

برآورد مجموع در آمارگیری‌های دوچارچوبی

مرجان نورینی* و انور قیطولی

مرکز آمار ایران

چکیده: در آمارگیری‌های نمونه‌ای، چارچوبی که تقریباً تمام واحدهای جامعه‌ی هدف را پوشش دهد و نمونه‌گیری از آن کم‌هزینه باشد، کم‌تر یافت می‌شود. در حالی که ممکن است چارچوب‌های دیگری وجود داشته باشند که هر یک کسری از جامعه‌ی هدف را پوشش داده و در عین حال نمونه‌گیری از آن‌ها نیز کم‌هزینه باشد. در این مقاله نخست مفاهیم، مزایا و معایب چارچوب چندگانه تشریح شده و سپس روش‌های مختلف برای برآورد پارامترها مطرح شده است. در نهایت با استفاده از شبیه‌سازی مونت کارلویی، این روش‌ها با یکدیگر مقایسه شده‌اند.

واژگان کلیدی: آمارگیری نمونه‌ای؛ برآورد مجموع؛ چارچوب دوگان؛ شبیه‌سازی مونت کارلویی.

۱- مقدمه

یکی از مهم‌ترین مشکلات در اجرای آمارگیری‌های نمونه‌ای، چارچوب‌های فهرستی ناقص و یا تاریخ گذشته‌ای (Out of date) هستند که برای انتخاب نمونه مورد استفاده قرار می‌گیرند. در آمارگیری‌های نمونه‌ای که این نوع چارچوب‌ها مورد استفاده قرار می‌گیرند، برآورد پارامترهای جامعه دچار اریبی می‌شود. روزآمد کردن یک فهرست، کار پرهزینه و دشواری است که البته به یمن پیشرفت‌هایی که در مدیریت بانک‌های اطلاعاتی به وجود آمده، این کار تا حدودی امکان‌پذیر شده است. در هر حال، مهم‌ترین و پرهزینه‌ترین مرحله در روزآمدسازی یک فهرست، عملیات جمع‌آوری داده‌های مورد نیاز است که باید با استفاده از راهکارهای مناسب، هزینه را کاهش و در عین حال دقت را افزایش داد.

* نویسنده‌ی عهده‌دار مکاتبات

دریافت: ۱۳۹۱/۱۰/۲۶، پذیرش: ۱۳۹۲/۹/۴.

گاهی اوقات برای به دست آوردن فهرست واحدهای نمونه‌گیری (چارچوب)، از چند فهرست که در ارتباط با جامعه‌ی مورد بررسی است، استفاده می‌شود. در این حالت فرض می‌شود که این چارچوب‌های مختلف تقریباً کل جامعه را پوشش می‌دهند. مثلاً می‌توان با استفاده از منابع اداری، فهرست به دست آمده از یک سرشماری را روزآمد کرد. در این حالت چارچوب حاصل، بر اساس دو یا چند فهرست است. این روش را تنها وقتی می‌توان به کار برد که فهرست‌های مختلف، اطلاعات مورد نیاز برای تکمیل چارچوب را شامل شوند و اتصال داده‌ها و رکوردها به یکدیگر به خوبی انجام گیرد، به طوری که نتایج قابل اطمینانی را به دنبال داشته باشد. در غیر این صورت چارچوب حاصل ناقص بوده و واحدهای تکراری زیادی خواهد داشت. راه دیگر، استفاده از فهرست‌های مختلف به روش چارچوب چندگانه است. به عبارت دیگر، این فهرست‌ها ترکیبی از چارچوب‌های فهرستی و یا ترکیبی از چارچوب فهرستی به همراه چارچوب ناحیه‌ای هستند. در این روش، از ترکیب برآوردهای حاصل از واحدهای نمونه‌گیری متعلق به بخش غیر متداخل چارچوب‌ها با برآوردهای حاصل از واحدهای نمونه‌گیری متعلق به بخش متداخل چارچوب‌ها، یک برآوردها به دست می‌آید.

۲- تاریخچه

شاید بتوان گفت اولین شالوده‌ی آمارگیری‌های چندچارچوبی در سال ۱۹۴۹ با آمارگیری از فروشگاه‌های خرده‌فروشی که مجری آن دفتر سرشماری آمریکا بود، گذاشته شد. در این نمونه‌گیری که از دو چارچوب ناحیه‌ای و فهرستی استفاده شده بود، واحدهای آماری مشترک در دو چارچوب، شناسایی و از چارچوب ناحیه‌ای حذف شدند. بدین ترتیب با وجود پیچیدگی طرح، برآورد مجموع جامعه از ترکیب برآوردهای حاصل از دو چارچوب به دست آمد. از آنجایی که گاه شناسایی واحدهای آماری مشترک در دو یا چند چارچوب در حین عملیات میدانی امکان‌پذیر نیست، انجام نمونه‌گیری و تکمیل پرسشنامه برای واحدهای آماری مشترک، اجتناب‌ناپذیر است. به‌منظور جبران ورود بیش از یک بار واحدهای مشترک در نمونه هنگام برآورد، استفاده از تعدیل وزنی ضروری است. یکی از این تعدیل‌های وزنی که اغلب مورد استفاده قرار می‌گیرد، انتساب اوزانی به واحدهای مشترک در نمونه است که با امید ریاضی تعداد دفعات انتخابشان رابطه دارد.

هارتلی [۹] نظریه‌ی مقدماتی در مورد چارچوب چندگانه را توسعه داد و اولین کسی بود که انتساب وزن به واحدهای مشترک را مد نظر قرار داد. وی در بررسی‌های خود استفاده از تنها دو چارچوب با نمادهای A و B را در نظر گرفت. روشی که او ارایه نمود، بر اساس انتساب وزن θ به واحدهای مشترک متعلق به چارچوب A و وزن $(1 - \theta)$ به واحدهای مشترک متعلق به چارچوب B بود. بدین ترتیب می‌توان دو چارچوب با تعدادی واحد مشترک را به دو طبقه‌ی مجزا تبدیل نمود.

به علاوه، وی فرمول‌هایی را برای برآورد مجموع، مقدار بهینه‌ی θ و نیز کسرهای نمونه‌گیری بهینه ارایه کرد. بعد از او نیز مقاله‌های دیگری در این زمینه ارایه شد که از آن جمله می‌توان به [۲]، [۵]، [۶]، [۱۱]، [۱۳]، [۱۵]، [۱۶] و [۱۷] اشاره کرد. اخیراً استفاده از نمونه‌گیری‌های دوچارچوبی، که یکی از چارچوب‌ها را فهرست شماری تلفن افراد و دیگری را یک چارچوب ناحیه‌ای تشکیل می‌دهد، در برخی از کشورها مرسوم شده است. چرا که استفاده از چارچوبی شامل فهرست تلفن‌ها و انجام نمونه‌گیری به صورت متمرکز و تکمیل پرسشنامه از طریق مصاحبه‌ی تلفنی و نظارت هم‌زمان پرسشنامه توسط کارشناس در حین تکمیل پرسشنامه، هم خطای نمونه‌گیری را کاهش داده و هم هزینه‌ی نمونه‌گیری را به نحو چشمگیری تقلیل می‌دهد. از چنین طرح‌هایی می‌توان به طرح ملی بررسی آماری جنایات در آمریکا که در سال ۱۹۸۵ انجام شد، اشاره کرد. کسیدی و دیگران [۳] بررسی‌هایی را در مورد استفاده از چارچوب‌های دوگان انجام دادند که در آن‌ها از ترکیب چارچوب ناحیه‌ای که یک چارچوب کامل است و یک چارچوب تلفنی استفاده کردند. همچنین [۱۴] تأثیر نسبت‌های مختلف عدم پاسخگویی در نتایج را آزمون کرد. این بررسی‌ها توسط [۷] تعمیم و توسعه داده شدند. آن‌ها همچنین در مورد مزایای طرح‌های «چندچارچوبی - چند روشی» (Multiframe-Multimode)، با توجه به ساختار اجرایی (اداری) متفاوتی که برای استفاده از دو روش جمع‌آوری داده‌ها به طور هم‌زمان، مورد نیاز است، کاوش‌ها و بررسی‌هایی را انجام دادند. در مجموع بررسی‌های تلفنی دارای نرخ پاسخ پایین‌تری هستند. عقیده بر این است که بعضی از انواع پاسخ‌گویان (افراد مسن و کسانی که در گفتار و شنوایی دارای مشکلاتی هستند) در مصاحبه‌های تلفنی نسبت به مصاحبه‌های حضوری در جواب دادن به سؤالات، مشکلات بیشتری دارند. این مسئله میزان پرسش‌ها و ابهامات پیرامون کیفیت داده‌های جمع‌آوری شده تحت دو روش را افزایش می‌دهد.

