

کاربست طبقه‌بندی بهینه فضایی در آمارگیری هزینه و درآمد خانوار برای ارائه برآورد به روش نمونه‌گیری طرح‌مبنای فضایی

لیدا کلهری ندرآبادی^۱، روشنگ علی‌اکبری صبا^۲، آسیه عباسی^۳

تاریخ دریافت: ۱۴۰۰/۰۶/۱۵

تاریخ پذیرش: ۱۴۰۰/۱۲/۲۶

چکیده:

آمارگیری هزینه و درآمد خانوار یکی از مهم‌ترین آمارگیری‌های مرکز آمار ایران است که پارامترهای اصلی آن همبسته فضایی هستند. وقتی همبستگی فضایی میان واحدهای جامعه وجود دارد، انتخاب نمونه‌های مستقل به روش کلاسیک به دلیل برقرار نبودن شرط استقلال واحدهای جامعه با چالش‌های بسیاری مواجه است. استفاده از نمونه‌گیری فضایی راه‌حلی برای مواجهه با این مشکل است. به‌کارگیری نمونه‌گیری فضایی به دلیل دسترسی نداشتن به چارچوب مناسب، در آمار رسمی کمتر مورد توجه واقع شده است. در این مقاله یک روش طرح‌مبنای مدل‌یار برای طبقه‌بندی بهینه فضایی جامعه هدف مرور می‌شود. در حال حاضر اطلاعات مکانی واحدهای جامعه در چارچوب نمونه‌گیری آمارگیری هزینه و درآمد وجود ندارد، اما دسترسی به اطلاعات مکانی برخی از واحدهای نمونه توسط مرکز آمار ایران برای این مطالعه محقق شده است. تولید داده‌های مکانی یکی از مؤلفه‌های اصلی در مدرن‌سازی نظام آماری است و مورد توجه مرکز آمار ایران قرار دارد. بنابراین در این مقاله با شبیه‌سازی چارچوب فضایی نمونه‌گیری بر اساس الگوی داده‌های هزینه و درآمد، کاربرد طبقه‌بندی بهینه فضایی بر اساس فاصله تعمیم‌یافته با استفاده از خطای پیشگویی انجام می‌شود. نتایج نشان‌دهنده افزایش کارایی روش نمونه‌گیری با این طبقه‌بندی در مقایسه با نمونه‌گیری تصادفی ساده در سطح نواحی جغرافیایی است. همچنین نتایج شبیه‌سازی با شبکه‌هایی با اندازه‌های مختلف و میزان همبستگی متفاوت حاکی از کارایی این روش در مقایسه با روش فعلی آمارگیری هزینه و درآمد است.

واژه‌های کلیدی: طبقه‌بندی فضایی، فاصله تعمیم‌یافته، آمارگیری هزینه و درآمد خانوار

۱ مقدمه

استفاده شود. به این طریق می‌توان از گردآوری اطلاعات تکراری تا حدودی در نواحی نمونه‌گیری شده وجود دارد، بدون اینکه از قابلیت اعتماد برآوردها کاسته شود جلوگیری کرد و در هزینه‌های نمونه‌گیری صرفه‌جویی نمود. در نتیجه علاوه بر افزایش دقت نمونه‌گیری در برآورد پارامترهای جامعه، می‌توان الگوهای فضایی را نیز شناسایی و در صورت نیاز بر اساس آن‌ها پیشگویی انجام داد. خاورزاده و همکاران [۱۰] طرح نمونه‌گیری فضایی متعادل دومرحله‌ای برای پیشگویی میدان‌های تصادفی را مورد بررسی قرار داده‌اند. به‌کارگیری نمونه‌گیری فضایی در طیف گسترده‌ای از مطالعات کاربردی از جمله بررسی خاک، محیط‌زیست، هواشناسی و مطالعات پزشکی مورد توجه قرار گرفته ولی

یکی از ارکان جمع‌آوری اطلاعات برای دستیابی به نتایج آماری معتبر، انتخاب طرح نمونه‌گیری متناسب با شرایط جامعه است. در مواجهه با داده‌هایی که همبستگی آن‌ها ناشی از موقعیت قرارگیری در یک فضای معین است، یعنی دارای همبستگی فضایی هستند استفاده از رویکردهای مرسوم نمونه‌گیری که در آن‌ها همبستگی فضایی نادیده گرفته می‌شود، می‌تواند اعتبار نتایج را تحت تأثیر قرار دهد (محمدزاده، [۱]). از لحاظ شهودی واضح است وقتی اثر یا نشانه‌ای از ساختار همبستگی فضایی در پدیده‌ی فضایی که قرار است نمونه‌گیری شود وجود داشته باشد، مطلوب است که این اطلاعات در طرح نمونه‌گیری

^۱ استادیار پژوهشکده‌ی آمار (نویسنده مسئول): kalhori@src.ac.ir

^۲ استادیار پژوهشکده‌ی آمار، r_saba@src.ac.ir

^۳ کارشناس مرکز آمار ایران، asieh_abasi@yahoo.com

همبستگی و استفاده از پیش‌گویی فضایی می‌توان روش نمونه‌گیری طبقه‌بندی طرح‌مبنای فضایی را به کار گرفت که در آن طبقه‌بندی بر اساس فاصله‌ی موقعیت‌ها و الگوی همبستگی صورت می‌گیرد. با توجه به اینکه چارچوب فعلی آمارگیری هزینه و درآمد خانوار اطلاعات مربوط به موقعیت جغرافیایی واحدهای نمونه را شامل نمی‌شود، دستیابی به اطلاعات کمی برای شناسایی الگوی همبستگی داده‌ها حائز اهمیت است. استفاده از متغیر همبسته با درآمد که امکان دستیابی به آن از سایر منابع برحسب موقعیت‌های جغرافیایی موجود باشد توسط کلهری [۹] بررسی شده است. در بخش سوم، شبیه‌سازی چارچوب آمارگیری هزینه و درآمد خانوار و کاربری طبقه‌بندی فضایی با استفاده از این چارچوب مورد بررسی قرار گرفته است. بخش چهارم به بحث و نتیجه‌گیری اختصاص یافته است.

۲. ارائه‌ی برآوردهای طرح‌مبنای مبتنی بر بهینه‌سازی طبقه‌بندی فضایی با استفاده از خطای پیشگویی

طبقه‌بندی بهینه توسط بسیاری از آماردانان مورد توجه قرار گرفته است، ولی فرض وجود همبستگی فضایی در آن‌ها در نظر گرفته نشده است. ارائه‌ی برآورد میانگین متغیر هدف در یک ناحیه‌ی مشخص بر اساس اطلاعات حاصل از یک آمارگیری نمونه‌ای، همواره مسئله مهمی بوده است. معمولاً انگیزه‌های اقتصادی برای ارائه‌ی این برآوردها وجود دارند و بنابراین بهینه‌سازی برآوردها و کمی‌سازی عدم حتمیت از اهمیت ویژه‌ای برخوردار است. روش‌های طرح‌مبنا بر این اساس توسعه یافته‌اند و نمونه‌گیری تصادفی طبقه‌بندی شده یکی از روش‌های کارا و متداول است. تعداد بهینه‌ی طبقات و شکل آن‌ها و چگونگی تخصیص تعداد نمونه به هر طبقه از جمله مراحلی است که در طراحی نمونه‌گیری مورد توجه قرار می‌گیرد. با توسعه‌ی تصاویر ماهواره‌ای، اطلاعات کمی به صورت جغرافیایی به طور گسترده‌ای در دسترس قرار گرفته‌اند. وجود اطلاعات جغرافیایی امکان پیشگویی فضایی متغیر هدف را فراهم کرده است که در طراحی نمونه‌گیری قابل استفاده است.

