

تعیین طبقات در نمونه‌گیریهای چندمنظوره یا طبقه‌بندی شده

عباس گرامی^۱ محمدباقر سخاوت^۲ محمدباقر حقیقی انارکی^۳

چکیده

سودبخشی نمونه‌گیری طبقه‌بندی شده در مقایسه با نمونه‌گیری تصادفی ساده، در صورتی بیشتر است که واحدهای آماری درون هر طبقه دارای همگنی بسیار زیاد و واحدهای مربوط به طبقات متفاوت از همگنی کمتری برخوردار باشند. در ادبیات مربوطه، وقتی که وسیله تعیین مرز طبقات تنها یک متغیر کمکی باشد، مطالبی ارائه شده است که علاقه‌مندان را به [۱] رجوع می‌دهیم. در این مقاله، شرایطی بررسی شده است که به جای استفاده از یک متغیر کمکی، به روش‌های مختلف چندین متغیر به کار برد شده‌اند و سپس در یک مطالعه موردنی، قابلیت این روش‌های پیشنهادی مورد مقایسه قرار گرفته‌اند و نشان داده شده است که روش‌های مذکور در مقایسه با روش‌های استاندارد، از کارآیی بالاتری برخوردارند. حتی در شرایطی که بخواهیم تنها یک متغیر کمکی را وسیله تعیین مرز طبقات قرار دهیم، روش‌های پیشنهادی دارای مزیت بیشتری می‌باشند.

واژه‌های کلیدی: نمونه‌گیری چندمنظوره، نمونه‌گیری طبقه‌بندی شده، تحلیل خوشای، تحلیل مؤلفه‌های اصلی، ضریب تغییرات و معیارهای مقایسه.

۱. مقدمه

نمونه‌گیریها همواره کوشش می‌شود از طرح‌هایی استفاده شود که با

حجم نمونه معین، دقت بیشتری به همراه داشته باشند.

از جمله طرح‌های بسیار معمول و مورد استفاده در نمونه‌گیریها من جمله نمونه‌گیریهای کشاورزی، نمونه‌گیری طبقه‌بندی شده است که در آن یک متغیر کمکی یا ترکیب خطی ساده‌ای از چند متغیر کمکی (مثلًاً مجموع) به عنوان معیار طبقه‌بندی و تعیین طبقات مورد استفاده قرار می‌گیرد. کارآیی و قابلیت این کار، زمانی است که متغیرهای کمکی مورد استفاده در طبقه‌بندی با متغیرهای مورد بررسی

سالیانه مقادیر زیادی از نیرو و امکانات کشور صرف تهیه و تولید آمارهای مورد نیاز دستگاههای برنامه‌ریزی، اجرایی و نیز مدیریت جامعه در بخش‌های دولتی و خصوصی می‌گردد. اکثریت آمارهای مورد نیاز از طریق تهیه و اجرای طرح‌های نمونه‌گیری حاصل می‌شود و دقت آمارهای به دست آمده از این طریق، به عوامل زیادی از جمله طرح مورد استفاده بستگی دارد. مسلماً طرح‌های با حجم نمونه نسبتاً کم و دقت بالا در صرفه‌جویی امکانات تأثیر بسزایی دارد و بنابراین در

^۱ عضو هیئت علمی پژوهشکده آمار

^۲ آمارشناسان مرکز آمار ایران

یک متغیر کمکی، منطقی به نظر می‌رسد ولی این سوال مطرح می‌شود که اگر قرار باشد تنها از یک متغیر کمکی برای تعیین طبقات استفاده شود آیا مجموع ساده سطح زیر کشت در آخرین سرشماری می‌تواند بهترین متغیر کمکی باشد؟ بدیهی است که استفاده از تنها یک متغیر کمکی برای ساختن طبقات در نمونه‌گیریهای چند منظوره نمی‌تواند کارآیی بالایی به همراه داشته باشد. بنابراین تعداد مطلوب متغیر کمکی و نحوه تعیین آنها سؤال دیگری است که مطرح می‌شود. استفاده از مقادیر همه متغیرهای مورد بررسی براساس اطلاعات آخرین سرشماری به عنوان متغیرهای کمکی یک انتخاب معقول به نظر می‌رسد. با این وجود ممکن است حجم زیاد محاسبات در تعیین طبقات، کاربرد آن را مشکل و یا حتی غیر ممکن سازد. در این صورت باید تعداد کمتری متغیر کمکی را مورد توجه قرار داد. در این بخش پیشنهاداتی در مورد تعیین متغیر یا متغیرهای کمکی برای ساختن طبقات به شرح زیر ارائه می‌شود.

- انتخاب متغیری از بین متغیرهای کمکی ممکن که دارای بیشترین تغییرات است.
- استفاده از تحلیل مؤلفه‌های اصلی^۴، PCA، برای ساختن متغیرهای کمکی.