در ایالت‌های متحده، بررسی‌های اولیه برای استفاده از دو چارچوب در طرح‌های نمونه‌گیری کشاورزی، از اواسط دهه‌ی ۱۹۵۰ مورد توجه قرار گرفت و نمونه‌گیری مبتنی بر این روش در بسیاری از نمونه‌گیری‌های کشاورزی در این کشور در دهه‌ی ۱۹۶۰ صورت گرفت. بدین ترتیب می‌توان مشاهده کرد که استفاده‌ی توأم از دو چارچوب در طرح‌های آماری کشاورزی از سال‌ها پیش در ایالت‌های متحده رایج شده است. لیکن به دلیل نیاز به فن‌آوری پیشرفته برای ایجاد یک چارچوب ناحیه‌ای، تاکنون تعداد اندکی از کشورها از این روش بهره جسته‌اند. سازمان خواربار و کشاورزی ملل متحد، FAO، به دلیل مزایایی که نمونه‌گیری‌های دوچارچوبی دارند، استفاده از این روش را به سایر کشورها توصیه کرده است. البته استفاده از طرح دوچارچوبی در طرح‌های آماری کشاورزی مستلزم وجود چارچوب ناحیه‌ای و فهرستی است. برای این منظور، تهیه‌ی یک چارچوب ناحیه‌ای کشاورزی مستلزم صرف هزینه‌ی بسیار، نیروی کارشناسی بسیار ماهر و فن‌آوری پیشرفته در زمینه‌ی سیستم‌های ماهواره‌ای و کامپیوتری است. لذا تهیه‌ی یک چارچوب ناحیه‌ای کشاورزی برای بار اول در هر کشور نیاز به صرف هزینه، زمان و دقت بسیار دارد.

۳- دلایل استفاده از چارچوب چندگانه

چارچوب آماری اساس و مبنای یک طرح آمارگیری نمونه‌ای را تشکیل می‌دهد. در صورت وجود مشکلاتی از قبیل عدم پوشش کامل، وجود واحدهای آمارگیری مشابه و پوشش فراتر از حد لازم در یک چارچوب آماری، خطای برآورد پارامترها افزایش پیدا می‌کند. در مرحله‌ی برآورد که لازم است نتایج را از نمونه به جامعه تعمیم دهیم در صورت وجود مشکل یا نقص در چارچوب، برآوردهای ارایه‌شده اریب خواهند شد، بنابراین باید به چارچوب آماری مورد استفاده توجه کافی معطوف شود.

از روش چارچوب چندگانه در اجرای آمارگیری‌هایی استفاده می‌شود که در آن‌ها چارچوب‌هایی با پوشش ناقص، وجود دارند. این چارچوب‌های ناقص تحت عنوان «چارچوب‌های مکمل» نیز شناخته می‌شوند. دستورالعمل چارچوب چندگانه این امکان را ایجاد می‌کند که چارچوب‌های ناقص با یکدیگر تداخل داشته باشند. هیچ تلاشی برای حذف عناصری از یک چارچوب ناقص که در عین حال به چارچوب ناقص دیگر تعلق دارد صورت نمی‌گیرد و نمونه‌ها از همه‌ی چارچوب‌ها انتخاب می‌شوند.

به‌طور کلی موارد استفاده بیش از یک چارچوب بنا به دلایلی است که در زیر به به آن‌ها اشاره می‌کنیم:

۱. گاهی ممکن است چارچوبی که همه‌ی واحدهای جامعه‌ی مورد مطالعه را پوشش می‌دهد در دسترس نباشد، اما امکان دستیابی به پوشش کامل با تلفیقی از دو یا چند چارچوب فراهم شود. در چنین حالتی به‌منظور دسترسی به پوشش مناسب، از دو چارچوب یا بیش‌تر به‌طور هم‌زمان استفاده می‌شود.
۲. گاهی نیز ممکن است یک چارچوب، پوشش کامل را برای جامعه‌ی مورد مطالعه فراهم کند اما چارچوب ناقص دیگری موجود باشد که هزینه‌ی آمارگیری از آن کم‌تر از هزینه‌ی آمارگیری از چارچوب کامل باشد. در این شرایط به‌دلیل پایین‌تر بودن هزینه‌ی آمارگیری از چارچوب ناقص، می‌توان با هزینه‌ای مشخص و ثابت از دو چارچوب استفاده و اندازه‌ی نمونه را بزرگ‌تر کرده و واریانس برآورد را کاهش داد.
۳. گاهی ممکن است صفت مورد بررسی مربوط به جامعه‌ی خاصی باشد که نسبت به کل جامعه در اقلیت است. به‌عنوان مثال می‌توان به جامعه‌ی ناشنویان، جامعه‌ی افرادی که مبتلا به یک نوع بیماری خاص هستند و یا دسترسی به میزان تولید محصولی خاص که سطح زیر کشت نسبتاً کمی را از کل سطح زیر کشت در محدوده‌ی مورد بررسی به خود اختصاص داده است، اشاره کرد. در اصطلاح به چنین جوامعی، جوامع کمیاب گفته می‌شود. بدیهی است که با استفاده از یک چارچوب ناحیه‌ای که پوشش کاملی را ارائه می‌دهد، احتمال دسترسی به این واحدهای خاص بسیار کم بوده، لذا می‌توان با بهره‌گیری از یک چارچوب دیگر که فهرستی ناقص از واحدهای آمارگیری جامعه‌ی مورد مطالعه است، آمارگیری را انجام داد. بسته به هدف، هزینه و دقت برآورد می‌توان چارچوب فهرستی را سرشماری یا نمونه‌گیری کرد. برای مثال در صورت بررسی مربوط به جامعه‌ی ناشنویان، می‌توان فهرستی از افراد ناشنوا را از مدارس ناشنویان، مؤسسه‌ی ارائه‌ی خدمات به ناشنویان و ... تهیه کرد. در این شرایط کسری از جامعه‌ی ناشنویان که از امکانات فوق استفاده نمی‌کنند در چارچوب فهرستی قرار نمی‌گیرند، لذا به‌منظور برطرف کردن این نقصان از چارچوب ناحیه‌ای که همه‌ی افراد جامعه را در برمی‌گیرد نیز استفاده می‌شود.
۴. دلیل دیگر لزوم استفاده از دو یا بیش از دو چارچوب، حتی در صورت دسترسی به یک چارچوب کامل، تغییرپذیری بسیار زیاد متغیر یا متغیرهای مورد مطالعه در

جامعه است که منجر به چولگی شدید می‌شود. مثلاً در طرح‌های کشاورزی ممکن است تعداد کمی از بهره‌برداران کشاورزی، مالک بهره‌برداری‌های بسیار بزرگ و تعداد بسیار زیادی از بهره‌برداران، مالک بهره‌برداری‌های کوچک باشند. لذا در صورت شناسایی و فهرست کردن بهره‌برداری‌های بزرگ، با تلفیق دو چارچوب ناحیه‌ای و فهرستی می‌توان از مزایای هر دو چارچوب استفاده کرده و کارایی طرح را بالا برد.

۵. گاه ممکن است یک چارچوب فهرستی کامل در دسترس باشد اما عملاً با گذشت زمانی نسبتاً طولانی به دلیل بروز تغییرات فراوان در آن، منبعی برای بروز خطاهای غیرنمونه‌گیری شود. از طرفی چون یک فهرست ناحیه‌ای، کم‌تر در معرض تغییرات است، با ترکیب این دو چارچوب می‌توان از مزایای هر دو چارچوب بهره برد.

در سال‌های اخیر به دلیل کاربرد روزافزون مصاحبه‌های تلفنی، مطلوبیت طرح‌های چارچوب چندگانه افزایش چشمگیری پیدا کرده است. همان‌طور که اشاره شد مصاحبه‌های تلفنی عموماً ارزان‌تر از مصاحبه‌های غیرتلفنی بوده و همچنین تمرکز پرسش‌گران در یک مکان، امکان نظارت بیش‌تری را فراهم می‌نماید. بنابراین با وجود این‌که چارچوب‌های تلفنی اغلب، کل جامعه‌ی آماری را پوشش نمی‌دهند، دارای مزایای زیادی هستند. به‌عنوان مثال، طرح ملی آمارگیری تندرستی که هر ساله توسط مرکز ملی آمارهای تندرستی ایالت‌های متحده از طریق تلفن اجرا می‌شود و در آن جامعه‌ی آماری شامل همه‌ی خانوارهای آمریکایی است. چارچوب مورد استفاده در این طرح چارچوب شماره‌ی تلفن‌های ثابت است. البته افرادی که در خانوارهای فاقد تلفن ثابت زندگی می‌کنند از شمول طرح خارج می‌شوند. در این طرح، پوشش کامل را می‌توان با استفاده از یک طرح با چارچوب دوگان، متشکل از یک چارچوب تلفنی و یک چارچوب ناحیه‌ای تأمین نمود.

۴- آمارگیری دو چارچوبی

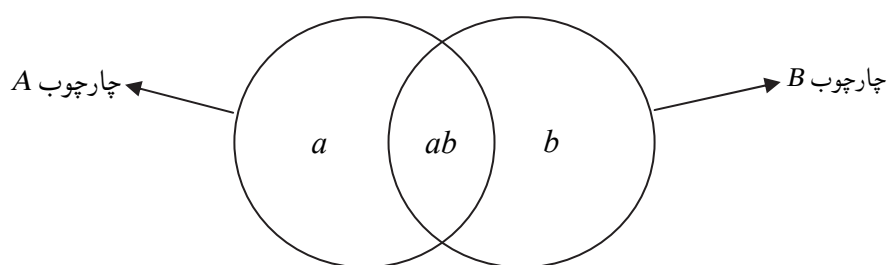
دو چارچوب A و B را در نظر گرفته و فرض می‌کنیم که از هر یک از آن‌ها نمونه‌ای استخراج شده است. می‌توان برای هر یک از چارچوب‌ها، طرح نمونه‌گیری متفاوتی استفاده نمود اما لازم است فرض‌های زیر برقرار باشند:

۱- کامل بودن، یعنی هر واحد در جامعه‌ی مورد مطالعه حداقل باید متعلق به یکی از چارچوب‌ها باشد.

