

خوشه‌بندی چند متغیره داده‌های فضایی بر مبنای مقیاس‌های منطقه‌ای همبستگی فضایی

ناصر طوسی^۱ - ابراهیم زینوند لرستانی^۲

^۱ آمار، دانشکده علوم، دانشگاه گلستان

^۲ گرواکولوژی، دانشگاه علوم کشاورزی و منابع طبیعی گرگان

چکیده: تاکنون معمولاً برای خوشه‌بندی داده‌های فضایی از خوشه‌بندی کلاسیک استفاده شده است، در حالی که خوشه‌بندی کلاسیک فقط بر مبنای اطلاعات چندمتغیره داده‌هاست، لذا برای خوشه‌بندی داده‌های فضایی نیاز به رویکرد و الگوریتمی می‌باشد که علاوه بر استفاده از اطلاعات چندمتغیره داده‌ها، توزیع فضایی داده‌ها را نیز مد نظر قرار دهد. در این مقاله یک روش تقسیم‌بندی شناخته شده به نام k -means را که با ساختار فضایی داده‌ها همگام‌سازی شده است، مطرح می‌شود. در این راستا برای لحاظ کردن موقعیت مکانی داده‌ها در خوشه‌بندی، از مقیاس‌های منطقه‌ای «خود همبستگی فضایی» و شاخص‌های موران و ضریب «گری» و آماره‌ی جی عمومی استفاده شده است.

واژه‌های کلیدی: تحلیل فضایی، خوشه‌بندی، شاخص‌های موران و گری، آماره جی عمومی، k -means

۱ مقدمه

فرض کنید n مشاهده در s مکان موجود باشد که از هر واحد p متغیر اندازه‌گیری شده است. الگوریتم خوشه‌بندی فضایی به دنبال آن است تا هر واحد را به یکی از k خوشه تخصیص دهد. در عمل داده‌ها بی‌وجود دارد که به موقعیت جغرافیایی وابسته اند. خود همبستگی فضایی برای نشان دادن شباهت و نزدیکی مکان داده‌ها استفاده می‌شود.

وجود خود همبستگی فضایی، بین مکان‌هایی که در همسایگی هم هستند را می‌توان با نشانگر خود همبستگی فضایی کل تعیین کرد (Cliff and Ord, 1981). همچنین نیاز باعث شده است که علاوه بر یافتن خود همبستگی کل، خود همبستگی به صورت محلی نیز بدست آورده شود. اگر خود همبستگی کل می‌تواند خلاصه‌ای از همبستگی

فضایی کل منطقه مورد مطالعه را ارائه دهد، خود همبستگی محلی، میزان انحراف از همبستگی کل را نشان داده و در یافتن داده‌های پرت مفید می‌باشد (Boots, 2002). خوشه‌بندی یکی از مهمترین مباحث آنالیز آماری می‌باشد که در اکثر موارد داده‌ها از هم مستقل فرض می‌شوند، در حالی که ممکن است به لحاظ مکانی به هم وابسته باشند و این وابستگی باید در خوشه‌بندی لحاظ شود.

در بخش دوم تحلیل مدل خودهمبستگی فضایی^۱ بیان می‌شود. در این راستا مفاهیم کلی و شاخص‌های خودهمبستگی فضایی معرفی می‌شود. در بخش سوم الگوریتم خوشه‌بندی فضایی پیشنهادی که با لحاظ کردن موقعیت مکانی داده‌ها به خوشه‌بندی می‌پردازد، بیان می‌شود. در بخش چهارم داده‌های به کار برده شده معرفی می‌شوند و در بخش نهایی عملکرد الگوریتم بر روی داده‌های واقعی مشاهده می‌شود.