۱.۲ انتخاب متغیر کمکی با بیشترین تغییرات

چنانچه یک جامعه آماری را که برای هر واحد آن، اطلاعات p متغیر کمکی در اختیار داریم بخواهیم با استفاده از یک متغیر طبقه‌بندی کنیم، انتظار داریم که متغیر انتخابی، آن متغیری باشد که بیشترین واریانس یا ضریب تغییرات را داشته باشد، گرچه استفاده از این متغیر برای طبقه‌بندی لزوماً منجر به یک طرح با خواص بهینه نخواهد بود. در این ارتباط می‌توان به این نکته نیز اشاره کرد که اگر بخواهیم یک جامعه آماری را که برای هر واحد آن m متغیر کمکی در دست است به کمک تعداد کمتری متغیر، مثلاً $m < p$ (متغیر، طبقه‌بندی کنیم)، در این صورت نیز مناسب به نظر می‌رسد تا از m متغیری که بیشترین واریانس یا بیشترین ضریب تغییرات را دارند استفاده کنیم.

در نمونه‌گیری، دارای همبستگی بالا باشند. اما ممکن است در عمل چنین وضعیتی وجود نداشته باشد. بنابراین لازم است روش‌های دیگری مورد استفاده قرار گیرد که در صورت برقرار نبودن وضعیت مورد بحث، از کارآیی بیشتری برخوردار باشند.

در نمونه‌گیریهای طبقه‌بندی شده، برای تعیین طبقات، تعداد نمونه، انتساب نمونه‌ها به طبقات و انتخاب واحدهای نمونه، موارد زیر باید بررسی شوند

الف - پیدا کردن یک یا چند متغیر کمکی مناسب برای طبقه‌بندی واحدهای جامعه

ب - روش ساختن طبقات

ج - تعیین تعداد طبقات

د - تعیین تعداد نمونه

ه - روش انتساب نمونه به طبقات

در مقاله حاضر بررسی بندهای الف و ب از بندهای فوق یعنی تعیین طبقات مورد توجه است. در این راستا ابتدا نحوه تعیین متغیر یا متغیرهای کمکی برای طبقه‌بندی واحدهای جامعه، مورد بحث قرار می‌گیرند و سپس روش‌هایی برای ساختن طبقات پیشنهاد می‌شود. در انتهای با معرفی چند معیار، روش‌های کلاسیک موجود برای ساختن طبقات با روش‌های پیشنهادی مقایسه می‌شوند.

۲. متغیرهای کمکی مناسب برای ساختن طبقات

نمونه‌گیریها اغلب چند منظوره هستند، بدین مفهوم که برآورده پارامترهای مربوط به چندین متغیر یا چندین پارامتر از یک متغیر از جامعه هدف و یا تلفیقی از این دو نوع را مورد نظر دارند. در صورتی که از طرح نمونه‌گیری طبقه‌بندی شده استفاده شود، موضوع مهم، تعیین متغیر یا متغیرهای کمکی است که بر اساس آن طبقات تعریف می‌شوند. بدیهی است که باید منطبقاً بین متغیرهای کمکی با متغیرهای مورد نظر در طرح، نوعی وابستگی وجود داشته باشد، تا طبقات تشکیل شده مو جب افزایش کارآیی برآوردها شود. در بعضی نمونه‌گیریها مثلاً نمونه‌گیریهای کشاورزی، استفاده از مجموع سطح زیر کشت محصولات مورد آمارگیری در آخرین سرشماری کشاورزی به عنوان

مورد استفاده قرار داد. به علاوه حتی زمانی که بخواهیم بر اساس تنها یک ترکیب خطی، طبقه‌بندی را انجام دهیم، باز هم استفاده از مهمنترین مؤلفه اصلی (اولین مؤلفه اصلی)^۷ در مقایسه با ترکیب خطی ساده مجموع صفات، اطلاعات بیشتری را انتقال خواهد داد. برای توضیح بیشتر پیرامون این تحلیل، به منبع [۵] مراجعه شود.

۳. معرفی روش‌های ساختن طبقات

روش‌های ساختن طبقات که در منابع نمونه‌گیری مورد بحث قرار می‌گیرند، عمدها بر مبنای یک متغیر کمکی است، مثل سطح زیر کشت محصولات کشاورزی در آخرین سرشماری کشاورزی که قبل از توضیح داده شد. استفاده از تنها یک متغیر کمکی برای ساختن طبقات از مزایای سادگی در محاسبات و تعیین (حدود) طبقات برخوردار است. اما مسلماً نمی‌توان گفت که این متغیر به تنها بخوبی طبقه‌بندی را بدست می‌دهد. در بخش قبل متغیر یا متغیرهای کمکی که می‌توانند برای طبقه‌بندی استفاده شوند، مورد بحث قرار گرفت. در صورتی که تنها یک متغیر کمکی برای طبقه‌بندی موجود باشد یا تنها یک متغیر مناسب باشد، روش کلاسیک دالینیوس^۸ [۱] مورد استفاده قرار می‌گیرد. گرچه سادگی این روش و مناسب بودن آن در شرایط خاص، از مزایای آن محسوب می‌شود، اما سؤال مهم این است که آیا این روش، بهترین طبقه‌بندی ممکن را بدست می‌دهد؟ گرچه پاسخ روشی برای این سؤال وجود ندارد، اما امکان استفاده از این روش برای شرایط باشند از یک متغیر کمکی، وجود ندارد. در زیر به معرفی بعضی تکنیکهای تحلیل خوشه‌ای که امکان طبقه‌بندی جامعه بر اساس بیش از یک متغیر کمکی را نیز فراهم می‌کند پرداخته می‌شود.

