



سنجش از دور & GIS ایران



سنجش از دور و GIS ایران سال دهم، شماره سوم، پاییز ۱۳۹۷
Iranian Remote Sensing & GIS Vol.10, No.3, Autumn 2018

۱۷-۳۲

بهبود خوشه‌بندی تصاویر فراطیفی با به‌کارگیری دیورژانس اطلاعات طیفی

حمید عزت‌آبادی پور*

مربی دانشکده مهندسی عمران، دانشگاه صنعتی سیرجان

تاریخ پذیرش مقاله: ۹۷/۴/۹

تاریخ دریافت مقاله: ۹۶/۶/۲۰

چکیده

الگوریتم خوشه‌بندی K-Means یکی از پرکاربردترین روش‌های طبقه‌بندی نظارت نشده در پردازش تصاویر سنجش از دور است. در الگوریتم K-Means استاندارد، از معیار عدم شباهت فاصله اقلیدسی، به منظور اندازه‌گیری عدم شباهت بین داده‌ها و خوشه‌ها استفاده می‌شود. فاصله اقلیدسی، یک معیار عدم شباهت قطعی است که بردار طیفی پیکسل‌ها و مراکز خوشه‌ها را به صورت نقاطی در یک فضای چندبعدی در نظر می‌گیرد و فاصله هندسی بین آن‌ها را اندازه‌گیری می‌کند. تصاویر فراطیفی همواره دارای عدم قطعیت هستند، به همین دلیل استفاده از یک معیار عدم شباهت آماری (غیرقطعی)، جهت خوشه‌بندی آن‌ها مناسب‌تر به نظر می‌رسد. بر این اساس در این مقاله، با به‌کارگیری یک معیار عدم شباهت آماری، یک روش نظارت نشده جدید برای خوشه‌بندی تصاویر فراطیفی طراحی و پیاده‌سازی شده است. روش خوشه‌بندی پیشنهادی، برای برآورد عدم شباهت بین مرکز خوشه‌ها و پیکسل‌ها، از یک معیار عدم شباهت آماری، به نام دیورژانس اطلاعات طیفی، به جای فاصله اقلیدسی استفاده می‌کند. دیورژانس اطلاعات طیفی، توزیع احتمال طیف‌ها را از طریق نرمال کردن امضای طیفی، مدل‌سازی می‌کند. سپس فاصله بین توزیع احتمال طیف یک پیکسل و توزیع احتمال طیف هر مرکز خوشه را برآورد می‌کند. آزمون‌های انجام‌شده بر روی داده‌های تصویری فراطیفی واقعی حاصل از سه سنجنده HyMap، HYDICE و Hyperion نشان می‌دهد که روش خوشه‌بندی پیشنهادی مبتنی بر دیورژانس اطلاعات طیفی، نتایج طبقه‌بندی را بهبود می‌بخشد، به طوری که ضریب کاپای نتایج طبقه‌بندی تصاویر فراطیفی مورد استفاده به ترتیب، حدود ۰/۷٪، ۵۶٪ و ۱۰٪ افزایش یافته است.

کلیدواژه‌ها: خوشه‌بندی، معیار عدم شباهت، دیورژانس اطلاعات طیفی، تصاویر فراطیفی

* نویسنده مکاتبه‌کننده: سیرجان، ابتدای جاده بافت، دانشگاه صنعتی سیرجان، کد پستی: ۷۸۱۳۷۳۳۳۸۵، تلفن ثابت: ۰۳۴-۴۲۲۴۲۱۱۶ - ۰۹۱۳۳۷۹۲۹۴۱

۱- مقدمه

نمی‌توانند دانش قبلی راجع به شکل یا اندازه کلی خوشه‌ها را در نظر بگیرند (Tsai et al., 2002). روش‌های خوشه‌بندی سلسله‌مراتبی پویا نیست، یعنی در این گروه خوشه‌بندی، داده‌هایی که در مراحل اولیه به یک خوشه ویژه تعلق می‌گیرند، نمی‌توانند به خوشه دیگر منتقل شوند. از این رو، این قبیل روش‌های خوشه‌بندی، برای مجموعه داده‌های پیچیده که خوشه‌های آن‌ها خیلی همگن نباشند، مناسب نیستند (Brereton, 1992). به‌علاوه، روش‌های سلسله‌مراتبی به حافظه و زمان محاسباتی زیادی نیاز دارند.