۲ تحلیل مدل خودهمبستگی فضایی

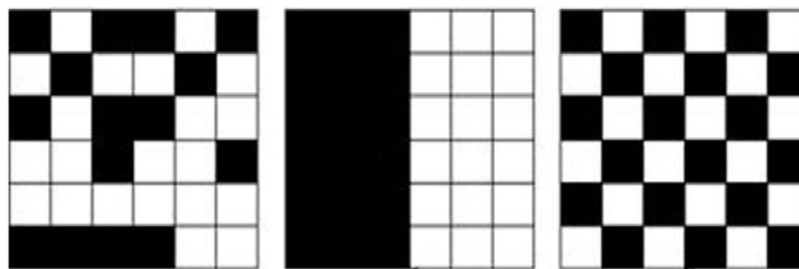
در طبقه‌بندی الگوهای فضایی، خواه خوشه‌ای، پراکنده و تصادفی - می‌توان بر چگونگی نظم و ترتیب قرارگیری واحدهای ناحیه‌ای متمرکز شد. می‌توان مشابهت نبود مشابهت هر جفت از واحدهای ناحیه‌های مجاور را اندازه گرفت. وقتی این مشابهت نبود مشابهت‌ها برای الگوهای فضایی تعیین شود، خود همبستگی فضایی شکل می‌گیرد (Oldand, 1988). خود همبستگی یا همبستگی سریالی، به ارتباط باقی مانده‌های معادله رگرسیونی اشاره دارد. به وسیله خود همبستگی، وضعیتی را توصیف میکنیم که در آن هر باقی مانده یا ضریب خطا مرتبط به ضریب‌های قبلی است (Clark, 1986).

خود همبستگی فضایی قوی بدین مفهوم است که ارزش صفات پدیده‌های جغرافیا بی به طور قوی با یکدیگر رابطه دارند (مثبت یا منفی). ضریب ویژگی توزیع پدیده‌های جغرافیایی مجاور، ارتباطات و نظم ظاهری مختلفی دارد که گفته می‌شود دارای ارتباط فضایی ضعیف، قوی و یا دارای الگوی تصادفی می‌باشند (شکل شماره ۱).

مطالعه همبستگی فضایی، پیشنهادهای ضمنی مهمی برای کاربرد تکنیک‌های آماری در تحلیل اطلاعات فضایی در بر دارد. برای اندازه‌گیری همبستگی فضایی، آماره‌هایی وجود دارد که به ما اجازه می‌دهند با نقاط یا پلی‌گونها (سطوح نواحی) کار کنیم.

این روشها ممکن است برای اندازه‌گیری تعامل فضایی داده‌های عددی و فاصله‌ای، نسبی به کار روند. بخصوص داده‌های شمارشی پیوسته می‌تواند برای تعامل فضایی در میان پلی‌گونه‌ای با داده‌های عددی دوتایی استفاده شود. برای داده‌های فاصله‌ای

^۱ Spatial Autocorrelation



شکل ۱: انواع خودهمبستگی فضایی، همبستگی منفی (قاب راست)، همبستگی مثبت (قاب وسط)، عدم همبستگی فضایی (قاب چپ)

نسبی، شاخص موران (Moran's I) و ضریب گری (Gray Ratio) و شاخص محلی (G-Statistics) به کار می‌روند.

۳ انواع معیارهای تعامل فضایی

مدل‌های متفاوتی برای اندازه‌گیری آماره‌های تعامل فضایی وجود دارد. اگر صفت‌های فضایی یا متغیرهای مورد مطالعه با مقیاس اسمی (Nominal) و دوتایی (Binary) باشند (به عنوان نمونه صفتها فقط دو ارزش ممکن صفر و یک دارند)، پس آماره محاسبات عددی، تعداد اتصالاتها (Joint Count) می‌تواند استفاده شود. اگر متغیرهای فضایی اندازه‌گیری شده، دارای مقیاس فاصله‌ای یا نسبی باشند، آماره‌های ارتباط فضایی مناسب شاخص موران (Moran's I) و ضریب گری (Gary Ratio) می‌باشند و گزینه ممکن دیگر، آماره G عمومی (G-Statistic) است.

یک توزیع یا یک الگوی فضایی می‌تواند در شرایط مختلف از نظر فضایی ناهمگن باشد. برای توصیف ناهمگنی خود همبستگی فضایی، باید بر معیارهایی متکی باشیم که میتوانند خود همبستگی فضایی را در مقیاس محلی کشف کنند. شاخص محلی تمرکز فضایی (Local Indicator of Spatial Association-LISA) و آماره G محلی (Local G-Statistics) برای این هدف مورد استفاده قرار می‌گیرد.