یک مشکل اساسی که در پیشگویی‌ها وجود دارد میزان عدم حتمیت آن‌ها است که می‌تواند به صورت جدی باعث بیش‌برآوردی یا کم‌برآوردی در برآورد واریانس متغیر هدف در هر طبقه باشد. سؤالی که مطرح

در حوزه‌ی آمار رسمی کمتر به آن پرداخته شده است (آبی، [۲]). با توجه به اینکه درآمد خانوارها همبسته فضایی هستند و هزینه‌های خانوارها نیز با میزان درآمد آن‌ها همبسته است، انتظار می‌رود به‌کارگیری نمونه‌گیری فضایی در طرح آمارگیری هزینه و درآمد خانوار کیفیت نتایج را بهبود ببخشد.

یکی از مهم‌ترین ابزارهای مورد نیاز برای انتخاب نمونه چارچوب نمونه‌گیری است. ساخت و مشخصه‌سازی چارچوب نمونه‌گیری از جوامع بشری می‌تواند به علت نبود اطلاعات در خصوص بعضی از واحدهای جامعه یا در دسترس نبودن ویژگی‌های اجتماعی-اقتصادی آن‌ها با چالش‌های بسیار همراه باشد. در نمونه‌گیری احتمالی ناحیه‌ای، واحدهای نمونه‌گیری اولیه انتخاب می‌شوند و سپس خانوارهای نمونه در این واحدها به یکی از روش‌های کلاسیک نمونه‌گیری انتخاب می‌شوند. مشکل این روش آن است که اطلاعات به دست آمده تجمیع می‌شوند و تعیین توزیع جمعیت در هر واحد نمونه دشوار است. یک راهکار برای تعیین توزیع جغرافیایی جمعیت، استفاده از نقشه‌های اسکن زمین است. اما به دلیل اینکه چنین نقشه‌هایی برای مطالعه و بررسی آمارگیری هزینه و درآمد در دسترس نیست استفاده از این روش در حال حاضر میسر نیست و در صورتی که در آینده چنین نقشه‌هایی تهیه شوند امکان استفاده از آن‌ها وجود خواهد داشت. یک راه‌حل در وضعیت فعلی، شبکه‌بندی نقشه‌های موجود است اما به دلیل این‌که اطلاعات سرشماری به صورت مکانی وجود ندارند امکان تعیین توزیع جمعیت در هر پیکسل وجود ندارد و نمی‌توان انتخاب نمونه متناسب با جمعیت پیکسل‌ها و پس‌از آن فرآیند وزن‌دهی را به انجام رساند. از طرف دیگر در نمونه‌گیری‌های فضایی که به صورت ناحیه‌ای انجام می‌شوند الگوی همبستگی فضایی واحدها اغلب نادیده گرفته می‌شود. الگوی همبستگی فضایی در طرح‌های نمونه‌گیری مدل‌مبنا لحاظ می‌شوند اما با توجه به اینکه مراکز آماری در دنیا روش‌های نمونه‌گیری طرح‌مبنا را به کار می‌گیرند و با توجه به اینکه روش‌های متداول معرفی شده برای نمونه‌گیری فضایی مدل‌مبنا هستند، نمونه‌گیری فضایی در مراکز آماری دنیا استفاده نشده است. در این مقاله کاربری روش نمونه‌گیری طبقه‌بندی طرح‌مبنای فضایی در آمارگیری هزینه و درآمد خانوار با مطالعات شبیه‌سازی و بررسی داده‌های طرح آمارگیری از هزینه و درآمد خانوار ارزیابی می‌شود.

در بخش دوم مقاله، روش ارائه‌شده توسط دی‌گروتر و همکاران [۷] مرور می‌شود که یک روش طرح‌مبنای فضایی است و با مشکلات طرح‌های ناحیه‌ای مواجه نیست. علاوه بر آن، این روش الگوی همبستگی فضایی را در طبقه‌بندی ناحیه‌ی مورد مطالعه لحاظ می‌کند. با شناسایی الگوی

می‌شود این است که اگر مقادیر پیشگویی شده به صورت جدی فاقد اطلاع باشند، چگونه می‌توان طبقه‌بندی انجام داد. در این شرایط و با فرض وجود همبستگی فضایی دی‌گروتر و همکاران [۷] روشی ارائه کرده‌اند که در ادامه خواهد آمد. در این روش برای انجام طبقه‌بندی از سه اطلاع شامل مقادیر پیشگویی، واریانس خطای پیشگویی و فاصله‌ی جغرافیایی بین موقعیت‌ها استفاده می‌شود. این سه اطلاع برای ساختن اندازه‌ای از فاصله (ماتریس عدم شباهت) استفاده می‌شود. این رهیافت در وضعیت‌هایی قابل استفاده است که نمونه‌گیری تصادفی ساده در طبقه قابل اجرا باشد و هزینه‌ی نمونه‌گیری از واحدها تقریباً یکسان باشد.

۱.۲ طبقه‌بندی فضایی بهینه و تخصیص نمونه

در مطالعات غیر فضایی معمولاً برای طبقه‌بندی از روش ریشه‌ی تجمعی دالینوس و هوجز [۶] یا سایر روش‌های معرفی شده توسط بیلگرن و ریوست [۳]، بالین و بارکارلی [۴]، کوزاک [۱۱] و هارگون [۸] استفاده می‌شود. این روش‌ها را می‌توان در مطالعات فضایی نیز مورد استفاده قرار داد ولی مشکلی که ایجاد می‌شود آن است که فرض می‌شود پیشگویی‌ها خطای ناچیزی دارند یا این‌که اثر معنی‌داری بر بهینگی طبقه‌بندی ندارند. ولی در عمل این فرض نمی‌تواند صحیح باشد. دی‌گروتر و همکاران [۷] طبقه‌بندی بر اساس مقادیر پیشگویی در موقعیت‌های جغرافیایی را با در نظر گرفتن پیشگویی ارائه نمودند. در رهیافت آن‌ها اطلاعات موقعیت‌ها، پیشگویی و واریانس پیشگو به‌عنوان یک اندازه‌ی واحد ادغام شده است و این اندازه‌ی واحد به‌عنوان فاصله تعمیم‌یافته معرفی شده است. دی‌گروتر و همکاران [۷]، روشی برای طبقه‌بندی ناحیه‌ی آمارگیری در فضایی اقلیدسی به‌منظور کمینه کردن واریانس برآوردگر هورویتز تامپسون ارائه نمودند و آن را طبقه‌بندی بهینه‌ی فضایی^۴ (OSPATS) نام نهادند. آن‌ها با فرض این‌که ناحیه‌ی مورد مطالعه مشبکه منظم از N موقعیت متناهی است از اطلاعات مربوط به موقعیت‌های جغرافیایی، مقادیر پیشگویی و واریانس پیشگویی برای معرفی یک فاصله جدید استفاده نمودند. همان‌طور که پیش‌ازین ذکر شد، هدف دی‌گروتر و همکاران [۷] کمینه کردن واریانس برآوردگر هورویتز تامپسون در ارائه‌ی برآوردهای جامعه‌ای با همبستگی فضایی بوده است. در ادامه، بدون اینکه از کلیت مسئله کاسته شود فرض شده است که پارامتر مورد نظر میانگین جامعه باشد، بنابراین واریانس نمونه‌گیری به صورت زیر تعریف می‌شود

$$\text{Var}(\hat{z}) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \cdot \text{Var}(\hat{z}_h) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \cdot \frac{S_h^2}{n_h} \quad (1)$$