اما در روش خوشه‌بندی جزءبندی، اغلب از روش‌های بهینه‌سازی تناوبی، برای تشکیل خوشه‌ها استفاده می‌شود. ولی طبیعت تکراری این روش‌ها، آن‌ها را به مقاداردهی اولیه و کمینه‌های محلی حساس می‌کند. خوشه‌بندی جزءبندی، نیازمند تعیین تعداد بهینه خوشه‌هاست و به نویز و مشاهدات خطا حساس است. این نوع خوشه‌بندی، از طریق نمونه‌های اولیه^۴ و اندازه‌گیری فاصله، می‌تواند دانش مربوط به شکل یا اندازه خوشه‌ها را در خود جای دهد (Guha et al., 2001). برخلاف خوشه‌بندی سلسله‌مراتبی، روش خوشه‌بندی جزءبندی پویاست، یعنی در خوشه‌بندی جزءبندی، داده‌ها می‌توانند در تکرارهای متوالی از یک خوشه به خوشه دیگر جابجا شوند.

روش‌های خوشه‌بندی مورد استفاده در حوزه سنجش از دور، عمدتاً از نوع جزءبندی هستند که معروف‌ترین و پرکاربردترین آن K-Means است. الگوریتم خوشه‌بندی K-Means، برای اولین بار توسط مک‌کوئین در سال ۱۹۶۷ مورد استفاده قرار گرفت و در سال ۲۰۰۱ توسط دودا و همکارانش به صورت کامل‌تری ارائه شد

روش‌های طبقه‌بندی نظارت نشده، یکی از مهمترین روش‌های تفسیر تصاویر و استخراج اطلاعات برای کاربردهای گوناگون سنجش از دور است. این روش‌ها، فقط متکی بر داده‌های تصویری بوده و اغلب به صورت خودکار انجام می‌شوند. اگرچه روش‌های نظارت نشده، نسبت به روش‌های نظارت شده معمولاً دارای دقت پایین‌تری هستند، اما نیازمند هزینه و اطلاعات کم‌تری بوده و به همین دلیل توجه بسیاری از محققان را به خود جلب کرده‌اند. در میان روش‌های طبقه‌بندی نظارت نشده، روش‌های خوشه‌بندی از اهمیت ویژه‌ای برخوردارند. در این روش‌ها، هدف عبارتست از طبقه‌بندی داده‌ها به شکلی که دو داده در یک خوشه، تا حد امکان به هم شبیه و در دو خوشه متفاوت تا حد امکان از هم متمایز باشند (Timm et al., 2004). روش‌های خوشه‌بندی به طور کلی به دو گروه سلسله‌مراتبی^۱ و جزءبندی^۲ تقسیم می‌شوند (Jain and Dubes, 1988).

خوشه‌بندی سلسله‌مراتبی، عبارت است از تشکیل متوالی گروه‌هایی که عضوهای آن‌ها بیشترین شباهت را به هم داشته یا جداسازی متوالی گروه‌هایی که عضوهای آن‌ها بیشترین اختلاف را با هم دارند. در خوشه‌بندی سلسله‌مراتبی، مشکلات ناشی از مقاداردهی اولیه و کمینه‌های محلی وجود ندارد (Tsai et al., 2002). در این نوع خوشه‌بندی، نویز و مشاهدات خطا^۳ در خوشه‌های جداگانه قرار گرفته و خوشه‌های دیگر را تحت تأثیر خود قرار نمی‌دهند (Tran et al., 2003). در عوض، روش‌های سلسله‌مراتبی به معیارهای مشخصی، جهت تعیین تعداد بهینه خوشه‌ها نیاز داشته و

^۴Outliers

^۳Prototypes

Hierarchical

Partitioning

آماري، ديورژانس اطلاعات طيفي^۴ (SID)، معيار عدم شباهت شناخته شده‌ايست که برای اولین بار توسط چانگ در سال ۲۰۰۰ به منظور اندازه‌گیری عدم شباهت بين دو طيف در تصاویر فراطيفي ارائه شد (Chang, 2000).

معيار عدم شباهت SID، تاکنون مورد توجه بسياری از محققان قرار گرفته و در زمينه‌های مختلفی جهت آناليز تصاویر فراطيفي استفاده شده است (Du et al., 2004; Chen et al., 2009; Galal et al., 2012; Adep et al., 2016; Palsson et al., 2017; Erudel et al., 2017; Gholizadeh et al., 2018). اما در زمينه خوشه‌بندی تصاویر فراطيفي با استفاده از معيار عدم شباهت SID، تحقیقات زیادی انجام نشده است. از این رو، ضروری است در این زمينه، مطالعه بیشتری صورت پذیرد. در همین راستا، ایده اصلی در این مقاله، طراحی و ابداع یک روش نظارت نشده جدید، برای خوشه‌بندی تصاویر فراطيفي، با به‌کارگیری معيار عدم شباهت SID است. با توجه به مطالب فوق، انتظار می‌رود که روش خوشه‌بندی جدید پیشنهادی، نتیجه خوشه‌بندی تصاویر فراطيفي را بهبود بخشد.