۱.۳ شاخص موران و ضریب گری

شاخص‌های موران و گری^۲ مشخصه‌های مشترکی دارند، اما خواص آماری آنها متفاوت است. اکثر تحلیلگران با شاخص موران موافق‌ترند، که اساساً به خاطر توزیع مشخصاتش، مطلوب‌تر است (Cliff and Ord, 1973, 1981).

^۲ Moran I & Gary Ratio C

هنوز هر دو روش بر مقایسه ارزش‌های همسایگی واحدهای ناحیه‌های متکی هستند. اگر واحدهای ناحیه ای همسایگی در طول ناحیه ارزش‌های مشابهی داشته باشند، آماره ها (مدل‌ها) بر یک خود همبستگی فضایی قوی دلالت داشته‌اند. اگر واحدهای ناحیه‌ای همسایگی ارزش‌های خیلی نامشابهی داشته باشند، آماره‌ها باید یک خودهمبستگی فضایی منفی خیلی قوی را نشان دهند. به هر حال دو مدل، روش‌های متفاوتی را برای مقایسه ضریب همسایگی‌ها به کار می‌گیرند.

۱.۱.۳ شاخص موران

شاخص موران^۳، به شرح ذیل است:

$$I = \frac{N \sum_{i=1}^n \sum_{j=1}^n w_{ij} (x_i - \bar{x})(x_j - \bar{x})}{(\sum_{i=1}^n \sum_{j=1}^n w_{ij}) \sum_{i=1}^n (x_i - \bar{x})^2} \quad (1)$$

که در آن N تعداد کل داده‌ها، X_i داده مکان i ام، X_j داده مکان j ام، w_{ij} شاخص وزنی که مکان i و j را مرتبط می‌سازد.

ضریب موران بین -1 تا 1 متغیر است. -1 برابر تعامل فضایی منفی و 1 برابر تعامل فضایی مثبت به کار می‌رود. اگر تعامل فضایی وجود نداشته باشد، ضریب‌های مورد انتظار موران، برابر است با:

$$E(I) = -\frac{1}{n-1} \quad (2)$$

وقتی شاخص موران محاسبه می‌شود، ماتریس‌های وزنی فضایی مورد استفاده ماتریس‌های دوتایی و تصادفی می‌باشند. اگر شاخص دوتایی (زوجی) استفاده شود، W در مخرج کسر اساساً دو برابر مرزهای مشترک در کل ناحیه مورد مطالعه خواهد بود یا $2T$. به هر حال، این امکان وجود دارد که انواع دیگری از ماتریس‌های وزنی را به کار ببریم.

۲.۱.۳ ماتریس وزن^۴

یکی از مهمترین ابزارهای آمار فضایی ماتریس وزن می‌باشد، چرا که میزان ارتباط فضایی بین داده‌ها به وسیله ماتریس وزن به داده‌ها نسبت داده می‌شود و ما با توجه به ساختار فضایی داده‌ها ماتریس متناسب را روی داده‌ها اعمال می‌کنیم. متداولترین

^۳ Moran's I

^۴ Weight matrix

حالت‌های ماتریس وزن عبارت است از: ماتریس معکوس فاصله، زمان گردش، فاصله ثابت و K همسایگی مجاور و...

در موضوع مورد مطالعه، ماتریس وزنی که لحاظ شده است، ماتریس وزن، براساس معکوس فاصله می‌باشد به این شکل که با توجه به مختصات داده‌ها فاصله آنها را از هم محاسبه و این فاصله را عکس می‌کنیم تا بدین ترتیب هر چه داده‌ها از هم دورتر شوند وزن کمتری بگیرند و بالعکس

$$W_{ij} = \frac{1}{d_{ij}}$$

۳.۱.۳ ضریب «گری»^۵

مشابه روش شاخص موران برای اندازه‌گیری خود همبستگی فضایی، ضریب C گری می‌تواند یک عبارت حاصل ضرب ضریب موران را سازگار کند. ضریب «گری» به صورت فرمول زیر است:

$$C = \frac{(n-1) \sum_i \sum_j W_{ij} (X_i - X_j)^2}{2 \sum_i \sum_j W_{ij} (X_i - \bar{X})^2} \quad (۳)$$

شبهه شاخص موران، ضریب «گری» می‌تواند با هر نوع ماتریس وزنی فضایی به کار رود؛ گرچه عمومی‌ترین آنها، ماتریس‌های دوتایی و تصادفی می‌باشند.