که در آن

H : تعداد طبقات

N_h : تعداد نقاط مشبکه در طبقه h ام، (اندازه طبقه h)

N : تعداد کل نقاط مشبکه در ناحیه‌ی مورد مطالعه، (اندازه‌ی جامعه)،

\hat{z}_h : برآورد میانگین در طبقه h ام،

S_h^2 : واریانس z در طبقه h ،

n_h : اندازه نمونه در طبقه h ، است. تعداد نمونه در هر طبقه با استفاده از

تخصیص نیمن تعیین می‌شود. اگر اندازه‌ی نمونه‌ی تعیین شده n باشد و

هزینه‌ی نمونه‌گیری هر واحد در همه‌ی طبقات یکسان باشد آنگاه اندازه

نمونه در هر طبقه از رابطه‌ی زیر تعیین می‌شود

$$n_h = n \cdot \frac{N_h S_h}{\sum_{h=1}^H N_h S_h} \quad (2)$$

در حالتی که هیچ اطلاعی وجود نداشته باشد از الگوریتم k - میانگین استفاده می‌شود و بر اساس فاصله‌ی اقلیدسی بین موقعیت‌ها طبقه‌بندی انجام می‌شود. اما اگر پیشگویی متغیر هدف در تمام نقاط در دسترس باشد می‌تواند برای طبقه‌بندی مورد استفاده قرار بگیرد. ماتریس فاصله بر اساس اطلاعات مقادیر پیشگویی، واریانس خطای پیشگویی و فاصله‌ی جغرافیایی ساخته می‌شود. بنابراین اطلاعات مورد نیاز برای طبقه‌بندی شامل مختصات جغرافیایی، مقادیر پیشگویی متغیر هدف و واریانس خطای پیشگویی هستند و پیشگویی فضایی با روش‌های آماری مانند رگرسیون فضایی یا کریگیدن انجام می‌شود و واریانس خطای پیشگویی را فراهم می‌کند. با جایگذاری n_h در $\text{Var}(\hat{z})$ و قرار دادن $a_h = \frac{N_h}{N}$ داریم

$$\text{Var}(\hat{z}) = \frac{1}{n} \left\{ \sum_{h=1}^H a_h S_h \right\}^2 \quad (3)$$

و به این ترتیب برای هر نمونه به اندازه n معیار بهینگی به صورت $\bar{O} = \left\{ \sum_{h=1}^H a_h S_h \right\}$ تعریف می‌شود و واریانس بهینه با کمینه کردن آن به دست می‌آید. واریانس در هر طبقه بر اساس تفاضل مقدار پاسخ در هر موقعیت و میانگین طبقه h به دست می‌آید. بنابراین می‌توان به صورت مستقیم از تفاضل داده‌ها در موقعیت‌های هر طبقه استفاده کرد، یعنی در هر طبقه h برای $i, j = 1, \dots, N_h$

$$((z_i - \bar{z}_h) - (z_j - \bar{z}_h))^2 = (z_i - z_j)^2.$$

بنابراین واریانس در طبقه h به صورت زیر تعریف می‌شود

$$S_h^2 = \frac{1}{2N_h} \sum_{i=1}^{N_h} \sum_{j=1}^{N_h} (z_i - z_j)^2 = \frac{1}{N_h} \sum_{i=1}^{N_h-1} \sum_{j=i+1}^{N_h} (z_i - z_j)^2 \quad (4)$$

⁴Optimize Spatial Stratification

حداقل ثابت بماند. بنابراین هر اندازه که فاصله کمتر باشد، کواریانس بیشتر می‌شود و در نتیجه D_{ij}^y کمتر می‌شود و احتمال قرار گرفتن واحدهای نمونه مشابه در یک طبقه افزایش می‌یابد. پس از محاسبه D_{ij}^y ها طبقه‌بندی به روش تخصیص تکرارشونده انجام می‌شود. در اولین گام، یک طبقه‌بندی اولیه انجام می‌شود و در گام‌های بعدی یک روش تکراری با انتقال موقعیت‌ها از یک طبقه به طبقه‌ی دیگر انجام می‌شود و این روند تا زمانی ادامه می‌یابد که کاهش بیشتری در تابع هدف رخ ندهد. میزان کاهش تابع هدف به ازای انتقال یک موقعیت از یک طبقه به طبقه‌ی دیگر محاسبه می‌شود. سهم طبقه‌ی A در معیار O به صورت زیر است:

$$O(A) = \left\{ \sum_{i=1}^{N_A-1} \sum_{j=i+1}^{N_A} D_{ij}^y \right\}^{\frac{1}{2}} \quad (9)$$

و سهم طبقه‌ی B نیز به صورت مشابه محاسبه می‌شود. بعد از انتقال موقعیت مکانی st از طبقه‌ی A به طبقه‌ی B ، سهم طبقه‌ی A از رابطه زیر به دست می‌آید

$$O(A-t) = \left\{ \sum_{i=1}^{N_A-1} \sum_{j=i+1}^{N_A} D_{ij}^y - \sum_{i=1}^{N_A} D_{it}^y \right\}^{\frac{1}{2}} \quad (10)$$

و سهم طبقه‌ی B از رابطه زیر به دست می‌آید

$$O(B+t) = \left\{ \sum_{i=1}^{N_B-1} \sum_{j=i+1}^{N_B} D_{ij}^y + \sum_{i=1}^{N_B} D_{it}^y \right\}^{\frac{1}{2}} \quad (11)$$

میزان کاهش Δ در معیار O به ازای این انتقال به صورت

$$\Delta = O(A-t) - O(A) + O(B+t) - O(B) \quad (12)$$

محاسبه می‌شود. زمانی که جابجایی موقعیت‌ها در طبقات منجر به تغییر Δ نشود، طبقه‌بندی نهایی تلقی می‌شود.

۲.۲ تعیین اندازه‌ی نمونه و محاسبه‌ی واریانس نمونه‌گیری

پس از اتمام طبقه‌بندی و تعیین اندازه‌ی نمونه‌ی کل، اندازه‌ی نمونه در هر طبقه و واریانس نمونه‌گیری به صورت زیر محاسبه می‌شوند. تعیین اندازه‌ی نمونه‌ی بهینه‌ی هر طبقه برحسب سهم تغییرپذیری در آن طبقه از کل تغییرپذیری تابع هدف O تعیین می‌شود. در صورتی که واریانس متغیر پاسخ در هر طبقه مشخص باشد، اندازه‌ی نمونه‌ی بهینه در هر طبقه از رابطه $n_h = n \cdot \frac{N_h S_h}{\sum_{h=1}^H N_h S_h}$ به دست می‌آید.