۲- شناسایی‌پذیری، یعنی بتوان تشخیص داد که هر کدام از واحدهای انتخاب‌شده در نمونه به کدام یک از چارچوب‌ها تعلق دارند.

هنگامی که یکی از چارچوب‌ها، چارچوب ناحیه‌ای باشد، فرض کامل بودن معمولاً برقرار است. جوامعی مانند مزارع (در طرح‌های کشاورزی) به‌طور طبیعی با قطعات زمین که واحدهای چارچوب ناحیه‌ای را تشکیل می‌دهند، مرتبط هستند. با وجود ساده بودن فرض نظری شناسایی‌پذیری، اجرای یک طرح آمارگیری مبتنی بر چارچوب چندگانه دارای مشکلات زیادی می‌باشد. تعیین این‌که آیا یک واحد نمونه‌گیری‌شده در یک چارچوب به چارچوب دیگری نیز تعلق دارد، کار ساده‌ای نیست.

هارتلی [۹] نظریه‌ی بنیادی نمونه‌گیری مبتنی بر چارچوب چندگانه را توسعه داد و [۴] این نظریه را گسترش داد. هارتلی جامعه را با استفاده از چارچوب‌های نمونه‌گیری و تداخل بین آن‌ها به حوزه‌های به‌طور توأم ناسازگار تقسیم کرد. برای مثال، دو چارچوب نمونه‌گیری A و B ، سه حوزه‌ی ممکن را به‌صورت زیر تشکیل می‌دهند:



شکل ۱- مثالی از پوشش جامعه توسط دو چارچوب

در شکل ۱:

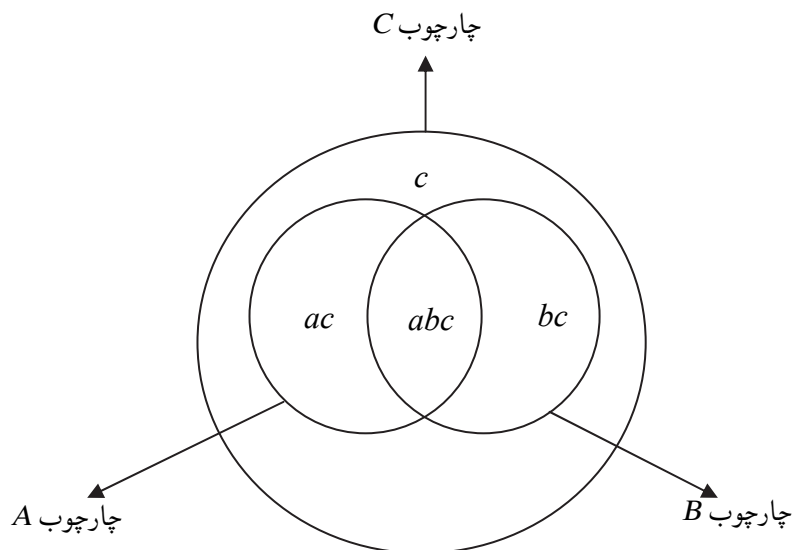
۱- حوزه‌ی a شامل واحدهایی است که فقط به چارچوب A تعلق دارند،

۲- حوزه‌ی b شامل واحدهایی است که فقط به چارچوب B تعلق دارند و

۳- حوزه‌ی ab شامل واحدهایی است که به هر دو چارچوب تعلق دارند.

فرض کنید می‌خواهیم یک بررسی در مورد آماردانان ایالت‌های متحده انجام دهیم. یک روش، استفاده از فهرست اعضای انجمن آمار آمریکا (ASA) به‌عنوان چارچوب نمونه‌گیری و گرفتن یک نمونه‌ی احتمالاتی از فهرست اعضا می‌باشد. اکثر افرادی که در

این بررسی آماری می‌شوند، آماردان هستند اما خیلی از آماردانان در ایالت‌های متحده عضو ASA نیستند پس اگر تنها از این فهرست به‌عنوان چارچوب استفاده شود با کم پوششی جامعه‌ی آماردانان مواجه خواهیم شد. برای بهبود پوشش چارچوب باید یک نمونه‌ی تکمیلی از چارچوب نمونه‌گیری دیگری نظیر فهرست اعضای انجمن آمار ریاضی (IMS) گرفته شود. بنابراین ما دو نمونه گرفته‌ایم؛ یک نمونه‌ی احتمالاتی از چارچوب A (فهرست ASA) و یک نمونه‌ی احتمالاتی از چارچوب B (فهرست اعضای IMS) که مستقل از نمونه‌ی چارچوب A است. دو چارچوب متداخل می‌باشند، چون خیلی از آماردانان عضو هر دو انجمن هستند. همان‌طور که در شکل ۱ نشان داده شده است ۳ حوزه‌ی a ، b و ab در جامعه وجود دارند. چارچوب A و B هر دو با هم پوشش بهتری برای جامعه‌ی آماردانان ایالت‌های متحده ایجاد می‌کنند اما هنوز شامل همه‌ی آماردانان نیستند. در این حالت، یک چارچوب سومی (C) می‌تواند استفاده شود که شامل جمعیت افراد بزرگسال است. ساختار جامعه با ۳ چارچوب در شکل ۲ نشان داده شده است.



شکل ۲- پوشش جامعه‌ی آماردانان آمریکا توسط سه چارچوب A ، B و C

چارچوب C کل جامعه را پوشش می‌دهد و دو چارچوب متداخل A و B را نیز در بر دارد. در این طرح، حوزه‌ی C بخشی از جامعه در چارچوب کامل C است که در دو فهرست A و B نیستند.

چون چارچوب C کل افراد جامعه را در بر دارد امکان طرح این سؤال وجود دارد که با وجود کامل بودن پوشش چارچوب C ، چرا تنها از آن به‌عنوان چارچوب نمونه‌گیری استفاده نمی‌شود. لازم به ذکر است که چارچوب C شامل کل جامعه‌ی بزرگسالان است ولی نمونه‌گیری از آن پرهزینه است. چون اکثر بزرگسالان، آماردان نیستند و بسیاری از نمونه‌های حاصل از چارچوب C آماردان نیستند. نمونه‌گیری از چارچوب‌های A و B خیلی ارزان‌تر است اما آن‌ها کل جامعه‌ی مورد مطالعه را پوشش نمی‌دهند. با ترکیب نمونه از هر سه چارچوب می‌توان از کم هزینه بودن دو چارچوب A و B و کامل بودن چارچوب C سود برد. دلایل مختلفی برای استفاده از چارچوب‌های چندگانه در جمع‌آوری داده‌ها وجود دارد. چارچوب‌های چندگانه به‌خصوص در مواردی همچون نمونه‌گیری از جوامع کمیاب به‌طور چشمگیری کارایی جمع‌آوری داده‌ها را بهبود می‌بخشد. یک جامعه‌ی کمیاب یک زیرگروه از جامعه‌ی مورد مطالعه است که تنها بخش کوچکی از جامعه (معمولاً ۱۰ درصد یا کم‌تر) را در بر دارد. در مثال بالا، آماردانان بخش کوچکی از جامعه‌ی بزرگسالان ایالت‌های متحده هستند. اگر نمونه‌های حاصل از فهرست‌های اعضا را به‌عنوان مکمل نمونه از کل جامعه در نظر بگیریم، تعداد آماردانان در مجموعه‌ی داده‌ها افزایش می‌یابد. فرض کنید می‌خواهیم افراد با یک بیماری خاص را مورد بررسی قرار دهیم. در این حالت می‌توان با نمونه‌گیری از چارچوب‌های اضافی که شامل نسبت بالایی از این افراد هستند نمونه‌ی بزرگ‌تری بدست آورد. البته برای نمونه‌گیری از جوامع کمیاب روش‌های دیگری نیز وجود دارد که از این بحث خارج است.

در صورت استفاده از آمارگیری‌های چندچارچوبی می‌توان از روش‌های مختلفی برای جمع‌آوری داده‌ها استفاده کرد. به‌عنوان مثال در شکل ۱، چارچوب A ممکن است شامل شماره‌های تلفن ثابت و چارچوب B شامل شماره‌های موبایل باشد.

هارتلی [۹] هنگام معرفی نظریه‌ی چارچوب چندگانه، دلیل اصلی بررسی آمارگیری چندچارچوبی را کاهش هزینه‌های جمع‌آوری داده‌ها معرفی کرد. «در یک آمارگیری ممکن است یک چارچوب با پوشش کامل داشته باشیم اما آمارگیری از آن پرهزینه است در

حالی که چارچوب‌های دیگری موجود هستند که نمونه‌گیری از آن‌ها ارزان‌تر است و بخشی از جامعه‌ی هدف را پوشش می‌دهند».

حالت دیگری که اغلب در نمونه‌گیری چندچارچوبی مورد استفاده قرار می‌گیرد و توسط هارتلی بررسی شده در شکل ۳ نشان داده شده است. در این حالت، چارچوب A یک چارچوب ناحیه‌ای است که پوشش کاملی از جامعه دارد و چارچوب B اغلب یک چارچوب فهرستی است.