۲.۳ آماره G عمومی^۶

شاخص محلی دیگر خودهمبستگی فضایی آماره‌ی G عمومی است (Getis and Ord 1992). آماره G عمومی محلی برای هر واحد ناحیه‌ای محاسبه می‌شود و بر این دلالت دارد که چگونه ارزش واحد ناحیه مورد مطالعه مرتبط به ارزش‌های واحدهای ناحیه ای مجاور، از طریق آستانه مسافت (d) تعریف شده می‌باشد. و از نظر فرمولی، به شرح زیر است:

$$G_i = \frac{\sum_j w_{ij} x_i}{\sum_j x_j} \quad i = 1, \dots, n \quad (۴)$$

در اینجا نیز بهتر است آماره را در بطن آماره استاندارد شده تفسیر کنیم. برای به دست آوردن آماره استاندارد شده، به دانستن امید ریاضی و واریانس آماره نیاز است. امید

^۵ Geary's Ratio C

^۶ General G-Statistic

ریاضی واریانس به شرح زیر بدست می آید:

$$\text{var}(G_i) = \frac{w_i(n-w_i)s^2}{n^2(n-1)\bar{x}^2}, E(G_i) = \frac{w_i}{n} \quad (5)$$

که در آن $\bar{x} = \sum_i x_i/n$ و $w_i = \sum_j w_{ij}$ و $s^2 = \sum_i (x_i - \bar{x})^2/n$ حال با توجه به (۵) آماره استاندارد شده G عمومی به صورت زیر حاصل می شود که در الگوریتم خوشه بندی که در ادامه مطرح می شود از این آماره استفاده شده است.

$$z(G_i) = \frac{\sum_{j=1}^n w_{ij}x_j - \bar{x}w_{ij}}{\sqrt{\frac{s^2}{n-1}(n(\sum w_{ij}^2) - w_i^2)}} \quad (6)$$

۴ تشخیص خوشه بندی فضایی بر مبنای آماره های محلی برای خود همبستگی فضایی

خوشه داده ها در مباحث آنالیز آماری از پیشینه و سابقه طولانی برخوردار است. تحلیل خوشه ای در جستجوی یافتن روندی برای گروه بندی داده ها می باشد. اساس تحلیل خوشه ای، یافتن مبنایی از شباهت ها و عدم شباهت ها می باشد که در گروه بندی لحاظ می شود.

فرض کنید n مشاهده در s مکان داشته باشیم که از هر واحد p متغیر اندازه گیری شده است. الگوریتم خوشه بندی به دنبال آن است تا هر واحد را به یکی از k خوشه تخصیص دهد. بیشتر الگوریتم های خوشه بندی را می توان به دو دسته تقسیم بندی نمود.

(۱) الگوی پارتیشن بندی واحدها را به تعداد خوشه های از پیش تعیین شده تقسیم می کند،

(۲) روش های سلسله مراتبی خوشه هایی به صورت پیشامدهای تصادفی تولید می کند.

در این مقاله از روش اول یعنی الگوی پارتیشن بندی با تعداد خوشه های از پیش تعیین شده استفاده شده است. یکی از متداول ترین روش ها روش k -means می باشد که مراحل این الگوریتم برای k خوشه به شکل زیر می باشد:
 گام اول: ابتدا k نمونه انتخاب کرده و در k نقطه ی مرکزی قرار می دهیم.
 گام دوم: فاصله $n - k$ نقطه باقیمانده را نسبت به مرکزها حساب کرده و این واحدها

به خوشه‌هایی که فاصله یشان کمترین است قرار می‌گیرد. گام سوم: دوباره مرکزها را حساب و فواصل واحدها را تا مرکزها محاسبه کرده و این عملیات تا جایی ادامه پیدا میکند که به یک پایداری در مرکزها برسد.

