h	طبقات
۹۸۴۹	۱۱۲۵۶	۱۴۰/۷۵	۱۲۵	۷۰	۸۰	شمال- شهری
۲۸۱۴۰	۳۳۷۶۸	۲۸۱/۴۰	۲۵۰	۱۰۰	۱۲۰	شمال- روستایی
۳۴۱۵۵	۳۸۷۰۹	۲۲۷/۷۷	۲۳۶	۱۵۰	۱۷۰	جنوب- شهری
۳۸۵۵۶	۷۷۱۱۲	۲۱۴/۲۰	۲۲۲	۱۸۰	۲۶۰	جنوب- روستایی
۱۱۰۷۰۰	۱۶۰۸۴۵			۵۰۰	۷۳۰	کل

بنا بر این برآورد نسبت خانوارهایی که به مراقبت‌های اولیه بهداشت دسترسی دارند برابر است با:

$$\hat{p}_f = \frac{\sum w_f t_h}{\sum w_f r_h} = \frac{110700}{160845} = 0.688 \text{ یا } 68.8\%$$

همان‌طور که نتایج جدول نشان می‌دهد، با وزن‌های تعديل شده از روش پساطبقه‌بندی، تعداد نمونه وزنی برای مناطق شمالی و جنوبی به ترتیب $(11256 + 33768) = 45.24$ و $(38709 + 77112) = 115821$ است که به اعداد کنترلی که در بالا به آن اشاره شده است، بسیار نزدیک است.

۶- افزایش واریانس به دلیل وزن دهی

با وجود این که وزن دهی باعث کاهش اربیی برآوردها می‌شود، می‌تواند باعث افزایش واریانس برآوردها نیز شود. برای روشن شدن مطلب، طرح نمونه‌گیری یک مرحله‌ای

طبقه‌بندی شده، با احتمال برابر در هر طبقه را در نظر می‌گیریم. اگر واریانس طبقات (واریانس درون واحدهای طبقه) در هر طبقه یکسان نباشد، در نظر گرفتن وزن‌های نابرابر در طبقات، دقت برآوردهای طرح را بالا می‌برد و در صورتی که واریانس طبقات در هر طبقه یکسان باشد، در نظر گرفتن وزن‌های نابرابر منجر به افزایش واریانس برآوردها می‌شود. بنا بر این اثر وزن‌دهی، در افزایش واریانس جامعه برآورد شده به وسیله عامل زیر اندازه‌گیری می‌شود.

$$L = n \times \frac{\sum_h n_h w_h^r}{(\sum_h n_h w_h)^r}$$

که در آن n حجم نمونه قابل قبول یا واجد شرایط، w_h وزن نهایی و n_h حجم واقعی نمونه در طبقه h است. اکنون در مثال زیر به بررسی این مطلب می‌پردازیم.

مثال ۳: محاسبه فاکتور افزایش واریانس با استفاده از اطلاعات مثال قبل با وزن‌های نهایی w_h و n_h به عنوان حجم نمونه طبقه h ام.

جدول ۴- اطلاعات موردنیاز برای محاسبه فاکتور افزایش واریانس با استفاده از جدول قبلی

$w_f^r r_h$	$w_f r_h$	w_f	r_h	طبقات
۱۵۸۴۲۸۷	۱۱۲۶	۱۴۰/۷۵	۸۰	شمال- شهری
۹۵۰۲۳۱۵	۲۳۷۶۸	۲۸۱/۴۰	۱۲۰	شمال- روستایی
۸۸۱۴۰۳۹	۲۸۷۰۹	۲۲۷/۷۷	۱۷۰	جنوب- شهری
۱۶۵۱۷۳۹۰	۷۷۱۱۲	۲۱۴/۲۰	۳۶۰	جنوب- روستایی
۳۶۴۹۸۰۲۶	۱۶۰۸۴۹	۷۳۰		کل

$$L = 730 \times \frac{36418026}{(160849)^r} = 1/0.3$$

با توجه به فرمول فوق، فاکتور افزایش واریانس برابر است. به این معنا که به دلیل استفاده از روش وزن‌دهی به منظور جبران بی‌پاسخی و عدم پوشش، افزایش واریانس در برآوردهای آماری تقریباً ۳ درصد است.

۷- پیرایش وزن‌ها

زمانی که وزن‌دهی برای جبران عیوب و ویژگی‌های مربوط به طرح‌های آمارگیری محاسبه می‌شود، بررسی توزیع وزن‌های تعديل شده نیز لازم است. وزن‌های خیلی بزرگ، اگر تنها روی بخش کوچکی از اقلام واحدهای نمونه‌گیری تأثیر بگذارند، می‌توانند باعث افزایش قابل توجه‌ای در واریانس براوردهای آماری شوند. بنا بر این لازم است روی وزن‌های خیلی بزرگ پیرایش صورت گیرد، تا از این طریق بتوانیم واریانس براوردهای مربوطه را کنترل کنیم. پیرایش وزن‌ها، اکثرًا بعد از مرحله تعديل وزن بی‌پاسخی صورت می‌گیرد.