چارچوب ناحیه‌ای (مجموعه‌ی نواحی جغرافیایی)، چارچوبی شامل واحدهای ناحیه‌ای است و هر عنصر جامعه متعلق به تنها یکی از ناحیه‌ها است و آن عنصر بعد از ارتباط با آن واحد ناحیه‌ای می‌تواند مشخص شود. به‌علت کامل بودن چارچوب‌های ناحیه‌ای، می‌توان از آن‌ها در موارد بسیاری استفاده کرد. برای مثال زمانی که:

- ۱- یک چارچوب کامل دیگر در دسترس نباشد.
- ۲- واحدهای نمونه‌گیری در چارچوب به‌سرعت تغییر کنند.
- ۳- چارچوب موجود، کهنه شده باشد.
- ۴- چارچوب موجود، از یک سرشماری با پوشش ضعیف بدست آمده باشد.
- ۵- یک چارچوب چندمنظوره برای برآورد صفت‌های مختلف (کشاورزی، زیست محیطی و غیره) مورد نیاز باشد.

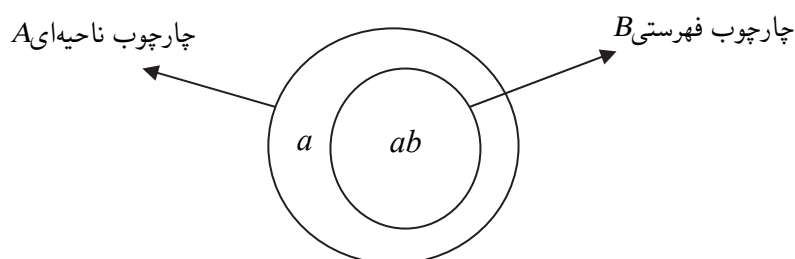
در طرح‌هایی که مبتنی بر چارچوب ناحیه‌ای است می‌توان بدون انجام مصاحبه، در ارتباط با خصوصیات قابل مشاهده بر روی زمین، برآوردهایی به‌دست آورد. به‌همین علت، خطاهای غیرنمونه‌گیری از جمله خطای بی‌پاسخی در مصاحبه‌ها کاهش می‌یابد. البته طرح‌های نمونه‌گیری مبتنی بر چارچوب ناحیه‌ای دارای معایبی به شرح ذیل هستند:

- ۱- پرهزینه بودن انجام و اجرای نمونه‌گیری؛
- ۲- ضرورت استفاده از تجهیزات کارتوگرافیکی؛
- ۳- حساسیت به داده‌های دورافتاده؛
- ۴- بی‌ثباتی و ناستواری برآوردها.

متداول‌ترین راه برای اجتناب از بی‌ثباتی برآوردها و بهبود دقت آن‌ها استفاده‌ی هم‌زمان از چارچوب‌های ناحیه‌ای و فهرستی است. مسأله‌ی اصلی در این روش تعیین واحدهای نمونه‌گیری ناحیه‌ای است که در فهرست نیز موجود می‌باشند. وقتی نمونه‌های انتخاب‌شده از چارچوب ناحیه‌ای را نتوان در فهرست تعیین کرد، برآوردگر مجموع جامعه

دچار آریبی بالا می‌شود. نقش و وظیفه‌ی چارچوب ناحیه‌ای در چارچوب‌های چندگانه حل مشکلات مربوط به کامل نبودن فهرست و همچنین برآورد آن می‌باشد. به‌عنوان مثالی از این نوع چارچوب‌های چندگانه در آمارگیری‌های کشاورزی، فرض کنید چارچوب B فهرست زمین‌های کشاورزی بزرگ در یک آبادی است. این فهرست ممکن است کهنه شده باشد و یا شامل زمین‌هایی که اخیراً به‌وجود آمده‌اند، نباشد. علاوه بر آن، این فهرست زمین‌های کشاورزی کوچک را هم در بر ندارد. هر چند که روش‌های کارایی برای جمع‌آوری داده می‌تواند در چارچوب B استفاده شود. اگر فقط از چارچوب B استفاده کنیم پوشش جامعه ناقص است، اما اگر از یک چارچوب ناحیه‌ای مثل چارچوب A استفاده کنیم پوشش کاملی از زمین‌های کشاورزی به وجود خواهد آمد. در یک چارچوب ناحیه‌ای، کشور به نواحی جغرافیایی تقسیم می‌شود و بخش‌هایی در این نواحی نمونه‌گیری می‌شوند. آمارگیران به هر بخش انتخاب‌شده رفته و از هر فعالیت کشاورزی درون این محدوده، آمارگیری می‌کنند. چارچوب ناحیه‌ای شامل همه‌ی زمین‌های کشاورزی کشور است اما جمع‌آوری داده از آن پرهزینه‌تر از جمع‌آوری داده از یک چارچوب فهرستی است.

آمارگیری از هزینه‌های مشتری در آمریکا (The U.S. survey of consumer finances) مثالی از یک آمارگیری دو چارچوبی مطابق با شکل ۳ است که در آن یک نمونه (نمونه از چارچوب A) از جامعه‌ی آمریکا با استفاده از چارچوب ناحیه‌ای با نمونه‌گیری طبقه‌بندی چندمرحله‌ای انتخاب می‌شود. ولی این نمونه ممکن است شامل همه‌ی افراد مورد مطالعه نباشد بنابراین یک نمونه دومی از چارچوب فهرستی (چارچوب B) گرفته می‌شود (چارچوب B شامل رکوردهای مالیاتی است).



شکل ۳- مثالی از دو چارچوب ناحیه‌ای و فهرستی

۵- برآوردهای مجموع در طرح‌های چارچوب دوگانه

در این بخش، روش‌های برآورد برای طرح‌های چارچوب دوگان، هنگامی که نمونه‌های حاصل از دو چارچوب به‌طور مستقل استخراج شوند، مورد بررسی قرار می‌گیرند. فرض کنید که مجموعه‌ی U دارای N واحد باشد. مجموع جامعه برای ویژگی y ، به‌صورت $Y = \sum_{i=1}^N y_i$ است. اگر در یک طرح نمونه‌گیری احتمالاتی برای یک نمونه‌ی S که از یک چارچوب استفاده می‌شود، $\pi_i = p (i \in S)$ احتمال شمول واحد i در نمونه‌ی S باشد، برآوردگر هورویتز-تامپسون مجموع جامعه برابر است با:

$$\hat{Y}_{HT} = \sum_{i \in S} w_i y_i$$

که در آن $w_i = \frac{1}{\pi_i}$ ، وزن نمونه‌گیری است.

اکنون حالتی را در نظر می‌گیریم که دو نمونه‌ی مستقل S_A و S_B از دو چارچوب متداخل A و B همانند شکل ۱ استخراج شود. چارچوب A شامل N_A واحد و چارچوب B شامل N_B واحد می‌باشند. واحدهای جامعه در حوزه‌ی متداخل ab می‌تواند از واحدهای یک چارچوب یا هر دو چارچوب انتخاب شوند. فرض می‌کنیم

$$\delta_i(a) = \begin{cases} 1 & \text{اگر واحد } i \text{ متعلق به حوزه‌ی } a \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

$$\delta_i(b) = \begin{cases} 1 & \text{اگر واحد } i \text{ متعلق به حوزه‌ی } b \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

و

$$\delta_i(ab) = \begin{cases} 1 & \text{اگر واحد } i \text{ متعلق به حوزه‌ی } ab \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

در این صورت می‌توان مجموع جامعه را به‌عنوان مجموع سه حوزه به‌صورت زیر نوشت:

$$Y = Y_a + Y_b + Y_{ab} = \sum_{i=1}^N \delta_i(a) y_i + \sum_{i=1}^N \delta_i(b) y_i + \sum_{i=1}^N \delta_i(ab) y_i$$

برآوردگر مجموع جامعه \hat{Y} ، می‌تواند با مجموع برآوردگرهای \hat{Y}_a ، \hat{Y}_b و \hat{Y}_{ab} بدست آید. فرض کنید احتمال شمول واحد i ام در نمونه‌ی S_A برابر $\pi_i^A = P\{i \in S_A\}$ بوده و S_A شامل n_A واحد نمونه باشد. به‌طور متناظر برای چارچوب B احتمال شمول واحد j ام در نمونه‌ی S_B برابر $\pi_j^B = P\{j \in S_B\}$ و شامل n_B واحد نمونه می‌باشد. فرض می‌کنیم w_i^A وزن نمونه‌گیری از چارچوب A و w_i^B وزن نمونه‌گیری از چارچوب B باشند. این وزن‌ها هم می‌توانند عکس احتمال‌های π_i^A و π_i^B و هم می‌توانند وزن‌های هایک [۸] با استفاده از این وزن‌ها، برآوردگرهای حوزه‌های a و ab از چارچوب A را به‌صورت زیر تعریف می‌کنیم:

$$\hat{Y}_a^A = \sum_{i \in S_A} w_i^A \delta_i(a) y_i$$

$$\hat{Y}_{ab}^A = \sum_{i \in S_A} w_i^A \delta_i(ab) y_i$$

به‌طور متناظر، برآوردگرهای حوزه‌های b و ab از نمونه‌ی چارچوب B به‌صورت زیر هستند:

$$\hat{Y}_b^B = \sum_{i \in S_B} w_i^B \delta_i(b) y_i$$

$$\hat{Y}_{ab}^B = \sum_{i \in S_B} w_i^B \delta_i(ab) y_i$$

تحت نتایج نظریه‌ی نمونه‌گیری استاندارد، هر کدام از برآوردگرهای مجموع در حوزه‌های a ، b و ab برای پارامتر میانگین جامعه تقریباً نارایب هستند. اندازه‌ی جامعه در هر یک از حوزه‌های a ، b و ab با مقداری $y_i = 1$ برای همه‌ی واحدها برآورد می‌شود و این برآوردگرها را با \hat{N}_a^A ، \hat{N}_b^B ، \hat{N}_{ab}^A و \hat{N}_{ab}^B نشان می‌دهیم. در یک طرح چارچوب دوگان که واحدهای مشترک از یک چارچوب غربال می‌شوند، حوزه‌ی ab تهی است و می‌توان چارچوب‌های A و B را به‌عنوان دو طبقه‌ی مجزا فرض کرد. بنابراین برآورد مجموع جامعه در یک طرح چارچوب دوگان غربالی به‌صورت زیر است:

$$\hat{Y} = \hat{Y}_a^A + \hat{Y}_b^B$$

در ادامه، روش‌هایی که برای برآورد مجموع جامعه، Y ، در حالت چارچوب دوگانه بیان شده، به صورت مختصر تشریح می‌شود.