۴. روش‌های ساختن طبقات با استفاده از تحلیل خوشه‌ای

هدف از تحلیل خوشه‌ای عبارت از تقسیم n فرد بین k خوشه است، به گونه‌ای که بر اساس معیاری که تعریف می‌شود، افراد داخل

۲-۲ تحلیل مؤلفه‌های اصلی و استفاده از آن در تعیین متغیرهای کمکی مناسب

هدف اصلی تحلیل مؤلفه‌های اصلی، کاهش حجم داده‌های چند متغیرهای است که دارایوابستگی‌های خطی درونی قابل توجهی هستند، به طوری که داده‌های کاهش یافته بیشترین اطلاعات داده‌های اصلی را منتقل کنند. این روش مبتنی بر ویژه بردارهای^۹ ماتریس واریانس - کوواریانس متغیرها یا ماتریس ضربه همبستگی مربوط به آنهاست. چنانچه در این تحلیل، مشاهدات به شکل اصلی خود مورد استفاده قرار گیرد، از ماتریس واریانس - کوواریانس و در حالتی که این مشاهدات به صورت استاندارد درآیند از ماتریس ضربه همبستگی استفاده می‌شود.

شایان ذکر است که تحلیل مؤلفه‌های اصلی، معمولاً گام آخر تجزیه و تحلیل داده‌ها را تشکیل نمی‌دهد، بلکه این تحلیل ابزاری جهت انجام تجزیه و تحلیل‌های بعدی نیز ارائه می‌کند که بعداً به رابطه آن با تحلیل خوشه‌ای^{۱۰} اشاره می‌کنیم. در صورتی که حجم داده‌های موجود از نظر تعداد متغیر بسیار زیاد باشد و بخواهیم از تحلیل خوشه‌ای استفاده کنیم، ممکن است ابزارهای محاسباتی در دسترس، قابلیت انجام محاسبات روی حجم بالی داده‌ها را نداشته باشند. در چنین حالتی به جای استفاده از کلیه متغیرها، از تعداد اندکی از مؤلفه‌های اصلی مهتم استفاده می‌شود که امکانات محاسباتی آن در اختیار باشد.

در پاره‌ای از نمونه‌گیریها، مخصوصاً نمونه‌گیریهای کشاورزی، از یک ترکیب خطی، مثلاً مجموع سطح کشت تمامی محصولات هر آبادی، به عنوان معیار طبقه‌بندی استفاده می‌شود. اگر قرار باشد تنها از یک ترکیب خطی استفاده کنیم، شاید این ترکیب که همان مجموع سطح زیر کشت هر آبادی است، معیار بدی نباشد. اما معمولاً استفاده از یک ترکیب خطی به تنها یک افت اطلاعاتی زیادی را به همراه دارد. اگر قرار باشد از دو ترکیب خطی یا بیشتر استفاده کنیم که افت اطلاعاتی کمتر شود، پیدا کردن این ترکیبات خطی بدون استفاده از مؤلفه‌های اصلی، مشکل و یا غیرممکن می‌نماید. در چنین موقعی می‌توان ترکیبات خطی متناظر با ویژه بردارهای ماتریس واریانس - کوواریانس یا مؤلفه‌های اصلی مهم را که از خواص آماری مناسبی برخوردار هستند

بیشتری ارائه نمی‌شود. توضیحات بیشتر را می‌توان در منابع [۲]، [۳] و [۶] پیدا کرد.

۴- روشهای توکیبی

این روشهای از لحاظ کاربرد بسیار معمول و شامل انواع متفاوتی هستند. متداولترین انواع این روشهای عبارتند از

۱- تک اتصالی^{۱۳} (SL)

۲- اتصال کامل^{۱۴} (CL)

۳- اتصال متوسط^{۱۵} (AL)

۴- روش مرکز هندسی^{۱۶}

۵- روش وارد^{۱۷}

روشهای فوق همگی یک فرایند مشترک دارند ولی در آنها از معیارهای مختلفی برای ادغام خوشها استفاده می‌شود. در این فرایند دنباله‌ای از ادغامها، روی داده‌ها صورت می‌پذیرد. اولین ادغام، روی^{۱۸} خوش‌انجام می‌گیرد که هر کدام شامل یک عنصر است و آخرین ادغام که صورت گیرد، یک خوش بدست می‌آید که شامل همه^{۱۹} عنصر می‌شود. در این فرایند اصول عملیات پایه برای تمامی این روشهای مشابه و بین ترتیب است که در شروع، هر خوش از یک عنصر تشکیل می‌شود. سپس

• مرحله ۱- نزدیکترین زوج از خوشها مجرا شناسائی و در هم ادغام می‌شوند و از تعداد خوشها، یکی کاسته می‌شود.