حال باتوجه به شناخت معیارهای خودهمبستگی فضایی و روش خوشه‌بندی k-means الگوریتم زیر برای خوشه‌بندی داده‌های فضایی پیشنهاد می‌شود:

(۱) ماتریس وزن W_{ij}

را برای داده‌ها محاسبه کرده که یک ماتریس $n \times n$ می‌باشد. برای هر متغیر، آماره z_j استاندارد شده z_j و آرد محاسبه می‌شود فرمول (۶) $z_j(x_i)$ را برای متغیر j ام ($j = 1, \dots, p$) در i امین واحد ($i = 1, \dots, n$) محاسبه کرده و در نهایت ماتریس Z با بعد $n \times p$ حاصل می‌شود. هر ستون از Z نشانگر مولفه‌های خودهمبستگی فضایی محلی برای یک متغیر است در حالی که هر سطر آن نمایانگر تمرکز خوشه‌های در اطراف هر واحد است.

(۲) الگوریتم خوشه‌بندی k-means را روی داده‌های جدید که در ماتریس Z ذخیره شده است، اعمال می‌شود. این مرحله باعث می‌شود خوشه‌بندی پیشنهادی هم براساس اطلاعات فضایی داده‌ها و هم براساس اطلاعات خود متغیرها انجام شود. الگوریتم k-means تا جایی ادامه پیدا میکند که تابع زیر کمینه گردد.

$$WSS = \sum_{k=1}^K \sum_{C(i)=k} \|z_i - \bar{z}_k\|^2 \quad (7)$$

که در آن z_i ، i امین سطر از ماتریس z و \bar{z}_k ($j = 1, \dots, k$) مرکز خوشه‌ها، فاصله اقلیدسی می‌باشد. $\|x - y\| = \sqrt{\sum_{j=1}^n (x_j - y_j)^2}$

(۳) مراحل بالا برای k مقدار محاسبه می‌شود و برای بدست آوردن تعداد خوشه‌های بهینه از آماره‌ی Gap استفاده می‌شود.

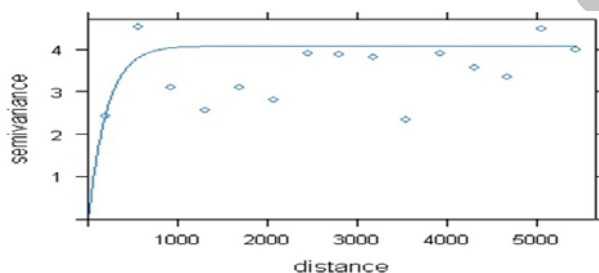
۵ داده‌های استفاده شده

داده‌ها بی که در این مطالعه مورد بررسی قرار گرفته است، از یک بررسی در سال 1389 در شهرستان گرگان (استان گلستان، حوزه قره سو) می‌باشد، که 79 مزرعه گندم به طور تصادفی از شمال، جنوب، غرب و شرق حوزه انتخاب و با استفاده از کادر ۰/۲۵ متر مربعی نمونه‌برداری شده‌اند. در هر کادر بیماری‌های قارچی لکه قهوه‌ای گندم (*Bipolaris sorokiniana*) و زنگ قهوه‌ای (*Puccinia recondita f. sp. tritici*)

با توجه به علائم موجود بر روی اندام گیاه، مشخص و فراوانی آن مورد شمارش قرار گرفته‌اند. سپس با توجه به تعداد کادر پرتاب شده در هر مزرعه میانگین فراوانی هر بیماری قارچی محاسبه گردیده است.

۶ نتایج، بحث و کارهای آینده

برای خوشه‌بندی در محدوده مورد مطالعه به این صورت عمل شد. با استفاده از نرم افزار R و با اعمال الگوریتم پیشنهادی، آماره (۶) برای داده‌ها محاسبه گردید. در این راستا نیاز به ماتریس وزن بوده که از روش معکوس فاصله، درایه‌های ماتریس محاسبه شد. با توجه به وریوگرام رسم شده (شکل ۲) مشاهده شد که تقریباً از فاصله ۵۰۰ متری ساختار فضایی از بین می‌رود، لذا درایه‌هایی از ماتریس وزن را که کمتر از (۱/۵۰۰) بود، برابر صفر قرار داده شد. با توجه به ماتریس جدید، مقادیر Z محاسبه و خوشه‌بندی انجام گرفت. نتایج خوشه‌بندی به صورت جدول شماره (۱) می‌باشد که نتایج گام دوم که دارای خطای کمتری بود، انتخاب گردید.