۱-۵- برآوردگر هارتلی

هارتلی اولین کسی بود که روش‌شناسی سیستماتیک را برای تحلیل طرح‌های آمارگیری چندچارچوبی پایه‌گذاری کرد. در آمارگیری‌های چارچوب دوگان، \hat{Y}_{ab}^A و \hat{Y}_{ab}^B ، هر دو Y_{ab} را برآورد می‌کنند. بنابراین به منظور اجتناب از چندبارگی (Multiplicity) برآورد، هارتلی در [۹] یک برآوردگر موزون به صورت زیر پیشنهاد داد:

$$(۱) \quad \hat{Y}(\theta) = \hat{Y}_a^A + \theta \hat{Y}_{ab}^A + (1 - \theta) \hat{Y}_{ab}^B + \hat{Y}_b^B,$$

که در آن θ پارامتر آمیختگی است؛ $0 \leq \theta \leq 1$. این برآوردگر، وزن هر واحد نمونه‌گیری شده در حوزه‌ی مشترک ab را به منظور جبران چندبارگی کاهش می‌دهد. وزن‌های جدید را به صورت زیر تعریف می‌کنیم:

$$\begin{aligned} \tilde{w}_i^A &= \delta_i(a) w_i^A + \theta \delta_i(ab) w_i^A \\ \tilde{w}_i^B &= \delta_i(b) w_i^B + (1 - \theta) \delta_i(ab) w_i^B \end{aligned}$$

بنابراین

$$\hat{Y}(\theta) = \sum_{i \in S_A} \tilde{w}_i^A y_i + \sum_{i \in S_B} \tilde{w}_i^B y_i.$$

هر کدام از برآوردگرهای مجموع در هر حوزه تقریباً نااریب است، بنابراین $\hat{Y}(\theta)$ یک برآوردگر تقریباً نااریب از مجموع جامعه‌ی Y است. چون چارچوب‌های A و B به طور مستقل نمونه‌گیری می‌شوند و θ ثابت است، واریانس برآوردگر برابر است با:

$$(۲) \quad \text{Var} [\hat{Y}(\theta)] = \text{Var} [\hat{Y}_a^A + \theta \hat{Y}_{ab}^A] + \text{Var} [(1 - \theta) \hat{Y}_{ab}^B + \hat{Y}_b^B].$$

برآوردگر واریانس رابطه‌ی (۱) برای محاسبه آسان است اما ممکن است نسبت به برآوردگرهای دیگر کارایی کمتری داشته باشد. θ در رابطه‌ی (۱) به گونه‌ای انتخاب می‌شود که واریانس $\hat{Y}(\theta)$ را مینیمم کند ([۹] و [۱۰]). به دلیل این که نمونه‌گیری از چارچوب‌ها

به‌طور مستقل انجام می‌شود، واریانس $\hat{Y}(\theta)$ برابر است با مجموع واریانس دو چارچوب که در رابطه‌ی (۲) نشان داده شده است. مقدار بهینه‌ی θ که واریانس فوق را مینیمم کند برابر است با:

$$(۳) \quad \theta_{opt} = \frac{V(\hat{Y}_{ab}^B) + \text{Cov}(\hat{Y}_b^B, \hat{Y}_{ab}^B) - \text{Cov}(\hat{Y}_a^A, \hat{Y}_{ab}^A)}{V(\hat{Y}_{ab}^A) + V(\hat{Y}_{ab}^B)}$$

با θ_{opt} ، واریانس مینیمم شده به‌صورت زیر حاصل می‌شود:

$$\text{Var}[\hat{Y}(\theta)] = \text{Var}(\hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}^B) - \theta_{opt}^2 [\text{Var}(\hat{Y}_{ab}^A) + \text{Var}(\hat{Y}_{ab}^B)]$$

با توجه به رابطه‌ی (۳)، بزرگ‌تر بودن $\text{Var}(\hat{Y}_{ab}^B)$ نسبت به $\text{Var}(\hat{Y}_{ab}^A)$ منجر به بزرگ‌تر شدن θ_{opt} می‌شود. نکته آن‌که اگر مقدار قدر مطلق $\text{Cov}(\hat{Y}_b^B, \hat{Y}_{ab}^B)$ یا $\text{Cov}(\hat{Y}_a^A, \hat{Y}_{ab}^A)$ بزرگ باشد در این صورت، این امکان وجود دارد که θ_{opt} کوچک‌تر از صفر یا بزرگ‌تر از یک شود. هنگامی که دو چارچوب A و B یکسان باشند (حوزه‌های a و b تهی هستند)، θ_{opt} بین صفر و یک است.

در عمل، واریانس‌ها و کوواریانس‌ها در رابطه‌ی (۳) نامعلوم بوده و باید برآورد شوند. فرض کنیم $\hat{\theta}_{opt}$ برآوردگر θ_{opt} باشد. با جایگذاری آن در وزن‌های تعدیل شده داریم:

$$\begin{aligned} \tilde{w}_{i,H}^A &= \delta_i(a) w_i^A + \hat{\theta}_{opt} \delta_i(ab) w_i^A \\ \tilde{w}_{i,H}^B &= \delta_i(b) w_i^B + (1 - \hat{\theta}_{opt}) \delta_i(ab) w_i^B \end{aligned}$$

۲-۵- برآوردگر فولر-بورمیستر (Fuller-Burmeister)

فولر و بورمیستر برآوردگر هارتلی را با افزودن اطلاعات بیشتر راجع به برآورد N_{ab} تعدیل و پیشنهاد کردند [۶]. این برآوردگر به‌صورت زیر است:

$$(۴) \quad \hat{Y}_{FB}(\beta) = \hat{Y}_a^A + \hat{Y}_b^B + \beta_1 \hat{Y}_{ab}^A + (1 - \beta_1) \hat{Y}_{ab}^B + \beta_2 (\hat{N}_{ab}^A - \hat{N}_{ab}^B).$$

همانند برآوردگر هارتلی، پارامترهای β_1 و β_2 به‌گونه‌ای انتخاب می‌شوند که واریانس $\hat{Y}_{FB}(\beta)$ را مینیمم کنند. مقادیر بهینه‌ی β_1 و β_2 عبارت‌اند از:

$$\begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} \text{Var}(\hat{Y}_{ab}^A - \hat{Y}_{ab}^B) & \text{Cov}(\hat{Y}_{ab}^A - \hat{Y}_{ab}^B, \hat{N}_{ab}^A - \hat{N}_{ab}^B) \\ \text{Cov}(\hat{Y}_{ab}^A - \hat{Y}_{ab}^B, \hat{N}_{ab}^A - \hat{N}_{ab}^B) & \text{Var}(\hat{N}_{ab}^A - \hat{N}_{ab}^B) \end{bmatrix}^{-1}$$

$$\times \begin{bmatrix} \text{Cov}(\hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}^B, \hat{Y}_{ab}^A - \hat{Y}_{ab}^B) \\ \text{Cov}(\hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}^B, \hat{N}_{ab}^A - \hat{N}_{ab}^B) \end{bmatrix}$$

واریانس‌ها و کوواریانس‌ها مانند برآوردگر هارتلی بهینه، باید از روی داده‌ها برآورد شوند. این نتایج در یک مجموعه‌ی متفاوت از وزن‌ها برای هر متغیر پاسخ استفاده می‌شود.

۳-۵- برآوردگر تک‌چارچوبی (Single-frame Estimator)

بانکیر [۲] و [۱۱] روش‌های تک‌چارچوبی را پیشنهاد دادند که مشاهدات را در یک مجموعه داده‌ها ترکیب می‌کند و سپس وزن‌ها در حوزه‌ی مشترک تعدیل می‌شوند. یک واحد مشاهده شده‌ی i در حوزه‌ی ab می‌تواند هم در S_A و هم در S_B انتخاب شود، بنابراین تعداد دفعات مورد انتظار انتخاب واحد i در حوزه‌ی ab برابر $\pi_i^A + \pi_i^B$ است. بنابراین برآوردگر تک‌چارچوبی [۱۱] از $1/(\pi_i^A + \pi_i^B)$ به‌عنوان وزن واحدهای مشاهده‌شده در حوزه‌ی ab استفاده می‌کند. بنابراین، اگر $w_i^A = 1/\pi_i^A$ و $w_i^B = 1/\pi_i^B$ ، در این صورت وزن تعدیل شده برای واحدهای نمونه‌گیری شده در چارچوب A به‌صورت زیر است:

$$\tilde{w}_{i,S}^A = \begin{cases} w_i^A & i \in a \\ (1/w_i^A + 1/w_i^B)^{-1} & i \in ab \end{cases}$$