• مرحله ۲- در صورتی که تعداد خوشها در این مرحله یک باشد، فرایند متوقف می‌شود. در غیر این صورت مرحله ۱ دوباره تکرار می‌شود.

در هر یک از مراحل، آن دو خوشاهای در هم ادغام می‌شوند و یک خوش را تشکیل می‌دهند که نزدیکترین (شبیه‌ترین) خوشها باشند. تفاوت بین روشهای مختلف، تفاوت در تعریف فاصله (یا شباهت) بین یک خوش شامل یک عنصر، با خوش دیگر شامل گروهی از عناصر، یا

یک خوش از لحاظ P متغیر مورد بررسی x_1, x_2, \dots, x_p ، که از این به بعد با بردار X نشان داده می‌شود، دارای تفاوت‌های اندکی با یکدیگر باشند، درحالی که افراد از خوشاهای مختلف، تفاوت‌های زیادی با هم داشته باشند. روشهای مختلفی برای این نوع تحلیل وجود دارد که قالب کلی آنها مختصرآ در زیر توضیح داده می‌شود.

معمولترین روشهای تحلیل خوشاهی را می‌توان در دو قالب کلی، روشهای سلسله مراتبی^{۲۰} و روشهای تقسیم‌بندی^{۲۱}، گروه‌بندی کرد.

۴- روشهای سلسله مراتبی

در این روشهای اساس کار بر ترکیب یا تقسیم متوالی افراد داخل خوشاهای مختلف یا خوشاهای یکسان است. این روشهای نیز همان گونه که از تعریف بر می‌آید خود به دو گروه اصلی تقسیم‌بندی می‌شوند.

الف- روشهای توکیبی یا تراکمی^{۲۲}

ب- روشهای تقسیمی^{۲۳}

در فرایند روشهای توکیبی، ابتدا هر فرد به تنها یک خوش را تشکیل می‌دهد، یعنی در بدو امر به تعداد افراد خوش وجود دارد. سپس در مراحل بعدی، در هر مرحله، این خوشاهای در هم ترکیب شده و خوشاهای بزرگتر را تشکیل می‌دهند، به گونه‌ای که رفتن از یک مرحله به مرحله بعدی، موجب کم شدن یک خوش نسبت به تعداد خوشاهای مرحله قبلی خواهد شد، و در آخرین مرحله کلیه افراد در داخل یک خوش قرار می‌گیرند.

در روشهای تقسیمی، بر عکس روشهای تراکمی، ابتدا همه افراد در یک خوش منظور می‌شوند و سپس در هر مرحله خوش جدیدی اضافه می‌شود به طوری که تعداد خوشاهای هر مرحله یکی بیشتر از تعداد خوشاهای مرحله بالا فاصله قبل از آن است. بدین ترتیب در آخرین مرحله، هر فرد در داخل یک خوش قرار می‌گیرد و تعداد خوشاهای در این مرحله برابر با تعداد افراد است. در خصوص این روشهای به دلیل اینکه کمتر در عمل مورد استفاده قرار می‌گیرند در اینجا توضیح

Single Linkage^{۱۳}

Complete Linkage^{۱۴}

Average Linkage^{۱۵}

Centriod Method^{۱۶}

Ward Method^{۱۷}

Hierarchical Techniques^{۱۸}

Partitioning Techniques^{۱۹}

Agglomerative Techniques^{۱۱}

Divisive Techniques^{۱۲}

که در روش‌های ۴ و ۵، در کلیه مراحل، اطلاعات، همه صفات برای کلیه عناصر مورد استفاده واقع می‌شود.

بین یک خوش با گروهی از عناصر با خوش‌ای دیگر است که آن نیز شامل گروهی از عناصر می‌باشد. بدین ترتیب که