این مقاله شامل چهار بخش است. در بخش اول، به مقدمه‌ای کوتاه درباره انگیزه، هدف و بیان مسأله تحقیق و روش کار پرداخته شده است. در بخش دوم، مبانی نظری روش‌ها تشریح می‌شود. بخش سوم، در برگیرنده پیاده‌سازی الگوریتم‌های خوشه‌بندی و ارزیابی آنها خواهد بود. در بخش چهارم و پایانی نیز، نتیجه‌گیری‌ها ارائه می‌شود.

(MacQueen, 1967; Duda et al., 2001). الگوریتم K-Means استاندارد، از معيار عدم شباهت^۱ فاصله اقلیدسی، جهت اندازه‌گیری اختلاف بين داده‌ها و خوشه‌ها استفاده می‌کند. این در حالیست که در آناليز تصاویر فراطيفي، معيارهای گوناگونی برای ارزیابی عدم شباهت، ارائه و به‌کار رفته است.

معيارهای عدم شباهت مورد استفاده برای تصاویر فراطيفي، با توجه به راهبرد اندازه‌گیری عدم شباهت بين دو طيف، به دو دسته قطعی^۲ و آماری^۳ تقسیم می‌شوند (Chang, 2003; Homayouni and Roux, 2006; van der Meer, 2006). معيارهای عدم شباهت قطعی (مانند فاصله اقلیدسی و نگاشت‌کننده زاویه طيفي)، هر طيف را به صورت یک بردار طيفي n -بعدی (n : تعداد باندهای طيفي) در نظر گرفته و عدم شباهت بين دو بردار طيفي را در فضای n -بعدی اندازه‌گیری می‌کنند (van der Meer, 2006). در مقابل، معيارهای عدم شباهت آماری، هر بردار طيفي را به صورت یک متغیر تصادفی در نظر گرفته و اختلاف بين توزیع احتمال دو بردار طيفي را اندازه‌گیری می‌کنند (Chang, 2003).

عواملی چون بی‌ثباتی سنجنده، تغییرات توپوگرافی سطح زمین و اثرات محیطی و جوی، همواره سبب بروز عدم قطعیت در داده‌های سنجنش از دور، به‌ویژه داده‌های فراطيفي می‌شوند (Chang, 2003; Shi, 2009). به همین دلیل در پردازش و آناليز تصاویر فراطيفي، معيارهای عدم شباهت آماری (غیرقطعی) برای اندازه‌گیری عدم شباهت مناسب‌تر بوده و به نتایج بهتری منجر می‌شوند (Chang, 2000; Chang, 2003; van der Meer, 2006). در میان معيارهای عدم شباهت

^۲Stochastic

^۳Spectral Information Divergence

^۱Dissimilarity

^۲Deterministic

۲- مواد و روش‌ها

۲-۱- الگوریتم خوشه‌بندی K-Means

الگوریتم خوشه‌بندی K-Means استاندارد، از معیار عدم شباهت فاصله اقلیدسی، به‌منظور اندازه‌گیری اختلاف بین داده‌ها و خوشه‌ها استفاده می‌کند و در آن هدف، کمینه کردن تابع هدف زیر است (Duda et al., 2001; Jain, 2010).

$$J(V) = \sum_{i=1}^K \sum_{\vec{x}_j \in C_i} \|\vec{x}_j - \vec{v}_i\|^2 \quad (1) \text{ رابطه}$$

در رابطه (۱)، K تعداد خوشه‌ها، C_i خوشه i ام، \vec{x}_j بردار طیفی i امین پیکسل، \vec{v}_i مرکز i امین خوشه، و $\|\vec{x}_j - \vec{v}_i\|^2$ مربع فاصله اقلیدسی بین \vec{x}_j و \vec{v}_i است. به $\|\vec{x}_j - \vec{v}_i\|^2$ مربع خطا^۱ نیز اطلاق می‌شود (Duda et al., 2001). شرط لازم برای کمینه‌سازی تابع هدف J ، به‌سادگی از طریق برابر با صفر قرار دادن گرادیان رابطه (۱)، نسبت به \vec{v}_i که می‌بایست بهینه شود، به دست می‌آید. با انجام این عمل، رابطه زیر حاصل می‌شود:

$$\vec{v}_i = \frac{\sum_{\vec{x}_j \in C_i} \vec{x}_j}{m_i} \quad \forall i = 1, 2, \dots, K \quad (2) \text{ رابطه}$$

در رابطه (۲)، m_i تعداد پیکسل‌های متعلق به C_i است. تابع هدف J را نمی‌توان به‌طور مستقیم کمینه کرد. از این رو، کمینه‌سازی آن به‌صورت تکراری انجام می‌شود و مراکز خوشه‌ها از طریق رابطه (۲) در هر تکرار، بهینه می‌شوند.