تغییرپذیری در هر طبقه h از رابطه $O(h) = \left\{ \sum_{i=1}^{N_h-1} \sum_{j=i+1}^{N_h} D_{ij}^y \right\}^{\frac{1}{2}}$

با قرار دادن $d_{ij}^y = (z_i - z_j)^2$ معیار بهینگی \bar{O} به صورت زیر بازنویسی می‌شود

$$\bar{O} = \frac{1}{N} \sum_{h=1}^H \left\{ \sum_{i=1}^{N_{h-1}} \sum_{j=i+1}^{N_h} (d_{ij}^y)^2 \right\}^{\frac{1}{2}} \quad (5)$$

با توجه به اینکه اندازه‌ی شبکه یعنی N ثابت است رابطه‌ی بالا به صورت

$$O = \sum_{h=1}^H \left\{ \sum_{i=1}^{N_{h-1}} \sum_{j=i+1}^{N_h} (d_{ij}^y)^2 \right\}^{\frac{1}{2}} \quad (6)$$

ساده می‌شود. d_{ij}^y ها در واقع توان دوم تفاضلات مقادیر واقعی میدان فضایی هستند که در عمل مقدار واقعی آن‌ها در دسترس نیست. بنابراین باید مقادیر واقعی و خطای واریانس آن‌ها پیشگویی شوند تا از طریق آن d_{ij}^y پیشگویی شوند. مقدار پیشگویی \tilde{z}_i در موقعیت مکانی s_i به صورت مجموع مقدار واقعی و خطای تصادفی e_i در موقعیت i ام یعنی $\tilde{z}_i = z_i + e_i$ نوشته می‌شود. فرض می‌شود پیشگویی‌ها طی فرآیند تصادفی ξ انجام می‌شوند و بنابراین مقادیر امید ریاضی و واریانس قابل محاسبه و درون‌یابی هستند. ویژگی‌های فرآیند ξ در ادامه آمده است:

- ۱- پیشگویی‌ها نارایب باشند یعنی امید ریاضی مؤلفه خطا صفر، $E_{\xi}(e_i) = 0$ ، با مقادیر واریانس $\text{Var}(e_1) \dots \text{Var}(e_N)$ باشند.
- ۲- \tilde{z}_i و e_i ناهمبسته باشند.

مقدار $(\tilde{z}_i - \tilde{z}_j)^2$ برای $i, j = 1, \dots, N_h$ در هر طبقه، می‌تواند یک پیشگوی ساده برای d_{ij}^y باشد ولی به دلیل لحاظ نکردن مؤلفه‌های خطای مربوط به پیشگویی، ممکن است برآورد حاصل نارایب نباشد. بنابراین امید ریاضی d_{ij}^y به شرط \tilde{z}_i تا \tilde{z}_N به صورت زیر به عنوان پیشگو محاسبه می‌شود

$$D_{ij}^y = E_{\xi}(d_{ij}^y | \tilde{z}_i : i = 1, \dots, N) \quad (7)$$

به عبارت دیگر

$$\begin{aligned} D_{ij}^y &= E_{\xi}(z_i - z_j)^2 = (\tilde{z}_i - \tilde{z}_j)^2 + E_{\xi}(e_i - e_j)^2 \\ &= (\tilde{z}_i - \tilde{z}_j)^2 + E_{\xi}(e_i^2 - e_j^2 - 2e_i e_j) \\ &= (\tilde{z}_i - \tilde{z}_j)^2 + \text{Var}(e_i) + \text{Var}(e_j) - 2\text{Cov}(e_i, e_j) \end{aligned} \quad (8)$$

فرض می‌شود که مقادیر پیشگویی و واریانس آن‌ها به عنوان اطلاعات اولیه در دسترس هستند. مؤلفه‌ی کواریانس هم از ساختار تابع اتوکواریانس خطای پیشگویی قابل دست‌یابی است. انتظار می‌رود کواریانس با افزایش تأخیر فضایی کاهش یافته یا این‌که

۳ طبقه‌بندی فضایی بهینه‌ی آمارگیری هزینه و درآمد با استفاده از خطای پیشگویی

در بخش ۲ یک روش جدید برای طبقه‌بندی فضایی بهینه‌ی ناحیه‌ی مورد مطالعه مرور شد که یک روش طرح‌مبنای مدلیار است. این روش از طریق تعریف یک فاصله جدید تعمیم‌یافته که ترکیبی از فاصله‌ی اقلیدسی و پیشگویی فضایی است، طبقه‌بندی را به‌گونه‌ای انجام می‌دهد که واریانس برآوردگر هورویتر تامپسون کمینه شود.

اطلاعات مکانی واحدهای جامعه در چارچوب آمارگیری هزینه و درآمد خانوار وجود ندارد. برای انجام این پژوهش اطلاعات مکانی بخشی از واحدهای نمونه سال ۱۳۹۶ در مناطق شهری از دفتر نقشه و اطلاعات مکانی مرکز آمار ایران دریافت شده است که به‌عنوان مجموعه داده واقعی در این مقاله مورد مطالعه قرار گرفته است. در بخش‌های بعدی مقاله، برآورد پارامترهای حاصل از این مجموعه داده به‌عنوان مقدار واقعی جامعه و به‌عنوان معیاری برای محاسبه توان دوم خطا در نظر گرفته شده‌اند. از سوی دیگر اطلاعات متغیرهای اصلی طرح هزینه و درآمد شامل درآمد، هزینه‌ی خوراکی و هزینه‌ی غیرخوراکی نیز در چارچوب نمونه‌گیری وجود ندارد و لازم است متغیر دیگری به‌عنوان متغیری کمکی انتخاب شود تا بتوان بر اساس آن پیشگویی و طبقه‌بندی مورد نیاز را انجام داد. برای رسیدن به این هدف لازم است طبقه‌بندی موضوعی بر اساس متغیری انجام شود که با متغیرهای اصلی طرح هزینه و درآمد خانوار یعنی هزینه‌ی کل، هزینه‌ی خوراکی، هزینه‌ی غیرخوراکی یا درآمد همبستگی داشته باشد، علاوه بر آن، اطلاعات این متغیر باید قابل دسترسی باشد. با بررسی اطلاعات موجود از آمارگیری هزینه و درآمد، ضریب همبستگی متغیر درآمد و متغیر اجاره‌بهای واحد مسکونی محل سکونت خانوار در حدود ۷۰٪ به دست آمد که در سطح اطمینان ۹۵ درصد، معنی‌دار است. بنابراین اجاره‌بهای واحد مسکونی محل سکونت خانوار می‌تواند به‌عنوان یک متغیر مناسب که اطلاعات آن از سایر منابع نیز در دسترس است، برای شبیه‌سازی چارچوب قابل استفاده باشد.

برای به‌کارگیری و بررسی کارایی روش طبقه‌بندی پیشنهادی دی‌گروتر و همکاران [۷] در آمارگیری هزینه و درآمد خانوار اطلاع از الگوی همبستگی فضایی متغیر اجاره‌بها به‌عنوان متغیر کمکی در دسترس، ضروری است. آماره موران متغیر لگاریتم اجاره‌بها در سطح اطمینان

محاسبه می‌شود. بنابراین تغییرپذیری کل از رابطه $O = \sum_{h=1}^H O(h)$ به دست می‌آید و طبق تخصیص نیمی اندازه بهینه نمونه هر طبقه با توجه به سهم آن از کل تغییرپذیری تابع O از رابطه $\tilde{n}_h = \tilde{n} \cdot \frac{N_h O_h}{\sum_{h=1}^H N_h O_h}$ تعیین می‌شود که با تقسیم کردن صورت و مخرج به N معادل با رابطه

$$\tilde{n}_h = \tilde{n} \cdot \frac{a_h O_h}{\sum_{h=1}^H a_h O_h} \quad (13)$$

است که \tilde{n} در آن اندازه نمونه از پیش تعیین شده است. به ازای هر اندازه نمونه $\tilde{n}_1, \dots, \tilde{n}_H$ تخصیص یافته به طبقات، واریانس نمونه‌گیری از رابطه زیر به دست می‌آید

$$\text{Var}(\hat{z}) = \frac{1}{N^2} \sum_{h=1}^H O_h^2 / n_h \quad (14)$$

کارایی این روش نمونه‌گیری از طریق مقایسه با واریانس نمونه‌گیری تصادفی ساده ارزیابی شده است. نتیجه‌ی مطالعات دی‌گروتر و همکاران [۷] نشان داده است روش پیشنهادی آن‌ها در مقایسه با روش ریشه‌ی فراوانی تجمعی و الگوریتم k - میانگین برای طبقه‌بندی فضایی از کارایی بالاتری برخوردار است.