۴- روش‌های تقسیم‌بندی

گرچه در تحلیل خوش‌ای با استفاده از هر روش، تعیین تعداد مطلوب خوش، k ، نیز خود یک هدف است، اما با توجه به آنچه در بالا راجع به روش‌های سلسله مراتبی توضیح داده شد، در فرآیند مربوط به روش‌های مذکور، داشتن عدد k ، یک ضرورت نیست. حال آنکه برعکس، در روش‌های تقسیم‌بندی، داشتن عدد k یک ضرورت است. این موضوع که مقدار k مطلوب چقدر باشد بحث مبسوطی را می‌طلبد که در این مقاله نمی‌گنجد و علاقه‌مندان می‌توانند از منابع و مراجع ذکر شده در فهرست منابع استفاده کنند. علاوه بر این اختلاف، تفاوت فاحش دیگر بین روش‌های سلسله مراتبی و تقسیم‌بندی، آن است که در روش‌های تقسیم‌بندی به مفهوم دقیق کلمه، به ازاء یک k مشخص، باید ابتدا کلیه حالات ممکن تقسیم‌بندی k فرد بین خوش به نحوی که هر خوش شامل حداقل یک فرد باشد درنظر گرفته شود و سپس از بین کلیه حالات ممکن، براساس معیارهایی، حالت مطلوب مشخص گردد. از همین توضیح مشخص است که تعداد حالات ممکن تقسیم‌بندی فوق، بسیار زیاد خواهد بود و این خود یک محدودیت جدی در استفاده کامل از این روش را به همراه دارد. یک روش تقریبی از این قبیل که به روش k میانگین^{۲۱} یا خوش‌بندی سریع^{۲۲} معروف است، در نرم افزارهای آماری وجود دارد. در تحلیلهای عددی که در این مقاله ارائه گردیده، از این روش استفاده شده است.

۵. معیارهای مقایسه روش‌های مختلف ساختن طبقات

به منظور سنجش و مقایسه روش‌های مختلف تعیین طبقات و تعیین روش یا روش‌های برتر، تعریف معیارهای منطقی ضرورت دارد. هنگامی که از سودبخشی روشی نسبت به روش دیگر صحبت به میان می‌آید، معمولاً شاخص یا شاخصهایی برای بیان میزان سودبخشی مطرح

- در روش تک اتصالی، فاصله بین دو خوش، فاصله نزدیکترین افراد، یکی از خوش اول و دیگری از خوش دوم است. به همین دلیل، این روش به روش نزدیکترین همسایه^{۱۸} نیز معروف است.

- در روش اتصال کامل، بر عکس روش تک اتصالی، فاصله بین دو خوش، فاصله بین دورترین افراد، یکی از خوش اول و دیگری از خوش دوم است. به همین دلیل، این روش به روش دورترین همسایه^{۱۹} نیز معروف است.

- در روش اتصال متوسط، فاصله بین دو خوش، میانگین فاصله کلیه زوج عناصر، یکی از خوش اول و دیگری از خوش دوم است.

- در روش مرکز هندسی، هر خوش‌ای که تشکیل می‌شود، بردار میانگین کل عناصر آن محاسبه می‌شود و برای تعیین فاصله آن خوش با خوش دیگر مورد استفاده قرار می‌گیرد.

- در این روش که توسط وارد^{۲۰}[۷] معرفی شد و به روش مینیم واریانس معروف است، ادغامها را به طریقی مورد بررسی قرار می‌دهد که افزایش واریانس حاصل از ادغام به حداقل برسد. در این روش، در هر مرحله ابتدا ادغامهای هر دو زوج از خوش‌ها مورد توجه قرار می‌گیرد و در نهایت آن دو زوجی با هم ادغام می‌شوند که منجر به حداقل افزایش در واریانس شود.

لازم به توضیح است که در روش‌های ۱ تا ۳، ابتدا یک ماتریس فاصله^{۲۱} (یا شابات) تشکیل می‌شود به گونه‌ای که درایه سطر i ام و ستون j ام آن، یک کمیت فاصله‌ای بین فرد i ام و j ام براساس مشاهدات D صفت، این دو فرد را نشان می‌دهد. سپس در کلیه مراحل فرایند، محاسبات لازم بر اساس این ماتریس صورت می‌گیرد. در حالی

^{۱۸} Nearest Neighbour

^{۱۹} Furthest neighbour

^{۲۰} Ward

^{۲۱} Distance Matrix

برای مقایسه سودبخشی دو روش باید $(S)F$ مربوط را محاسبه کرد. بدینه است روشی از بین دو یا چند روش، سودبخشتر است که آن روش نسبت به روشهای دیگر کوچکتر باشد.

۳-۵ دترمینان ماتریس S_W ، S_W

دترمینان ماتریس S_W ، S_W $\det(S_W)$ نیز یکی از ملاکهایی است که می‌تواند برای مقایسه سودبخشی مورد استفاده قرار گیرد. در نرم افزارهای آماری در تحلیل چند متغیره واریانس یک طرفه^{۲۴}، معیار دیگری به نام ضریب *Wilks-Lambda* محاسبه می‌شود که از رابطه $\Lambda = \det(S_T)/\det(S_W)$ به دست می‌آید. اما از آنجایی که برای یک مجموعه مشخص از داده‌ها $\det(S_T)$ مقدار ثابتی است، لذا مقدار این ضریب، Λ ، بستگی مستقیم به مقدار $\det(S_W)$ دارد، بنابراین روشی سودبخش‌تر است که ضریب Λ مربوط به آن روش، در مقایسه با روشهای دیگر کوچکتر باشد.