مادامی که پیرایش وزن‌ها منجر به کاهش واریانس براوردها شود، می‌تواند منجر به اریبی در براوردها نیز شود. در بعضی موارد، کاهش واریانس به دلیل پیرایش وزن‌های خیلی بزرگ، ممکن است باعث افزایش اریبی شود. بنا بر این پیرایش وزن‌ها، باید زمانی که ضرورت دارد، صورت گیرد. به عبارت دیگر زمانی باید از این روش استفاده کرد که اریبی به وجود آمده به‌واسطه پیرایش وزن‌ها از اهمیت کمتری نسبت به کاهش واریانس‌ها برخوردار باشد.

در هر طبقه‌بندی، فرایند پیرایش وزن‌ها باید در هر طبقه صورت گیرد. فرایند با تعیین یک کران بالا برای وزن‌های اصلی آغاز می‌شود و سپس تعديل کل وزن‌ها به‌گونه‌ای صورت می‌گیرد که جمع وزن‌های پیرایش شده مشابه جمع وزن‌های اصلی باشد. w_{hi} را وزن نهایی نامین واحد از طبقه h و w_{hb} را کران بالا برای وزن‌های مشخص شده در طبقه h در نظر می‌گیریم. وزن پیرایش شده برای نامین واحد نمونه‌گیری شده در طبقه h به صورت زیر تعریف می‌شود.

$$w_{hi(T)} = \begin{cases} w_{hi} & w_{hi} < w_{hb} \\ w_{hb} & w_{hi} \geq w_{hb} \end{cases}$$

وزن‌های پیرایش شده برای کل نمونه می‌تواند به این صورت که جمع آن‌ها دقیقاً مشابه جمع وزن‌های معمول است تعديل شود. این کار با محاسبه فاکتور زیر صورت می‌گیرد.

$$F_T = \frac{\sum_h n_h w_h}{\sum_h n_h w_{h(T)}}$$

جمع‌بندی در نسبت فوق روی همه طبقات و همه واحدها در نمونه است و وزن پیرایش تعديل شده به صورت زیر تعریف می‌شود.

$$w_{h(i)}^* = F_T \times W_{h(T)}$$

روشن است که، $\sum_h n_h w_{h(T)}^*$ است و این همان هدف مطلوب و مورد نظر

ما است. اکنون در مثالی به شرح روش پیرایش وزن‌ها می‌پردازیم. این مثال تنها برای فهم این مطلب بیان شده است و بر اساس یک وضعیت واقعی نیست.

مثال ۴: در جدول زیر دو ستون اول تعداد کل واحدهای نمونه‌گیری شده و وزن‌های نهایی در ۷ طبقه است. بیشترین وزن ۲۵۰ انتخاب شده است و همان‌طور که در ستون سوم جدول نشان داده شده است، وزن‌های اصلی در ۳۵۰ بریده شده است.

جدول ۵- محاسبه وزن مربوط به پیرایش وزن‌ها

$n_h w_{h(T)}^*$	$n_h w_{h(T)}$	$n_h w_h$	$W_{h(T)}$	w_h	n_h
۱۱۸۲۳/۰۰	۱۱۲۶-	۱۱۲۶-	۱۴۰/۷۵	۱۴۰/۷۵	۸۰
۱۵۷۷۶/۲۵	۱۵۰۲۵	۱۵۰۲۵	۱۵۰/۲۵	۱۵۰/۲۵	۱۰۰
۲۲۹۶۸/۷۵	۲۱۸۷۵	۲۱۸۷۵	۱۷۵/۰۰	۱۷۵/۰۰	۱۲۵
۳۱۵۰۰/۰۰	۳....	۳....	۲۰۰/۰۰	۲۰۰/۰۰	۱۵۰
۳۱۵۰۰/۰۰	۳....	۳....	۲۵۰/۰۰	۲۵۰/۰۰	۱۲۰
۳۱۵۰۰/۰۰	۳....	۳۳۱۵	۲۵۰/۰۰	۲۷۵/۱۳	۱۲۰
۴۴۶۲۵/۰۰	۴۲۵۰-	۴۸۵۱۸	۲۵۰/۰۰	۲۸۵/۴۰	۱۷۰
۱۸۹۶۹۳	۱۸۰۶۶۰	۱۸۹۶۹۳		۸۶۵	

فاکتور F_T با استفاده از اطلاعات جدول فوق به صورت زیر است:

$$F_T = \frac{\sum_i n_h w_{hi}}{\sum_i n_h w_{h(T)}} = \frac{189693}{18060} = 1/0.5$$

پس از این مرحله، وزن‌های پیراش شده از ضرب هر وزن در فاکتور $F_T = 1/0.5$ به دست می‌آید، به گونه‌ای که رابطه $\sum n_h w_h^* = 189693$ برقرار باشد.