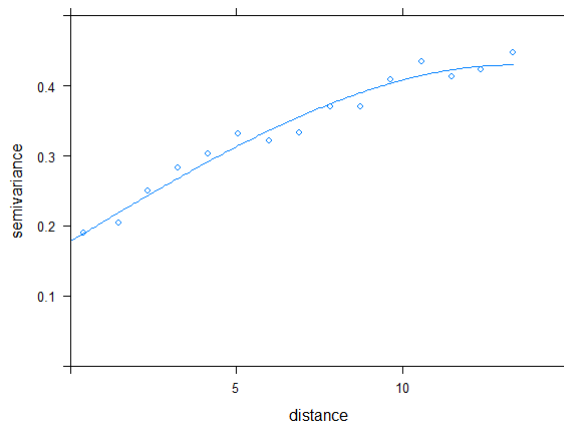
روش دی‌گروتر و همکاران [۷] فاصله‌ی تعمیم‌یافته بر اساس فاصله‌ی اقلیدسی و اطلاعات پیشگویی را در نظر می‌گیرد و طبقه‌بندی را به‌گونه‌ای انجام می‌دهد که واریانس برآوردگر هورویتر تامپسون کمینه شود، این یک روش طرح‌مبنا است که قابلیت بررسی برای داده‌های طرح هزینه و درآمد خانوار را دارد. بنابراین در ادامه‌ی این مقاله کاربست این روش برای داده‌های شبیه‌سازی شده با الگوی داده‌های هزینه و درآمد مورد بررسی قرار می‌گیرد.

همان‌طور که اشاره شد یکی از نکاتی که در به‌کارگیری روش *OSPAT* حائز اهمیت است تشخیص ساختار کوواریانس متغیر هدف است. بنابراین پیش از به‌کارگیری این روش بر داده‌های آمارگیری هزینه و درآمد خانوار لازم است ساختار کوواریانس متغیر هدف (به‌عنوان مثال درآمد خانوار) مشخص شود. با مدنظر قرار دادن این موضوع که اطلاعات متغیرهای هدف آمارگیری در چارچوب وجود ندارد، لازم است یک متغیر کمکی که با متغیر هدف همبستگی دارد و اطلاعات آن در چارچوب موجود یا از سایر منابع قابل حصول است، در نظر گرفته شود. میزان اجاره‌بهای واحد مسکونی محل سکونت خانوار می‌تواند به‌عنوان یک متغیر همبسته‌ی مناسب با درآمد در نظر گرفته شود.

در آمارگیری هزینه و درآمد میزان اجاره‌بها برای همه خانوارهای نمونه ثبت می‌شود. به این ترتیب که برای خانوارهای اجاره‌نشین میزان «اجاره‌بهای محل سکونت» و برای مالکین «برآورد اجاره‌بهای مسکن شخصی» ثبت می‌شود.

لگاریتم درآمد نیز مورد بررسی قرار گرفت. ضریب همبستگی موران برای متغیر لگاریتم درآمد 21% و در سطح 95% درصد، معنی‌دار است. برازش تغییر نگرهای مختلف برای این متغیر نیز نشان داد که تغییرنگار نمایی بهترین برازش را دارد و مقادیر واریانس و دامنه به ترتیب برابر با 23% و $6/22$ برآورد شد.

95% درصد با مقدار 27% معنی‌دار است. برازش تغییر نگرهای مختلف گاوسی، توانی، ماترون، کروی و نمایی نشان داد که تغییرنگار نمایی با مجموع توان‌های دوم خطا برابر با 1% ، بهترین برازش را دارد و مقدار پارامترهای واریانس و دامنه به ترتیب برابر با $45/0$ و $75/13$ برآورد شد. با توجه به همبستگی درآمد و اجاره‌بها، الگوی همبستگی



شکل ۱. واریوگرام تجربی اجاره‌بهای واحد مسکونی

ساده محاسبه می‌شود. مراحل بالا در ادامه به تفصیل ارائه شده‌اند.

گام اول: شبیه‌سازی چارچوب

فرض کنید $\mathbf{z} = (z(s_1), \dots, z(s_n))^T$ تحقق متغیرهای تصادفی $\mathbf{Z} = (Z(s_1), \dots, Z(s_n))^T$ در n موقعیت مجزای s_1, \dots, s_n باشند. برای سادگی در نوشتار فرض کنید z_i و Z_i به ترتیب نمایشگر $z(s_i)$ و $Z(s_i)$ باشند. فرض می‌کنیم بردار \mathbf{Z} متغیرهای پاسخ گاوسی باشند به طوری که $\mathbf{Z} | \beta, \sigma^2, \eta \sim N(\mu, \Sigma)$ و بردار تصادفی $\eta(s) = (\eta(s_1), \dots, \eta(s_n))$ نمایشگر میدان تصادفی گاوسی 5 (GRF) باشد. در ادامه $\eta(s_i)$ با η_i ، $i = 1, \dots, n$ نمایش داده می‌شود. مدل رگرسیون فضایی به صورت زیر تعریف می‌شود

$$\mu_i = \beta \mathbf{x}_i + \eta_i, \quad i = 1, \dots, n. \quad (15)$$

که در آن $\beta = (\beta_0, \beta_1, \dots, \beta_{(p-1)})$ و $\mathbf{x}_i = (x_{i0}, \dots, x_{i(p-1)})^T$ ترتیب بردار متغیرهای کمکی و ضرایب رگرسیون هستند. همان‌طور که قبلاً اشاره شد، درآمد یکی از متغیرهای اصلی آمارگیری هزینه و درآمد خانوار است که اطلاعات آن در چارچوب آمارگیری وجود ندارد. بنابراین برای کاربست نمونه‌گیری فضایی در شرایط واقعی لازم است تغییرنگار متغیر اجاره‌بها مورد توجه قرار بگیرد، بنابراین این نکته در شبیه‌سازی چارچوب نیز مدنظر قرار گرفت.

لگاریتم اجاره‌بها از توزیع نرمال پیروی می‌کند و با لگاریتم مساحت

با وجود اینکه مقادیر برآورد شده پارامترها برای متغیرهای لگاریتم درآمد و لگاریتم اجاره‌بها دقیقاً یکسان نیست، با در نظر گرفتن این نکته که هر دو از ساختار همبستگی نمایی پیروی می‌کنند و با توجه به همبستگی این دو متغیر، می‌توان در غیاب اطلاعات درآمد از اجاره‌بها به‌عنوان متغیر کمکی استفاده کرد (کلهری، [۹]).