۶. کاربرد تحلیل خوشه‌ای

جهت مقایسه کارآبی طبقه‌بندی واحدهای جامعه بر اساس چند صفت، یک مطالعه موردنی در سطح شهرستان لردگان از استان چهارمحال و بختیاری و در سطح استان همدان انجام شد. درین مطالعه، تعداد طبقات ۶ و ۷ در نظر گرفته شده که در نمونه‌گیریهای مرکز آمار ایران، معمول بوده است.

شاخصهای مربوط به طبقه‌بندی بر اساس متغیر کمکی مجموع سطح زیرکشت محصولات سالانه عده با استفاده از دو روش دالینیوس و K -میانگین محاسبه و در جداول ۱ و ۲ آورده شده است.

مقایسه شاخصهای مختلف در دو روش به خوبی نشانگر سودبخشی بیشتر روش طبقه‌بندی با استفاده از روش K -میانگین نسبت به روش دالینیوس می‌باشد. به عنوان مثال چنانچه با استفاده از روش K -میانگین آبادی‌های استان همدان را به ۶ طبقه تقسیم کنیم، حدود ۲۵ درصد کاهش واریانس نسبت به استفاده از روش دالینیوس خواهیم داشت و این کاهش واریانس در مورد شهرستان لردگان حدود ۵۶ درصد

می‌شود. اهم شاخصهایی که می‌توانند به نوعی سودبخشی روشهای مطرح شده در قسمتهای قبل برای ساختن طبقات را نسبت به یکدیگر ارزیابی کنند، در زیر توضیح داده می‌شود.

۱-۵ مجموع تغییرات متغیرها در داخل طبقات یا $tr(S_{W'})$

همان طور که قبلاً ذکر شد، در طبقه‌بندی واحدهای جامعه، هدف، کاهش هر چه بیشتر تغییرات داخل طبقات و افزایش آن در بین طبقات می‌باشد. تغییرات کل در یک طبقه‌بندی را می‌توان به صورت $S_T = S_B + S_W$ تقسیم‌بندی کرد که در آن S_T تغییرات کل، S_B تغییرات بین طبقات و S_W تغییرات داخل طبقات را نشان می‌دهد. بنابراین در مقایسه دو روش، روشی سودبخش‌تر است که S_W مربوط به آن روش در مقایسه با S_W' روش دیگر، کوچکتر باشد. در حالت کلی، یعنی وقتی بیش از یک متغیر تحت مطالعه می‌باشد، مجموع تغییرات کلیه متغیرها در داخل طبقات به عنوان ملاک مقایسه دو روش مطرح می‌باشد. بنابراین چنانچه ماتریس مجموع مربعات و مجموع حاصل‌ضربهای متغیرها در داخل طبقات را با S نشان دهیم، مجموع عناصر روی قطر این ماتریس یعنی $tr(S_W)$ ، به عنوان ملاک مقایسه دو روش، مورد استفاده قرار می‌گیرد.

۲-۵ شاخص *Jarque*

در این قسمت به معرفی معیاری از منبع [۴] می‌پردازیم که آن را بنا بر اسم نویسنده، شاخص *Jarque* نامیده‌اند. فرض کنیم که طبقه‌بندی جامعه بر اساس هر کدام از متغیرها به طور جداگانه انجام شود. در این صورت، این طبقه‌بندی‌ها که با علامت S_1^* و S_2^* و ... و S_p^* نشان داده می‌شود، هر کدام برای متغیر مربوطه واریانسی به دست می‌دهد که کمترین مقدار را داراست.

اگر MSE حاصل از طبقه‌بندی جامعه بر پایه متغیر زام را با $V_s^*(j)$ و MSE حاصل از به کار بردن روش مورد نظر را با $V_s(j)$ نشان دهیم، شاخص *Jarque* عبارت است از

$$F(s) = \sum_{j=1}^p \frac{V_s(j)}{V_s^*(j)}$$

بررسی روند تغییرات ($S_{W tr}$) در مقایسه با واریانس صفات مختلف نیز مؤید این موضوع است که اگر در طرح نمونه‌ای مورد نظر، کم کردن هر چه بیشتر واریانس صفات اهمیت داشته باشد، انتخاب صفتی که بیشترین واریانس را دارد می‌تواند منجر به کاهش ($S_{W tr}$) گردد و به عبارت دیگر کارآیی زیادی برای گروه‌بندی مناسب واحدهای جامعه داشته باشد.

می‌باشد. از آنجا که تعداد نمونه رابطه مستقیمی با واریانس دارد، لذا انتظار می‌رود که با استفاده از طبقه‌بندی به روش K -میانگین بجای روش دالینیوس، تعداد نمونه در شهرستان لردگان حدود ۵۶ درصد و در استان همدان حدود ۲۵ درصد کاهش یابد.

۹. نتیجه‌گیری و پیشنهادات

با توجه به نتایج به دست آمده، موارد زیر پیشنهاد می‌شود.