وزن‌های تعدیل شده برای واحدهای نمونه‌گیری شده در چارچوب B به‌طور مشابه به صورت زیر تعریف می‌شود:

$$\tilde{w}_{i,S}^B = \begin{cases} w_i^B & i \in b \\ (1/w_i^A + 1/w_i^B)^{-1} & i \in ab \end{cases}$$

بنابراین برآوردگر تک‌چارچوبی به‌صورت زیر است:

$$(5) \quad \hat{Y}_S = \sum_{i \in S_A} \tilde{w}_{i,S}^A y_i + \sum_{i \in S_B} \tilde{w}_{i,S}^B y_i$$

هنگامی که هر نمونه خودوزن است، همه‌ی وزن‌ها برای واحدهای نمونه‌گیری شده در S_A برابر w^A و همه‌ی وزن‌ها برای واحدهای نمونه‌گیری شده در S_B برابر w^B است، بنابراین

\hat{Y}_S یک حالت خاص از برآوردگر رابطه‌ی (۱) بوده و می‌توان آن را به صورت زیر بازنویسی کرد:

$$\hat{Y}_S = \hat{Y}_a^A + \theta_s \hat{Y}_{ab}^A + (1 - \theta_s) \hat{Y}_{ab}^B + \hat{Y}_b^B$$

به طوری که

$$\theta_s = \frac{\pi_i^A}{(\pi_i^A + \pi_i^B)} = [w^A (\frac{1}{w^A} + \frac{1}{w^B})]^{-1}$$

برخلاف برآوردگرهای هارتلی و فولر-بورمیستر، برآوردگرهای تک چارچوبی به برآورد کوواریانس‌های متغیر پاسخ بستگی ندارد. این برآوردگر از مجموعه‌ی وزن‌های یکسانی برای هر پاسخ مورد نظر، استفاده می‌کند. محاسبه‌ی وزن‌ها در حوزه‌ی ab ، نیازمند اطلاعات راجع به احتمال‌های شمول در دو چارچوب (نه فقط چارچوبی که یک واحد از آن انتخاب شده است) می‌باشد. اگر یکی از نمونه‌ها (برای مثال نمونه‌ی چارچوب A) خودوزن نباشد در این صورت π_i^A ممکن است برای یک واحد انتخاب‌شده در S_B نامعلوم باشد.

۴-۵- برآوردگر ماکسیمم درست‌نمایی نما

هنگامی که N_B و N_A معلوم باشند، [۱۶] یک برآوردگر دیگر که تعدیل برآوردگر فولر-بورمیستر است را برای نمونه‌های تصادفی ساده به منظور دستیابی به برآوردگر ماکسیمم درست‌نمایی نما (PML) که می‌تواند در طرح‌های پیچیده استفاده شود، پیشنهاد دادند. برآوردگر PML از مجموعه وزن‌های یکسانی برای همه‌ی متغیرهای پاسخ استفاده می‌کند و به صورت زیر است:

$$\hat{Y}_{PML} = \frac{N_A - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_a^A} \hat{Y}_a^A + \frac{N_B - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_b^B} \hat{Y}_b^B + \frac{\hat{N}_{ab}^{PML}(\theta)}{\theta \hat{N}_{ab}^A + (1-\theta) \hat{N}_{ab}^B} \left[\theta \hat{Y}_{ab}^A + (1-\theta) \hat{Y}_{ab}^B \right], \quad (۶)$$

که در آن $\hat{N}_{ab}^{PML}(\theta)$ کوچک‌ترین ریشه‌ی معادله‌ی درجه‌ی دو زیر است:

$$\left[\frac{\theta}{N_B} + \frac{1-\theta}{N_A} \right] x^2 - \left[1 + \frac{\theta \hat{N}_{ab}^A}{N_B} + \frac{(1-\theta) \hat{N}_{ab}^B}{N_A} \right] x + \theta \hat{N}_{ab}^A + (1-\theta) \hat{N}_{ab}^B = 0$$

فولر و بورمیستر استدلال کردند که وقتی یک نمونه‌ی تصادفی از هر یک از چارچوب‌ها گرفته می‌شود، برآوردگر رابطه‌ی (۶) می‌تواند با استفاده از اصول ماکسیمم درستنمایی حاصل شود و بنابراین به طور مجانبی کارا است [۶] و [۱۵]. برآوردگر PML ، برآوردگرهای سازگار از N_{ab} ، Y_{ab} ، Y_a و Y_b برای کمیت‌های متناظر در برآوردگر ماکسیمم درستنمایی که با استفاده از نمونه‌گیری تصادفی ساده استخراج می‌شود، جایگذاری می‌کند. برآوردگر PML برخلاف برآوردگر رابطه‌ی (۴)، تحت طرح‌های نمونه‌گیری پیچیده سازگار است. این برآوردگر به واریانس‌های \hat{Y}_{ab}^B و \hat{Y}_{ab}^A بستگی ندارد. اسکینر و راتو [۱۶] استفاده از مقدار $\theta_P = \theta$ که واریانس مجانبی $\hat{N}_{ab}^{PML}(\theta)$ را مینیمم می‌کند، پیشنهاد دادند:

$$(Y) \quad \theta_P = \frac{N_a N_b \text{Var}(\hat{N}_{ab}^B)}{N_a N_b \text{Var}(\hat{N}_{ab}^B) + N_b N_a \text{Var}(\hat{N}_{ab}^A)}$$

برآوردگر رابطه‌ی (۶) برآوردگرهای مجموع سه حوزه یعنی Y_a ، Y_b و Y_{ab} را با برآوردگر بهینه‌ی N_{ab} تعدیل می‌کند. نکته آن‌که محاسبه‌ی θ_P در رابطه‌ی (Y) نیازمند این است که هر سه حوزه‌ی a ، b و ab ناتهی بوده و واریانس‌های \hat{N}_{ab}^B و \hat{N}_{ab}^A مثبت باشند. در عمل، N_a ، N_b ، $\text{Var}(\hat{N}_{ab}^B)$ و $\text{Var}(\hat{N}_{ab}^A)$ از روی داده‌ها برآورد می‌شوند که در نتیجه‌ی آن یک برآوردگر $\hat{\theta}_P$ از θ_P در رابطه‌ی (۶) جایگزین می‌شود. در این صورت وزن‌های تعدیل یافته به صورت زیر هستند:

$$\tilde{W}_{i,S}^A = \begin{cases} \frac{N_a - \hat{N}_{ab}^{PML}(\hat{\theta}_P)}{\hat{N}_a^A} W_i^A & i \in a \\ \frac{\hat{N}_{ab}^{PML}(\hat{\theta}_P)}{\hat{\theta}_P \hat{N}_{ab}^A + (1 - \hat{\theta}_P) \hat{N}_{ab}^B} \hat{\theta}_P W_i^A & i \in ab \end{cases}$$

و

$$\tilde{W}_{i,S}^B = \begin{cases} \frac{N_b - \hat{N}_{ab}^{PML}(\hat{\theta}_P)}{\hat{N}_b^B} W_i^B & i \in b \\ \frac{\hat{N}_{ab}^{PML}(\hat{\theta}_P)}{\hat{\theta}_P \hat{N}_{ab}^A + (1 - \hat{\theta}_P) \hat{N}_{ab}^B} (1 - \hat{\theta}_P) W_i^B & i \in ab \end{cases}$$

$\hat{\theta}_P$ تنها به برآورد واریانس اندازه‌ی حوزه‌ی متداخل بستگی دارد. بنابراین در برآوردگر PML از مجموعه‌ی یکسان وزن‌ها برای هر متغیر پاسخ استفاده می‌شود. اسکینر و راتو

در [۱۶] دریافتند که برآوردگر PML دارای میانگین توان دوم خطای کوچک است و در طرح‌های آمارگیری خوب عمل می‌کند.

۶- شبیه‌سازی

فرض می‌کنیم که جامعه‌ی تحت مطالعه نامتناهی باشد. یک نمونه‌ی خوشه‌ای دو مرحله‌ای با \tilde{n}_A خوشه و m واحد از هر خوشه به‌عنوان نمونه‌ی حاصل از چارچوب A ($n_A = \tilde{n}_A \times m$) و یک نمونه‌ی تصادفی ساده با n_B واحد به‌عنوان نمونه‌ی حاصل از چارچوب B تولید می‌شود. با فرض نامتناهی بودن جامعه، $\frac{N_a}{N}$ و $\frac{N_b}{N}$ را با γ_a و γ_b جایگزین می‌کنیم. نمونه‌ی حاصل از چارچوب A شامل مقادیر

$$\{(y_{ij}, m_{ai}), i = 1, 2, \dots, \tilde{n}_A, j = 1, \dots, m\}$$

می‌باشد که در آن m_{ai} تعداد عناصر نمونه‌گیری شده از i امین خوشه‌ی نمونه، متعلق به حوزه‌ی a و y_{ij} نیز مقدار مرتبط با i امین عنصر نمونه در i امین خوشه‌ی نمونه می‌باشد. نمونه‌ی حاصل از چارچوب B شامل مقادیر نمونه‌ای

$$\{(y_j, n_b), j = 1, \dots, n_B\}$$

است که در آن n_b تعداد عناصر نمونه‌ای متعلق به حوزه‌ی b و y_j مقدار مرتبط با i امین عنصر نمونه است. برای بدست آوردن مقادیر نمونه‌ای حاصل از چارچوب A ، ابتدا m_{ai} را از یک توزیع بتا-دوجمله‌ای به‌صورت زیر تولید می‌کنیم. (۱) \tilde{y}_{ai} ، نسبت مورد انتظار مشاهداتی از خوشه‌ی i ام که در حوزه‌ی a قرار دارد و از توزیع بتا با پارامترهای b_1 و $b_2 = b_1(1 - \gamma_a - \gamma_b)/\gamma_a$ تولید می‌شود بنابراین

$$E(\tilde{y}_{ai}) = \frac{\gamma_a}{1 - \gamma_b} \quad \text{و} \quad \text{Var}(\tilde{y}_{ai}) = \frac{\gamma_a(1 - \gamma_a - \gamma_b)}{[b_1(1 - \gamma_b) + \gamma_a](1 - \gamma_b)^2}$$

در نتیجه مقدار کوچک‌تر b_1 ، تغییرپذیری \tilde{y}_{ai} از یک خوشه به خوشه‌ی دیگر را افزایش می‌دهد.