چالش دیگری که در حال حاضر برای کاربست طبقه‌بندی فضایی بر آمارگیری هزینه و درآمد وجود دارد، نبود اطلاعات مکانی در چارچوب آمارگیری است. لذا برای ارزیابی کارایی روش فعلی آمارگیری هزینه و درآمد خانوار، لازم است چارچوب بر اساس اطلاعات موجود و مشابه الگوی داده‌های طرح هزینه و درآمد شبیه‌سازی شود. بنابراین در گام اول، بر اساس مقادیر برآورد شده پارامترهای تغییرنگار بر لگاریتم اجاره‌بها، چارچوب نمونه‌گیری شبیه‌سازی می‌شود. در گام دوم با در دست داشتن الگوی همبستگی فضایی و چارچوب، پیش‌گویی فضایی در موقعیت‌های جدید با کریجیدن انجام می‌شود که یکی از اطلاعات مورد نیاز برای محاسبه‌ی فاصله‌ی تعمیم‌یافته و اجرای طبقه‌بندی بهینه است. در گام سوم طبقه‌بندی بهینه فضایی انجام می‌شود. در گام چهارم پس از تعیین اندازه‌ی نمونه، واحدهای نمونه در هر طبقه به روش تصادفی ساده انتخاب می‌شوند. در گام آخر، برآورد پارامترها به همراه ارزیابی نسبی و کارایی روش نمونه‌گیری نسبت به نمونه‌گیری تصادفی

⁵Gaussian Random Field

پیشگویی بر اساس صفحه‌ی 40×40 با پرش‌های چهارتایی در نظر گرفته می‌شود. در واقع اطلاعات یک چهارم از موقعیت‌های مشبکه حفظ می‌شوند و سپس در کل موقعیت‌ها پیشگویی با استفاده از کریگیدن عادی (کرس، [۵]) انجام می‌شود و مقادیر \tilde{z}_i برای $i = 1, \dots, n$ تولید می‌شود.

بنابراین در پایان گام دوم مقادیر واقعی z_i و مقادیر پیشگویی شده \tilde{z}_i و واریانس خطای پیشگویی $\text{Var}(e_i)$ در هر یک از موقعیت‌های s_1, \dots, s_n در دسترس هستند که مقادیر مورد نیاز برای انجام طبقه‌بندی فضایی OSPAT هستند.

گام سوم: طبقه‌بندی فضایی

در این مرحله الگوی طبقه‌بندی که در بخش ۲ معرفی شد در نرم‌افزار R کدنویسی شده و طبقه‌بندی داده‌های تولید شده در گام اول با استفاده از اطلاعات کمکی گام دوم، که در واقع به منزله طبقه‌بندی چارچوب نمونه‌گیری است، انجام می‌شود.

شکل ۲ داده‌های واقعی و مقادیر پیشگویی تولید شده در شبیه‌سازی را به همراه طبقه‌بندی نهایی در یک مشبکه با اندازه ۱۶۰۰ نمایش می‌دهد. همان‌طور که ملاحظه می‌شود پیش‌گویی میدان فضایی با الگوی داده‌های واقعی مطابقت دارد و همچنین طبقه‌بندی فضایی ناحیه‌ی مورد مطالعه با این دو الگو مشابهت دارد.

گام چهارم: استخراج نقاط نمونه از چارچوب

اندازه مشبکه شبیه‌سازی شده با الگوی داده‌های آمارگیری هزینه و درآمد خانوار ۱۶۰۰ است و اندازه طبقات پس از انجام طبقه‌بندی بهینه فضایی برابر با ۷۵۵، ۳۶۸ و ۴۷۷ به دست آمد. اندازه‌ی نمونه در هر طبقه با رابطه ۱۳ محاسبه و مقادیر آن به ترتیب برابر با ۵۲، ۱۳ و ۲۳ به دست آمد. واحدهای نمونه در هر طبقه به شیوه‌ی تصادفی ساده انتخاب شدند. وزن نمونه‌گیری در هر طبقه برابر با عکس احتمال انتخاب است.

گام پنجم: برآورد پارامترها

در روش فعلی آمارگیری هزینه و درآمد خانوار، انتخاب نمونه‌ها به صورت یک روش سه مرحله‌ای انجام می‌شود. واحد نمونه‌گیری مرحله‌ی اول یک خوشه از بین خوشه‌های نمونه‌ی پایه بوده که به صورت تصادفی انتخاب می‌شود. منظور از نمونه پایه، نمونه‌ای است که می‌توان از آن برای تأمین نیازهای چند آمارگیری یا چند دوره از یک آمارگیری، زیر نمونه‌هایی انتخاب کرد.

واحد مسکونی همبستگی دارد. بنابراین برای شبیه‌سازی چارچوب به‌گونه‌ای که در بردارنده اطلاع متغیر اجاره‌ها باشد، لازم است در ابتدا مدل مناسب بر داده‌های اجاره‌ها برازش شود و سپس این مدل برای شبیه‌سازی مقادیر این متغیر در چارچوب به کار گرفته شود. لذا در گام اول مدل رگرسیونی برای روند زدوده کردن لگاریتم اجاره‌ها بر داده‌های آمارگیری هزینه و درآمد خانوار برازش می‌شود و پارامترهای مدل با رهیافت ماکسیمم درستنمایی برآورد می‌شوند. بنابراین برای مدل $z_i = \beta_0 x_{i0} + \beta_1 x_{i1} + \eta_i$ که در آن لگاریتم اجاره‌ها و x_{i1} لگاریتم مساحت واحد مسکونی است، $(\hat{\beta}_0, \hat{\beta}_1) = (13/18, 0/48)$ به دست می‌آید.

فرض کنید اثرات فضایی η_i تحقق یک میدان تصادفی گاوسی با میانگین صفر و ماتریس کوواریانس Σ_η هستند که مؤلفه‌های ماتریس بر اساس فاصله موقعیت‌های جغرافیایی تعیین می‌شوند (ونگ و وال، [۱۲]). همان‌طور که پیش‌تر اشاره شد، تابع کوواریانس نمای بهترین برازش را بر داده‌های لگاریتم اجاره‌ها دارد. بنابراین تابع کوواریانس به صورت $\sigma_{kl}^2 = \gamma \exp(-|s_k - s_l|/\phi)$ تعریف می‌شود که در آن γ واریانس فضایی، ϕ دامنه و $|s_k - s_l|$ فاصله واحد نمونه k و l است. تابع درستنمایی مدل ۱۵ به صورت زیر تعریف می‌شود:

$$L(\beta, \phi, \gamma | z) = \int \prod_{i=1}^n f(z_i | \mathbf{x}, \beta, \phi, \gamma) \mathbf{f}(\mathbf{x} | \phi, \gamma) \mathbf{d}\mathbf{x}.$$

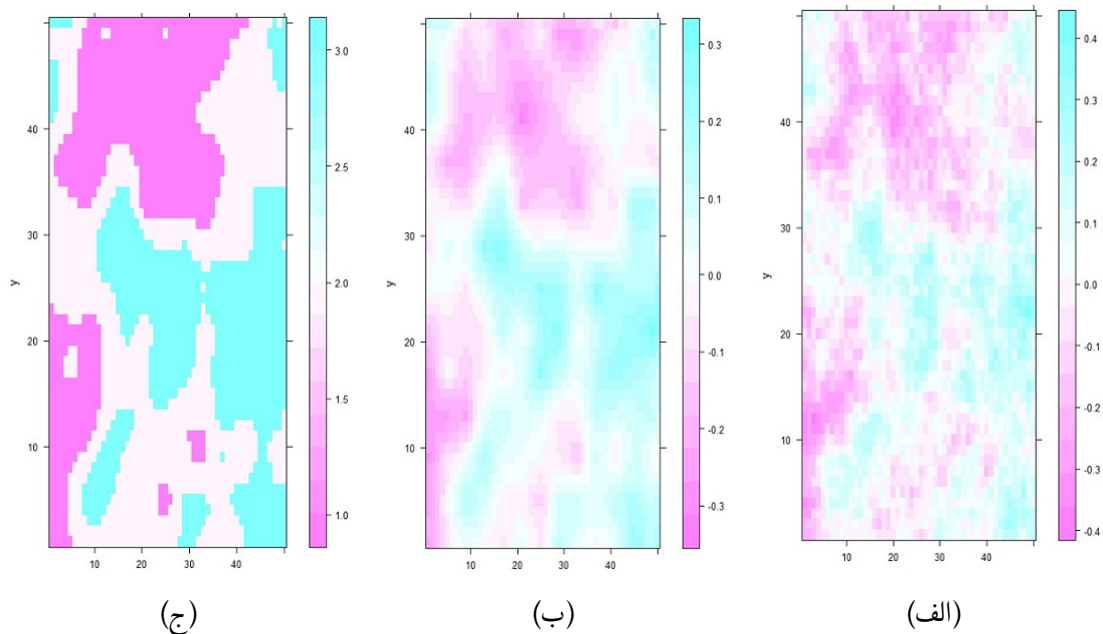
برای تولید مقادیر اثر تصادفی η_i که از توزیع $N(0, \Sigma_\eta)$ پیروی می‌کنند، مؤلفه ij ام ماتریس Σ_η به صورت $(\Sigma_\eta)_{ij} = \gamma \exp(-|s_i - s_j|/\phi)$ تعیین می‌شوند و مقادیر پارامترهای شبیه‌سازی بر اساس برآورد پارامترهای تغییرنگار تجربی برابر با $\gamma = 0/45$ و $\phi = 14$ در نظر گرفته می‌شوند. مقادیر متغیر پنهان فضایی با استفاده از تجزیه چولسکی بر مشبکه‌ای با اندازه 40×40 تولید می‌شوند.