- چنانچه در روش تجزیه و تحلیل خوشه‌ای امکانات محاسباتی لازم برای استفاده از کلیه صفات وجود داشته باشد، بهتر است که همواره از اطلاعات مربوط به تمامی صفات در گروه‌بندی واحدها استفاده شود.
- به کارگیری روش تحلیل مؤلفه‌های اصلی، زمانی توصیه می‌شود که امکان استفاده از خوشه‌بندی بر پایه کلیه صفات فراهم نباشد.
- در صورتی که به منظور استفاده از روش خوشه‌بندی ناچار باشیم که با توجه به امکانات محاسباتی موجود به جای p متغیر از $p < m$ متغیر استفاده کنیم، بهتر است برای این منظور m مؤلفه اصلی اول را به کار ببریم. هر چند در این بررسی عملاً در این باره تحقیقی صورت نگرفته است ولی با توجه به نظریه مؤلفه‌های اصلی می‌توان صحت این امر را پیش‌بینی کرد.
- چنانچه بنا باشد از یک ترکیب خطی از کلیه صفات برای گروه‌بندی واحدها استفاده کنیم، مؤلفه اصلی اول بر هر ترکیب خطی دیگر (مثلًاً مجموع) ارجحیت دارد.
- در صورتی که برای گروه‌بندی، از یک صفت یا ترکیب خطی از چند صفت استفاده به عمل می‌آید، توصیه می‌شود که به جای روش دالینیوس، یکی از تکنیکهای خوشه‌بندی مثلاً تکنیک K -میانگین به کار بردۀ شود.
- اگر بخواهیم فقط از یک صفت برای گروه‌بندی استفاده کنیم بهتر است صفتی را به کار ببریم که با توجه به هدف بررسی، بیشترین واریانس یا ضریب تغییرات را داشته باشد.

۷. کاربرد تحلیل مؤلفه‌های اصلی

بر اساس تئوری، اولین مؤلفه اصلی در مقایسه با هر ترکیب خطی دیگر، از جمله ترکیب خطی ساده مجموع متغیرها، اطلاعات بیشتری را ارائه می‌دهد. در نتیجه استفاده از اولین مؤلفه اصلی به عنوان تنها متغیر کمکی نیز موجب افزایش کارآیی می‌شود. برای نشان دادن کارآیی اولین مؤلفه اصلی در طبقه‌بندی نسبت به ترکیب خطی ساده مجموع، بر اساس اطلاعات سرشماری کشاورزی سال ۶۷ شهرستان لردگان و استان همدان، شاخصهای سه‌گانه محاسبه و در جداول ۳ و ۴ آورده شده است. همان طور که جداول ۳ و ۴ نشان می‌دهند، کارآیی مؤلفه اصلی اول در طبقه‌بندی، در مقایسه با ترکیب جمع، از لحاظ معیارهای ($S_{W tr}$ و Λ) بیشتر است ولی از لحاظ معیار ($F(s)$ این کارآیی قدری کمتر است.

۸. کاربرد صفت یا صفات با بیشترین تغییرات

در مثال حاضر، برای بررسی استفاده از صفت یا صفات با بیشترین واریانس یا ضریب تغییرات به عنوان ملاک مناسبی برای گروه‌بندی واحدهای جامعه، آبادی‌های استان همدان در یک مرحله بر پایه ۵ صفت که ضریب تغییرات متفاوتی دارند و در مرحله دیگر بر پایه ۵ صفت که واریانس متفاوتی دارند گروه‌بندی گردید که نتیجه محاسبات انجام شده در مورد شاخص یا شاخصهای مناسب برای سنجش کارآیی در جداول ۵ و ۶ آورده شده است.

مالحظه روند تغییرات شاخصهای Λ و $F(s)$ در مقایسه با ضریب تغییرات، یانگر این موضوع است که چنانچه برای کلیه صفات مورد بررسی اهمیت یکسانی قائل شویم، انتخاب صفتی که بیشترین ضریب تغییرات را داشته باشد، بیشترین کارآیی را برای طبقه‌بندی واحدهای جامعه به گروه‌های همگن به همراه دارد.

می‌گردد. از جناب آقای ابراهیم خدائی عضو گروه علمی سازمان سنجش آموزش کشور که در انجام محاسبات کامپیوتری کمکهای شایانی نمودند نیز قدردانی می‌شود.

۱۰. تشکر و سپاسگذاری

در خاتمه اضافه می‌نماید که تحقیق حاضر با پشتیبانی مالی و کارشناسی مرکز آمار ایران صورت گرفته است که نگارنده‌گان مقاله مراتب تشکر و قدردانی خود را از آن مرکز اعلام می‌دارند. ضمناً جناب آقای نادر فلاخ، عضو گروه علمی دانشگاه شاهد در تدوین و تنظیم مقاله حاضر کمک شایانی نمودند که بدین وسیله از ایشان نیز تشکر