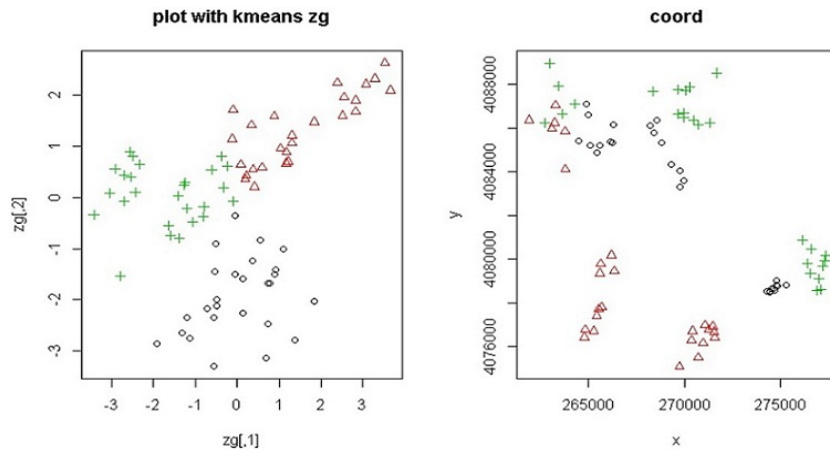


شکل ۲: وریوگرام بیماری قارچی لکه قهوه‌ای (*Bipolaris sorokiniana*)

جدول ۱: خطای درونی و خطای کل مربوط به سه مرحله اجرای الگوریتم k - means

گام	خطای درونی خوشه‌ها			خطای کل
اول	۳۵/۵۱۲۱۳	۳۳/۳۱۵۷۸	۴۹/۰۸۱۵۴	۱۱۷/۹۰۹۵
دوم	۲۶/۶۹۰۳۲	۵۸/۲۶۹۵۱	۳۲/۹۰۱۰۰	۱۱۷/۸۶۰۸
سوم	۳۳/۸۶۱۷۱	۳۲/۹۰۱۰۰	۵۱/۳۳۱۷۶	۱۱۸/۰۹۴۵

با توجه به شکل (۳)، روش پیشنهادی برای خوشه‌بندی داده‌های فضایی، موقعیت مکانی داده‌ها را به خوبی لحاظ کرده است به طوری که داده‌ها بی که در یک خوشه



شکل ۳: نمودار پراکندگی داده‌ها بر اساس آماره (۶) و خوشه‌بندی فضایی (قاب چپ)، نمودار پراکندگی داده‌ها بر اساس مختصات داده‌ها و خوشه‌بندی فضایی (قاب راست)

قرار گرفته اند تقریباً از لحاظ مکانی نزدیک هم می‌باشد. این موضوع از لحاظ زمین آماری بسیار مهم و مقرون به صرفه می‌باشد. همچنین لازم به ذکر است که در گذشته در این موارد از خوشه‌بندی کلاسیک استفاده شده که خوشه‌هایی با داده‌های پراکنده از لحاظ مکانی ایجاد کرده است. امید است در آینده با وارد کردن خوشه‌بندی فازی به این الگوریتم دقت این روش را افزایش دهیم.

مراجع

- Boots, B. (2002), Local measures of spatial association, *Ecoscience*, 9(2):168-176.
- Cliff A. D. and Ord, J. K. (1981), *Spatial processes - models and applications*, Pion, London.
- Getis, A. and Ord, J. K. (1992), The analysis of spatial association by use of distance statistics, *Geographical Analysis*, 24:189-206.
- Ord J.K., Getis A. (1995) Local spatial autocorrelation statistics: distributional issues and an application. *Geographical Analysis*, 27: 286-306.
- Scrucca, L. (2005), Clustering multivariate spatial data based on local measures of spatial autocorrelation. *Geographical Analysis*, 38: 34-59.