سپس مجموع مقادیر تولید شده برای متغیر پنهان فضایی و مقادیر تولید شده برای مؤلفه ثابت مدل $(\hat{\beta}_0 x_{i0} + \hat{\beta}_1 x_{i1})$ به عنوان لگاریتم اجاره‌ها در نظر گرفته شده و تبدیل وارون آن با تابع نمای برای مقادیر تولید شده اجاره‌ها در چارچوب قرار می‌گیرد.

گام دوم: کریگیدن عادی^۶

به منظور استفاده از روش OSPAT لازم است مقادیر میدان فضایی پیش‌گویی شوند و بر اساس مقادیر به دست آمده از کریگینگ فضایی، طبقه‌بندی انجام شود. به همین منظور لازم است کریگیدن بر اساس مقادیر واقعی متغیر پنهان صورت پذیرد. فرض می‌کنیم مقادیر تولید شده در گام اول، شبیه‌سازی مقادیر واقعی متغیر پنهان باشند و مشبکه

^۶Ordinary Kriging



شکل ۲. الف- داده‌های شبیه‌سازی ب- مقادیر پیشگویی، ج- پهنه‌بندی فضایی طبقات

با شیوه‌ی فعلی آمارگیری هزینه و درآمد، مقدار کارایی نسبی در برآورد میانگین درآمد محاسبه و مقدار آن برابر با $23/1$ محاسبه شد. با وجود اینکه شبیه‌سازی چارچوب بر اساس الگوی همبستگی متغیر اجاره‌بها انجام شده است، ولی کارایی روش طبقه‌بندی بهینه فضایی از نمونه‌گیری تصادفی ساده در سطح تقسیم‌بندی جغرافیایی، بیشتر است. بنابراین استفاده از طبقه‌بندی فضایی بهینه در خوشه‌های نمونه زمانی که همبستگی فضایی وجود دارد، منجر به بهبود برآوردها در آمارگیری هزینه و درآمد خانوار می‌شود.

با توجه به اینکه ممکن است میزان همبستگی فضایی در نواحی مختلف، به‌عنوان مثال استان‌های متفاوت کشور، تغییر کند یا اندازه زیرجامعه‌ها در نواحی مختلف با یکدیگر تفاوت داشته باشد، به‌منظور ارزیابی روش پیشنهادی گام‌های یک تا پنج که در بالا تشریح شد در یک مطالعه شبیه‌سازی برای شبکه‌هایی با اندازه‌های متفاوت و میزان مختلف همبستگی موران تکرار شده است. روش به کار گرفته شده برای نمونه‌گیری فضایی در ۱۰۰ تکرار شبیه‌سازی برای هر یک از شبکه‌ها با اندازه‌های ۱۰۰، ۴۰۰، ۹۰۰ و ۱۶۰۰ به ازای مقادیر مختلف ضریب همبستگی موران انجام شده است. در هر تکرار شبیه‌سازی میانگین توان دوم خطای برآورد پارامتر جامعه به روش نمونه‌گیری تصادفی ساده و نمونه‌گیری فضایی به دست آمده است که نسبت آن‌ها به‌عنوان معیار کارایی در نظر گرفته شده و مقادیر آن در جدول ۱ نمایش داده شده است.

هر خوشه‌ی نمونه شامل یک بلوک/آبادی، بخشی از یک بلوک/آبادی بزرگ یا در مواردی مجموعه‌ای از چند بلوک/آبادی کوچک است که بر اساس اطلاعات آخرین سرشماری عمومی نفوس و مسکن ساخته شده است. در مرحله‌ی دوم در هر خوشه‌ی نمونه، گروه‌های چرخش شش خانواری ساخته شده و سه گروه چرخش بر اساس الگوی چرخش، برای آمارگیری در هر سال تعیین شده‌اند. در مرحله‌ی سوم، داخل هر یک از گروه‌های چرخش، دو خانوار به روش نمونه‌گیری تصادفی ساده به‌عنوان خانوار نمونه‌ی اصلی انتخاب می‌شوند. در این بررسی به دلیل فقدان اطلاعات موقعیت‌های جغرافیایی واحدهای نمونه در نمونه پایه، امکان استفاده از چارچوب نمونه پایه وجود نداشت. تفاوت شیوه نمونه‌گیری فضایی پیشنهادی و روش فعلی نمونه‌گیری در آمارگیری هزینه و درآمد خانوار در این است که در روش فعلی، خانوارهای نمونه در هر خوشه به‌صورت تصادفی ساده انتخاب می‌شوند اما در نمونه‌گیری فضایی، خوشه انتخابی باهدف دستیابی به واحدهای نمونه که شباهت کمتری باهم داشته باشند به‌صورت فضایی طبقه‌بندی شده و واحدهای نمونه در هر طبقه به‌صورت تصادفی ساده انتخاب می‌شوند و برای نزدیک بودن طرح نمونه‌گیری پیشنهادی به طرح کنونی مرکز آمار ایران، از برآوردگر هورویتز تامپسون برای برآورد پارامترها استفاده شده است. در این گام پارامترهای جامعه با وزن‌دهی برآورد می‌شوند. وزن هر واحد نمونه‌گیری عکس احتمال انتخاب واحد در طبقه است.

برای ارزیابی کارایی نمونه‌گیری فضایی در برآورد پارامترها در مقایسه

جدول ۱. کارایی نمونه‌گیری فضایی نسبت به تصادفی ساده برای شبکه‌ها و میزان همبستگی فضایی مختلف در ۱۰۰ تکرار شبیه‌سازی

اندازه‌ی شبکه	میزان همبستگی موران	کارایی نسبی
۱۰۰	۱۲/۰	۱۵/۱
	۲۵/۰	۲۵/۱
	۴۰/۰	۶۱/۱
۴۰۰	۱۲/۰	۲۲/۱
	۲۵/۰	۲۷/۱
	۴۰/۰	۶۳/۱
۹۰۰	۱۲/۰	۲۳/۱
	۲۵/۰	۷۵/۱
	۴۰/۰	۷۷/۱
۱۶۰۰	۱۲/۰	۴۱/۱
	۲۵/۰	۸۶/۱
	۴۰/۰	۲۸/۲