جدول (۱) شاخصهای سه‌گانه در شهرستان لردگان

روش مورد استفاده	۶ طبقه			۷ طبقه		
	$tr(S_W)$	Λ	$F(s)$	$tr(S_W)$	Λ	$F(s)$
دالینیوس	۴/۷۸۱۵E۱۱	۰/۱۴۰۵۹	۲۰۸۱	۴/۶۴۵E۱۱	۰/۱۱۷۴۳	۵۲۰۰
- میانگین K	۲/۱۳۵۳E۱۱	۰/۰۰۲۲۴	۲۱۶۹	۱/۶۹۲۵E۱۱	۰/۰۰۱۰۹	۳۹۳۷

جدول (۲) شاخصهای سه‌گانه در استان همدان

روش مورد استفاده	۶ طبقه			۷ طبقه		
	$tr(S_W)$	Λ	$F(s)$	$tr(S_W)$	Λ	$F(s)$
دالینیوس	۴/۴۷۶۷E۱۲	۰/۱۴۱۷۱	۱۳۷۲۴	۴/۳۸۵۴E۱۲	۰/۱۲۲۸۶	۳۳۰۴۷
- میانگین K	۳/۳۶۷۷E۱۲	۰/۰۷۹۰۲	۱۳۷۹۳	۲/۷۱۴۳E۱۲	۰/۰۳۵۷۵	۳۳۲۲۱

جدول (۳) مقایسه شاخصهای سه‌گانه سنجش در شهرستان لردگان

روش	سهم از واریانس کل	۶ طبقه			۷ طبقه		
		$tr(S_W)$	Λ	$F(s)$	$tr(S_W)$	Λ	$F(s)$
مولفه اول	۵۳/۸۲	۲/۲ E ۱۱	۰/۰۰۷۰۶	۲۲۳۷	۲/۲ E ۱۱	۰/۰۰۲۲۸	۴۰۳۳
ترکیب جمع	۲۱/۱۶	۲/۸ E ۱۱	۰/۰۰۳۲۱	۲۲۲۹	۲/۷ E ۱۱	۰/۰۰۱۶۱	۳۷۲۹

جدول (۴) مقایسه شاخصهای سه‌گانه سنجش در استان همدان

روش	سهم از واریانس کل	۶ طبقه			۷ طبقه		
		$tr(S_W)$	Λ	$F(s)$	$tr(S_W)$	Λ	$F(s)$
مولفه اول	۷۷/۸۳	۲/۳ E ۱۲	۰/۰۳۷۳۹	۱۳۶۱۹	۲/۳ E ۱۲	۰/۰۳۵۴۲	۳۳۱۴۲
ترکیب جمع	۸/۳۶	۴/۰ E ۱۲	۰/۰۶۳۰۶	۱۳۶۰۰	۳/۸ E ۱۲	۰/۰۳۵۰۴	۳۳۱۲۱

جدول (۵) شاخصهای سنجش کارآیی برای طبقه‌بندی بر پایه ضریب تغییرات

نام صفت	ضریب تغییرات	Λ	$F(s)$
جو دیم	۱۶۷	۰/۰۴۱۹۸	۱۳۸۱۱
گندم دیم	۱۷۱	۰/۰۳۸۲۹	۱۳۷۷۳
هندوانه دیم	۸۱۹	۰/۰۰۸۶۹	۱۳۷۰۲
پیاز	۹۳۰	۰/۰۰۷۲۰	۱۳۷۶۰
برنج	۲۷۳۳	۰/۰۰۰۰۸	۲۴۸۰

جدول (۶) شاخصهای سنجش کارآیی برای طبقه‌بندی بر پایه ضریب تغییرات

نام صفت	واریانس	$tr(S_W)$
پنبه آبی	۲/۸۲۹۶ E_6	۸/۸۱۴۷ E_{12}
نخود آبی	۸/۲۸۴۰ E_8	۸/۳۱۱۹ E_{12}
عدس آبی	۱/۲۰۶۲ E_9	۸/۷۷۰۷ E_{12}
گندم آبی	۸/۰۴۴۳ E_{11}	۷/۱۲۵۰ E_{12}
گندم دیم	۶/۷۰۶۳ E_{11}	۲/۲۷۴۱ E_{12}

مراجع

- [1] Cochran, W.G., 1977. *Sampling Techniques*, 3rd ed, John Wiley & Sons, New Delhi.
- [2] Everitt, B.S., 1993. *Cluster Analysis*, Third Edition, John Wiley & Sons.
- [3] Dillon, W.R. and Goldstein M., 1984. *Multivariate Analysis. Methods and Applications*, John Wiley & Sons. Toronto.
- [4] Jarque, C.M., 1981. *A Solution to the Problem of Optimum Stratification in Multivariate Sampling*, Appl. Statist., Vol. 30, N0. 2, PP 163-169.
- [5] Jolliffe, I.I., 1986. *Principal Component Analysis*, Springer-Verlag, New York.
- [6] Mardia, K.V., Kent, J.T. and Bibby, J.M., 1982. *Multivariate Analysis*, Academic Press, London.
- [7] Ward, J.H., 1963. *Hierarchical grouping to optimize an objective function*, J.Amer.Statist., Vol. 58, PP 236-234.