(۲) به ازای مقدار \tilde{y}_{ai} ، m_{ai} را از یک توزیع دوجمله‌ای با پارامترهای m و \tilde{y}_{ai} تولید می‌کنیم.

(۳) مراحل ۱ و ۲ را مستقلاً تکرار کرده تا $\{m_{ai}, i = 1, 2, \dots, \tilde{n}_A\}$ به دست آید. برای هر خوشه به ازای مقدار m_{ai} مربوط، مقادیر y_{ij} از مدل خطای آشیانی (Nested Error Model) زیر تولید می‌شوند:

$$\begin{aligned} y_{ij} &= \mu_a + \alpha_{ai} + \varepsilon_{ij} & j &= 1, 2, \dots, m_{ai} \\ (۸) \quad y_{ij} &= \mu_{ab} + \alpha_{abi} + \varepsilon_{ij} & j &= m_{ai} + 1, \dots, m \end{aligned}$$

به ترتیب اثر تصادفی خوشه‌ی i ام در حوزه‌های a و ab می‌باشد که از یک توزیع نرمال دو متغیره با میانگین صفر و ماتریس کوواریانس 2×2 با عناصر قطری $\rho\sigma^2$ و عناصر غیر قطری $\delta\rho\sigma^2$ بدست می‌آیند. ε_{ij} ها هم متغیرهای تصادفی مستقل و هم‌توزیع با $(\alpha_{ai}, \alpha_{abi}, m_{ai})$ هستند. این مدل یک همبستگی یکسان $\sigma\rho$ را بین حوزه‌ها تضمین می‌کند.

برای بدست آوردن نمونه‌ی حاصل از چارچوب B ابتدا n_b را از یک توزیع دوجمله‌ای با پارامترهای n_B و $\gamma = \frac{\gamma_b}{1-\gamma_a}$ تولید می‌کنیم. به ازای مقدار n_b حاصل، مقادیر نمونه‌ای

$\{y_j, j = 1, \dots, n_B\}$ از چارچوب B را از مدل زیر تولید می‌کنیم:

$$\begin{aligned} y_j &= \mu_b + \delta_j & j &= 1, 2, \dots, n_b \\ (۹) \quad y_j &= \mu_{ab} + \delta_j & j &= n_b + 1, \dots, n_B \end{aligned}$$

δ_j ها متغیرهای تصادفی مستقل و هم‌توزیع با $(\delta_j \sim N(0, \sigma^2))$ می‌باشند. در اینجا از مقادیر $\mu_a = 9, \mu_{ab} = 10, \mu_b = 11, \rho = 0.1, \delta = 0.5$ و $\sigma^2 = 1$ استفاده می‌کنیم.

برای شبیه‌سازی، نمونه‌گیری را $R = 10000$ بار تکرار کرده و از هر نمونه برآوردهای $\hat{Y}_S(\hat{\theta}_S), \hat{Y}_{PML}(\hat{\theta}_P), \hat{Y}_{FB}(\hat{\beta}_{FB}), \hat{Y}_H(\hat{\theta}_H)$ ، بهینه‌ی $\hat{\theta}_S, \hat{\theta}_P, \hat{\beta}_{FB}, \hat{\theta}_H$ (این پارامترها در بخش ۵ آورده شده‌اند) و به ازای مقادیر مختلف $b_1, b_2, \gamma_a, \gamma_b, \tilde{n}_A$ و n_B محاسبه می‌کنیم. سپس میانگین توان دوم خطای تجربی هر برآوردها، متوسط توان دوم انحراف برآوردها از مقدار واقعی ($EMSE$) را به صورت زیر محاسبه می‌کنیم:

$$(۱۰) \quad EMSE = \frac{1}{R} \sum_{r=1}^R (\hat{Y}_r - Y)^2$$

که در آن عبارت است از \hat{Y} برای r امین تکرار نمونه‌گیری و R تعداد تکرارهای نمونه‌گیری می‌باشد. در رابطه‌ی (۱۰) چون Y معلوم و در دسترس نمی‌باشد در این صورت از امید ریاضی آن به صورت زیر استفاده شده است:

$$(۱۱) \quad \frac{1}{R} \sum_{r=1}^R (\hat{Y}_r - \bar{Y})^2$$

در این شبیه‌سازی، $EMSE$ به ازای $R = ۱۰۰۰۰$ نمونه محاسبه شده است. واریانس مونت کارلوی توان دوم خطا به صورت زیر محاسبه می‌شود:

$$(۱۲) \quad s^2(EMSE) = \frac{1}{R-1} \sum_{r=1}^R (\hat{Z}_r - \hat{Z})^2$$

که در آن $\hat{Z} = \sum_r \hat{Z}_r / R$ و $\hat{Z}_r = (\hat{Y}_r - Y)^2$ می‌باشد. با گرفتن جذر از رابطه‌ی فوق، خطای معیار مونت کارلو را می‌توان محاسبه کرد. تمام محاسبات با استفاده از نرم‌افزار S-plus صورت گرفته است.

جدول (۱) مقادیر $EMSE$ و خطای معیار مونت کارلوی متناظر را به ازای $\tilde{n}_A = ۱۰$ ، $\tilde{n}_A = ۲۰$ و $n_B = ۲۰۰$ ، $m = ۳۰$ ، $\tilde{n}_A = ۲۰$ و $n_B = ۱۰۰$ ، $m = ۳۰$ و $n_B = ۱۰۰$ ، $m = ۳۰$ و سه مقدار مختلف (b_1, b_2) و نیز مقادیر معین γ_a و γ_b نشان می‌دهد.

در جدول ۱ میانگین توان دوم خطای مونت کارلو ($EMSE$) برای ۱۰۰۰۰ بار تکرار نمونه‌گیری و در داخل دو کمان انحراف معیار مونت کارلو $s(EMSE)$ هر برآوردگر داده شده است.

همان‌طور که در جدول ۱ مشاهده می‌شود برآوردگر ماکسیمم درست‌نمایی نما (PML) همواره دارای $EMSE$ کوچک‌تری نسبت به برآوردگرهای دیگر است. در صورتی که $\gamma_a \leq \gamma_b$ یا n_A خیلی بزرگ‌تر از n_B باشد، برآوردگر تک‌چارچوبی (S) در مقایسه با برآوردگرهای دوچارچوبی (H, FB, PML) افزایش قابل ملاحظه‌ای را در $EMSE$ نشان می‌دهد.

جدول ۱- میانگین توان دوم خطا و خطای معیار مونت کارلوی براوردگرهای چارچوب دوگان (میانگین توان دوم خطا $(EMSE) \times 100$)

<i>S</i>	<i>PML</i>	<i>FB</i>	<i>H</i>		
$n_B = 100 \quad m = 30 \quad \tilde{n}_A = 10 \quad n_A = 300 \quad b_1 = 2, 2/5, 3$					
۶/۶۸(۰/۰۹)	۴/۴(۰/۰۶۱)	۵/۱۳(۰/۰۸۱)	۴/۷(۰/۰۶۵)	۰/۱	۰/۱
۱۱/۴۲(۰/۱۵۵)	۴/۵۷(۰/۰۶)	۵/۱۵(۰/۰۷۸)	۵/۰۵(۰/۰۷۱)	۰/۲	۰/۱
۵/۸(۰/۰۸۳)	۵/۴۶(۰/۰۷۹)	۵/۹(۰/۰۸۶)	۵/۷(۰/۰۸۱)	۰/۱	۰/۲
$n_B = 200 \quad m = 30 \quad \tilde{n}_A = 20 \quad n_A = 600 \quad b_1 = 2, 2/5, 3$					
۳/۵۲(۰/۰۴۸)	۲/۱۳(۰/۰۳۱)	۲/۴۷(۰/۰۳۹)	۲/۳۷(۰/۰۳۳)	۰/۱	۰/۱
۵/۴۵(۰/۰۷۶)	۲/۳(۰/۰۳۳)	۲/۴۳(۰/۰۳۷)	۲/۳۷(۰/۰۳۴)	۰/۲	۰/۱
۲/۹۸(۰/۰۴۲)	۲/۷(۰/۰۳۸)	۲/۸۷(۰/۰۴۱)	۲/۷۵(۰/۰۳۹)	۰/۱	۰/۲
$n_B = 100 \quad m = 30 \quad \tilde{n}_A = 20 \quad n_A = 600 \quad b_1 = 2, 2/5, 3$					
۷/۸۳(۰/۱۲)	۲/۸۵(۰/۰۴۱)	۳/۱۲(۰/۰۴۴)	۳/۲۱(۰/۰۴۶)	۰/۱	۰/۱
۱۳/۳۱(۰/۱۹)	۲/۷۴(۰/۰۴۱)	۲/۸۳(۰/۰۴)	۳/۰۷(۰/۰۴۳)	۰/۲	۰/۱
۶/۶(۰/۰۹۵)	۳/۹۲(۰/۰۵۵)	۴/۲۶(۰/۰۵۹)	۴/۲۵(۰/۰۵۹)	۰/۱	۰/۲