۸- نتیجه‌گیری

وزن دهی به عنوان بخش جدایی‌ناپذیر از تحلیل‌های آماری آمارگیری‌های خانوار در کشورهای در جال توسعه از اهمیت ویژه‌ای برخوردار است. در اغلب طرح‌ها و دستورالعمل‌های آمارگیری از روش وزن دهی استفاده می‌شود. دیدگاه مورد قبول عامه این بوده است که استفاده از وزن‌ها، تحلیل‌ها را پیچیده‌تر می‌کند. اما پیشرفت‌های زیادی در زمینه تکنولوژی کامپیوتر در دهه گذشته، این موضوع را رد کرده است و وجود سخت‌افزارها و نرم‌افزارهای متعدد و برنامه‌های کامپیوترا مشخص برای تحلیل داده‌ها که در بسیاری از کشورهای در حال توسعه در دسترس می‌باشد، کار محاسبات را بسیار آسان کرده است.

همان‌طور که قبلًا بحث شد، وزن دهی باعث کاهش اریبی ناشی از مشکلات عدم پوشش و بی‌پاسخی می‌شود. عدم پوشش و بی‌پاسخی انواع مختلفی از خطاهای ناشی از نقص یک آمارگیری طراحی شده برای تهیه اطلاعات از برخی واحدهای جامعه هدف می‌باشد. برای آمارگیری‌های خانوار در کشورهای در حال توسعه، عدم پوشش مشکل جدی‌تری نسبت به بی‌پاسخی است که در این مقاله به روش‌های تعديل آماری وزن‌های پایه برای جبران مشکلات غیر قابل اجتناب اشاره شده است.

وجود نرم‌افزارهای آماری که در دسترس همه می‌باشد، مسأله وزن دهی را به عنوان یک کار معمول در تحلیل‌های آماری مربوط به آمارگیری‌های خانوار حتی در کشورهای در حال توسعه درآورده است. توجه دقیق به مسأله وزن دهی به جهت وجود خطاهای جدی

۱- مقدمه

از سال ۱۹۵۰ به بعد که رایانه، در تحلیل و ذخیره‌سازی داده‌ها به کار رفت، حجم اطلاعات ذخیره شده در آن پس از حدود ۲۰ سال دو برابر شد و همزمان با پیشرفت فناوری اطلاعات، حجم داده‌ها در پایگاه داده‌ها هر دو سال یک بار، دو برابر شد و همچنان با سرعت بیشتری نسبت به گذشته حجم اطلاعات ذخیره شده بیشتر و بیشتر می‌شود. با وجود شبکه جهانی وب، سیستم‌های یکپارچه اطلاعاتی، سیستم‌های یکپارچه بازکنی، تجارت الکترونیکی و... لحظه به لحظه به حجم داده‌ها در پایگاه داده‌ها اضافه شده و باعث به وجود آمدن انبارهای (توده‌های) عظیمی از داده‌ها شده است، بهطوری که ضرورت کشف و استخراج سریع و دقیق دانش از این پایگاه داده‌ها را بیش از پیش نمایان کرده است.

شدت رقابت‌ها در عرصه‌های علمی، اجتماعی، اقتصادی، سیاسی و نظامی نیز اهمیت سرعت یا زمان دسترسی به اطلاعات را دو چندان کرده است. بنا بر این نیاز به طراحی سیستم‌هایی که قادر به اکتشاف سریع اطلاعات مورد علاقه کاربران با تأکید بر حداقل مداخله انسانی باشند از یک سو و روی آوردن به روش‌های تحلیل مناسب با حجم داده‌های جیجیم از سوی دیگر، به خوبی احساس می‌شود. در حال حاضر، داده‌کاوی مهم‌ترین فناوری برای بهره‌برداری مؤثر، صحیح و سریع از داده‌های جیجیم است و اهمیت آن رو به فزونی است.

با توجه به وجود اطلاعات ارزشمند در پایگاه‌های داده‌ای در اوخر دهه ۸۰ میلادی، تلاش برای استخراج و استفاده از اطلاعات پایگاه‌های داده‌ای شروع شد. داده‌کاوی فرایندی است که در آغاز دهه ۹۰ پا به عرصه ظهرور گذاشت و با نگرشی نو، به مسئله استخراج اطلاعات از پایگاه داده‌ها می‌پردازد. در سال ۱۹۸۹ و ۱۹۹۱ کارگاه‌های کشف دانش از پایگاه داده‌ها توسط پیاتسکی و همکارانش و در فاصله سال‌های ۱۹۹۱ تا ۱۹۹۴ کارگاه‌های فوق، توسط فایاد و پیاتسکی و دیگران برگزار شد. به طور رسمی اصطلاح داده‌کاوی برای اولین بار توسط «فیاض» در اولین کنفرانس بین‌المللی «کشف دانش و داده‌کاوی» در سال ۱۹۹۵ مطرح شد. از سال ۱۹۹۵ داده‌کاوی به صورت جدی وارد مباحث آمار شد [۸] و در سال ۱۹۹۶، اولین شماره مجله کشف دانش از پایگاه داده‌ها