آمارگیری‌های مرکز آمار ایران به شمار می‌رود به موضوع نمونه‌گیری فضایی پرداخته شده است. از میان روش‌های موجود طبقه‌بندی ناحیه‌ی مورد مطالعه بر اساس فاصله‌ی تعمیم‌یافته‌ای که بر مبنای فاصله‌ی واقعی و واریانس پیشگویی محاسبه می‌شود، برخی مشکلات موجود را مرتفع می‌کند و قابلیت بالاتری در اجرای عملیاتی نمونه‌گیری فضایی دارد. نتایج مطالعه‌ی داده‌های واقعی و شبیه‌سازی نشان می‌دهد که طبقه‌بندی ناحیه‌ی مورد مطالعه بر اساس فاصله‌ی تعمیم‌یافته می‌تواند در کاهش خطای برآوردها مؤثر باشد، اما به‌کارگیری دقیق‌تر این روش در عمل، نیازمند دسترسی به اطلاعات موقعیت مکانی واحدهای چارچوب است. این امر مستلزم تهیه سازوکار مناسب برای تولید داده‌های مکان محور است. نمونه‌گیری فضایی پیشنهادی فقط برای طبقه‌بندی خوشه‌ها ارائه شده است. در صورت وجود اطلاعات مکانی و همچنین اطلاعات یک متغیر همبسته با متغیر پاسخ که به صورت تجربی ثابت شود الگوی همبستگی فضایی مشابه با متغیر پاسخ دارد، امکان استفاده از نمونه‌گیری فضایی وجود دارد. مطالعات بیشتر در پیشبرد نمونه‌گیری فضایی طرح‌مبنا از الزاماتی است که پیش از استفاده نمونه‌گیری فضایی در یک آمارگیری ملی باید مورد توجه ویژه قرار بگیرد و هر روش جدید پیش از استفاده در سطح کل کشور نیازمند آزمایش و آسیب‌شناسی است. همچنین باید این نکته را مدنظر قرار داد که وجود تصاویر ماهواره‌ای از پهنه آمارگیری و در دست داشتن اطلاعات کمکی در بهبود روش‌ها کمک به سزایی خواهد کرد.

ملاحظه می‌شود با افزایش اندازه‌ی شبکه کارایی نسبی نمونه‌ی فضایی نسبت به نمونه‌گیری تصادفی ساده افزایش می‌یابد و علاوه بر آن در همه‌ی شبکه‌ها با افزایش میزان همبستگی فضایی، کارایی روش نمونه‌گیری فضایی افزایش می‌یابد.

۴ بحث و نتیجه‌گیری

در بسیاری از پدیده‌های طبیعی یا بشری همبستگی فضایی میان پاسخ‌ها وجود دارد و استفاده از روش‌های کلاسیک آماری به دلیل لحاظ نکردن این همبستگی برای بررسی این پدیده‌ها مناسب نیست. با وجود اینکه در سال‌های اخیر علاقه‌مندی به استفاده از آمار فضایی گسترش یافته است ولی استفاده عمده از آمار فضایی در مطالعات زیست‌محیطی بوده است و مشکلاتی در زمینه استفاده از آمار فضایی در سایر زمینه‌ها وجود دارد. یکی از مهم‌ترین کاستی‌ها در این زمینه، نبود یک چارچوب نمونه‌گیری کامل در آمارگیری‌های نمونه‌ای است و چگونگی نمایاندن جمعیت (به صورت نقطه‌ای، ناحیه‌ای)، اندازه و ترکیب چارچوب را بسیار تحت تأثیر قرار می‌دهد.

به دلیل موانع و مشکلاتی که به آن‌ها اشاره شد نمونه‌گیری فضایی در مباحث مربوط به آمار رسمی کمتر مورد توجه قرار گرفته است. در این مقاله با تمرکز بر آمارگیری هزینه و درآمد خانوار که یکی از مهم‌ترین

مراجع

- [۱] محمدزاده، م. (۱۳۹۸). آمار فضایی و کاربردهای آن، چاپ سوم، مرکز نشر آثار علمی دانشگاه تربیت مدرس، تهران.
- [2] Abi, N. (2019). *Spatially Balanced Sampling Methods in Household Surveys*. Ph.D Thesis, University of Canterbury.
- [3] Baillargeon, S. and Rivest, L.P. (2009). A General Algorithm for Univariate Stratification, *International Statistical Review*. **77**, 331-344.
- [4] Ballin, M., and Barcaroli, G. (2013). Joint Determination of Optimal Stratification and Sample Allocation Using Genetic Algorithm, *Survey Methodology*. **39**, 369-393.
- [5] Cressie, N. (1993). *Statistics for Spatial Data*. John Wiley and Sons, New York.
- [6] Dalenius, T. and Hodges, J.L. (1959). Minimum Variance Stratification, *Journal of the American Statistical Association*. **54**, 88-101.
- [7] De Gruijter, J.J., Minasny, B. and McBratney, A. B. (2015). Optimizing Stratification and Allocation for Design-Based Estimation of Spatial Means using Predictions with Error, *Journal of Survey Statistics and Methodology*. **3(1)**, 19-42.
- [8] Horgan, J.M. (2010). Choosing the Stratification Boundaries: The Elusive Optima, *Istanbul University Journal of the School of Business Administration*. **39**, 195-204.
- [9] Kalhori Nadrabadi, L. (2019). Utilization of Pearson Correlation to Grab Spatial Correlation in the of Absence of Information, *Proceedings of the 3rd Seminar on Spatial Statistics and Its Applications*. Zanjan, Iran, 69-77.
- [10] Khavarzadeh, R., Mohammadzadeh, M. and Mateu, J. (2018). A Simple Two-Step Method for Spatio-Temporal Design-Based Balanced Sampling, *Stochastic Environmental Research and Risk Assessment*. **32(2)**, 457-468.
- [11] Kozak, M. (2004). Optimal Stratification using Random Search Method in Agricultural Surveys, *Statistics in Transition*. **6**, 797-806.
- [12] Wang, F. and Wall, M. M. (2003). Generalized Common Spatial Factor Model, *Biostatistics*. **4**, 569-582.

Application of Optimal Spatial Stratification in Household Income and Expenditure Survey to Provide Estimates by Spatial Design-Based Sampling

Lida Kalhori Nadrabadi¹, Roshanak Aliakbari Saba² and Asiyeh Abbasi³

Abstract:

Household Income and Expenditure Survey (HEIS) is one of the most important surveys of the Statistical Center of Iran, the main parameters of which are spatially correlated. When there is a spatial correlation between the units of population, the classical way of selecting independent sampling units is challenging due to the lack of basic condition for the independence. Using spatial sampling is a solution to encounter this problem. Implementation of spatial sampling has received less attention in official statistics due to the lack of access to a suitable framework. In this paper we review a design-based model assisted method for optimal spatial stratification of the target population. At present, spatial information of population units are not available in the framework of HEIS, but access to spatial information of some sample units has been achieved by the Statistical Center of Iran for this study. The production of spatial data is one of the main components in the modernization of the statistical system which is considered by Statistical Center of Iran. In this paper, the sampling frame is simulated based on the HEIS data and then application of optimal spatial stratification based on a generalized distance is performed. The results demonstrate an increase in the efficiency of the mentioned sampling method compared to simple random sampling at the level of geographical areas. Also, simulation of grids with different sizes and correlations reflects the better performance of this method compared to the current method of HEIS.

Keywords: Spatial Stratification, Generalized Distance, Household Expenditure and Income Survey.

¹Statistical and Training Center, Tehran , Iran. kalhori@srtc.ac.ir

²Statistical and Training Center, Tehran , Iran. r_saba@srtc.ac.ir.

³ Statistical Center of Iran, Tehran , Iran. asieh_abasi@yahoo.com