برای مثال در جدول ۱ برای حالتی که $\gamma_a = 0/1$ ، $\gamma_b = 0/2$ و $n_A/n_B = 3$ داریم:
 $EMSE(\hat{Y}_S) = 5/45$ $EMSE(\hat{Y}_{PML}) = 2/3$ $EMSE(\hat{Y}_{FB}) = 2/43$
 $EMSE(\hat{Y}_H) = 2/37$

و وقتی که n_A/n_B به شش افزایش می‌یابد، داریم:

$EMSE(\hat{Y}_S) = 13/31$ $EMSE(\hat{Y}_{PML}) = 2/74$ $EMSE(\hat{Y}_{FB}) = 2/83$
 $EMSE(\hat{Y}_H) = 3/07$

حال اگر $\gamma_a = 0/2$ ، $\gamma_b = 0/1$ و $n_A/n_B = 3$ کارایی حاصل زیاد نخواهد بود. برای مثال:

$EMSE(\hat{Y}_S) = 2/98$ $EMSE(\hat{Y}_{PML}) = 2/7$ $EMSE(\hat{Y}_{FB}) = 2/87$
 $EMSE(\hat{Y}_H) = 2/75$

و وقتی که n_A/n_B به شش افزایش می‌یابد، داریم:

$$EMSE(\hat{Y}_S) = 6/6 \quad EMSE(\hat{Y}_{PML}) = 3/92 \quad EMSE(\hat{Y}_{FB}) = 4/26$$

$$EMSE(\hat{Y}_H) = 4/25$$

حال می‌خواهیم اثر افزایش اندازه‌های نمونه‌ای مختلف را بر میانگین توان دوم خطای برآوردگرهای دوچارچوبی H, FB, PML و برآوردگر تک‌چارچوبی (S) بررسی کنیم.

الف- زمانی که n_B ثابت و n_A افزایش یابد برای تمامی حالت‌های مختلف γ_a و γ_b ($\gamma_a > \gamma_b$ و $\gamma_a < \gamma_b$ ، $\gamma_a = \gamma_b$) میانگین توان دوم خطا برای برآوردگرهای دوچارچوبی PML, FB, H و کاهش و برای برآوردگر S افزایش می‌یابد.

ب- زمانی که n_B افزایش و n_A ثابت باشد برای تمامی حالت‌های مختلف γ_a و γ_b میانگین توان دوم خطای تمامی برآوردگرها کاهش می‌یابند. در این حالت $EMSE$ برآوردگر PML کم‌تر از $EMSE$ برآوردگرهای S, FB, H می‌باشد.

پ- زمانی که هر دو n_A و n_B افزایش یابند برای تمامی حالت‌های مختلف γ_a و γ_b میانگین توان دوم خطای تمامی برآوردگرها کاهش می‌یابند. در این حالت نیز $EMSE$ برآوردگر PML کم‌تر از $EMSE$ برآوردگرهای S, FB, H می‌باشد.

در نهایت، اثر افزایش γ_a و γ_b را بر میانگین توان دوم خطای برآوردگرهای دوچارچوبی PML, FB, H و برآوردگر تک‌چارچوبی (S) بررسی کنیم.

الف- زمانی که γ_a ثابت و γ_b افزایش یابد، در این صورت برای حالتی که $n_A = 300$ و $n_B = 100$ است، میانگین توان دوم خطای تمامی برآوردگرها افزایش می‌یابند، برای حالتی که $n_A = 600$ و $n_B = 100$ است، میانگین توان دوم خطای برآوردگرهای دوچارچوبی PML, FB, H کاهش و برآوردگر S افزایش می‌یابد و برای حالتی که $n_A = 600$ و $n_B = 200$ است، میانگین توان دوم خطای برآوردگرهای PML, FB و H تقریباً نزدیک به هم هستند اما میانگین توان دوم خطای برآوردگر S افزایش داشته است. در این حالت‌ها کم‌ترین $EMSE$ مربوط به برآوردگر PML است.

ب- زمانی که γ_a افزایش و γ_b ثابت باشد، در این صورت برای حالت‌های $n_A = 300$ و $n_B = 100$ ، $n_A = 600$ و $n_B = 100$ و $n_A = 600$ و $n_B = 100$ و $n_A = 600$ و $n_B = 200$ میانگین توان دوم خطای برآوردگرهای دوچارچوبی PML, FB, H و S کاهش می‌یابد. در این حالت با وجود آن‌که $EMSE$ برآوردگر S کاهش یافته ولی مقدار آن

بزرگ‌تر از $EMSE$ برآوردهای PML ، FB و H است و کم‌ترین $EMSE$ مربوط به برآوردهای PML می‌باشد.

۷- بحث و نتیجه‌گیری

آمارگیری‌های چندچارچوبی هنگامی که جوامع کمیاب هستند یا چارچوب کامل موجود برای نمونه‌گیری پرهزینه باشد، بسیار مفید و سودمند است.

در آمارگیری‌های چندچارچوبی، برآوردهای مختلفی برای برآورد مجموع جامعه پیشنهاد داده شده است. اما انتخاب برآوردهای بستگی به طراحی و پیچیدگی طرح آمارگیری دارد.

برآوردهای ماکسیمم درست‌نمایی نما به دلیل این‌که از وزن‌های یکسانی برای همه‌ی متغیرها استفاده می‌کند، نسبت به برآوردهای هارتلی و فولر-بورمیستر در بررسی‌هایی با متغیرهای متعدد، ارجحیت دارد.

با افزایش هم‌زمان اندازه‌های نمونه در دو چارچوب (n_A و n_B)، برای تمامی حالت‌های مختلف γ_a و γ_b ($\gamma_a > \gamma_b$ و $\gamma_a < \gamma_b$ ، $\gamma_a = \gamma_b$)، میانگین توان دوم خطای تمامی برآوردهای (S و H ، FB ، PML) کاهش می‌یابند.

مرجع‌ها

- [۱] نورینی، مرجان (۱۳۸۲). طرح‌های نمونه‌گیری چندچارچوبی. پایان‌نامه‌ی کارشناسی ارشد. دانشگاه صنعتی اصفهان، اصفهان.
- [2] Bankier, M.D. (1986). Estimators based on several stratified samples with applications to multiple frame surveys. *Journal of the American Statistical Association*, **81**, 1074–1079.
- [3] Casady, R.J., Snowden, C.B. and Sirken, M.G. (1981). A study of dual frame estimators for the national health interview survey. Proceedings of the Survey Research Section, *American Statistical Association*, 444–447.
- [4] Cochran, R.S. (1964). Multiple-frame Surveys, Proceedings at American Statistical Association Meeting, Univ. Wyoming, Dec.

- [5] Ford, B.L. and Bosecker, R.R. (1979). Multiple frame estimation with stratified overlap domain, in proceeding of the Social Statistics Section. *American Statistical Association*, 219–224.
- [6] Fuller, W.A. and Burmeister, L.F. (1972). Estimators for samples selected from two overlapping frames, in proceeding of the social statistics section. *American Statistical Association*, 245–249.
- [7] Groves, R.M. and Lepkowski, J.M. (1986). A mean squared error model for dual frame mixed model survey design. *Journal of the American Statistical Association*, **81**, 930–937.
- [8] Hájek, J. (1981). *Sampling from a Finite Population*. Marcel Dekker, New York.
- [9] Hartley, H.O. (1962). Multiple frame survey, In Proceeding of the Social Statistics Section, *American Statistical Association*, 203–206.
- [10] Hartley, H.O. (1974). Multiple frame methodology and selected application. *Sankhya*, Ser. C, **36**, 99–118.
- [11] Kalton, G. and Anderson, D.W. (1986). Sampling rare populations. *Journal of the Royal Statistical Society*, Series A **149**, 65–82.
- [12] Lohr, S.L. and Rao, J.N.K. (2006). Estimation in multiple-frame surveys. *Journal of the American Statistical Association*, **101**, 1019–1030.
- [13] Lund, R.E. (1968). Estimators in multiple frame surveys, in proceeding of the social statistics section. *American Statistical Association*, 282–288.
- [14] Sirken, M.G. and Cassady, R.J. (1982). Nonresponse in dual frame survey based on area/list and telephone frames. *American Statistical Association*, 151–153.
- [15] Skinner, C.J. (1991). On the efficiency of raking ratio estimation for multiple frame surveys. *Journal of the American Statistical Association*, **86**, 779–784.

- [16] Skinner, C.J. and Rao, J.N.K. (1996). Estimation in dual frame survey with complex designs. *Journal of the American Statistical Association*, **91**, 349–356.
- [17] Vogel, F.A. (1975). Survey with overlapping frame—problems in application, in proceedings of the social statistics. *American Statistical Association*, 694–699.

مرجان نورینی

کارشناس ارشد آمار

تهران، خیابان دکتر فاطمی، نیش رهی معیری، پلاک ۱، مرکز آمار ایران.

رایانشانی: mar_noorini@yahoo.com

انور قیطولی

کارشناس ارشد آمار

تهران، خیابان دکتر فاطمی، نیش رهی معیری، پلاک ۱، مرکز آمار ایران.

رایانشانی: a.ghaitoly@gmail.com