۲-۲- الگوریتم خوشه‌بندی SID-Based

در الگوریتم خوشه‌بندی پیشنهادی مبتنی بر SID (SID-Based)، به جای فاصله اقلیدسی، از دیورژانس اطلاعات طیفی (SID)، جهت اندازه‌گیری عدم شباهت استفاده شده است. SID از مفهوم دیورژانس در تئوری

اطلاعات، مشتق شده و اختلاف رفتارهای احتمالی بین طیف‌های دو بردار پیکسلی را اندازه‌گیری می‌کند (Chang, 2003). به عبارت دیگر، SID عدم شباهت بین دو بردار پیکسلی را بر اساس اختلاف بین توزیع‌های احتمال که از طریق نرمال کردن طیف به دست می‌آید، اندازه‌گیری می‌کند. پیش‌تر در شناسایی الگو و انتخاب باندهای بهینه نیز از دیورژانس استفاده شده است (Tou and Gonzalez, 1974; Jensen, 1996). در مقایسه با فاصله اقلیدسی که فاصله فضایی بین دو بردار پیکسلی را به دست می‌دهد، SID فاصله بین توزیع‌های احتمال ایجاد شده توسط طیف‌های دو بردار پیکسلی را اندازه‌گیری می‌کند. از این رو SID، در بهره‌گیری از تغییرات طیفی می‌تواند مؤثرتر عمل کند (Chang, 2003). دیورژانس اطلاعات طیفی بین بردار طیفی \vec{v}_i و بردار طیفی پیکسلی \vec{x}_j به صورت زیر تعریف می‌شود:

$$SID(\vec{v}_i, \vec{x}_j) = D(\vec{v}_i \parallel \vec{x}_j) + D(\vec{x}_j \parallel \vec{v}_i) \quad (3) \text{ رابطه}$$

$$D(\vec{v}_i \parallel \vec{x}_j) = \sum_{l=1}^n p_{il} D_l(\vec{v}_i \parallel \vec{x}_j) \quad (4) \text{ رابطه}$$

$$= \sum_{l=1}^n p_{il} (I_l(\vec{x}_j) - I_l(\vec{v}_i))$$

$$D(\vec{x}_j \parallel \vec{v}_i) = \sum_{l=1}^n q_{jl} D_l(\vec{x}_j \parallel \vec{v}_i) \quad (5) \text{ رابطه}$$

$$= \sum_{l=1}^n q_{jl} (I_l(\vec{v}_i) - I_l(\vec{x}_j))$$

در تئوری اطلاعات $D(\vec{v}_i \parallel \vec{x}_j)$ و $D(\vec{x}_j \parallel \vec{v}_i)$ را به ترتیب، انتروپی \vec{x}_j نسبت به \vec{v}_i و انتروپی \vec{v}_i نسبت به \vec{x}_j می‌نامند. بردارهای احتمال $\vec{p}_i = (p_{i1}, p_{i2}, \dots, p_{in})^T$ و $\vec{q}_j = (q_{j1}, q_{j2}, \dots, q_{jn})^T$ به ترتیب، توزیع احتمال \vec{v}_i و توزیع احتمال \vec{x}_j هستند. این توزیع‌های احتمال از طریق نرمال کردن طیف، به دست آمده و مولفه‌هایشان به صورت زیر محاسبه می‌شود:

^۱Entropy

^۱Squared-Error

$$\frac{\partial J(P)}{\partial \vec{p}_i} = \sum_{\vec{x}_j \in C_i} \left(\log(\vec{p}_i) - \log(\vec{q}_j) + \frac{1}{\vec{p}_i} (\vec{p}_i - \vec{q}_j) \right) = 0$$

$$\Rightarrow \sum_{\vec{x}_j \in C_i} \left(\log(\vec{p}_i) - \log(\vec{q}_j) + 1 - \frac{\vec{q}_j}{\vec{p}_i} \right) = 0$$

$$\Rightarrow m_i \log(\vec{p}_i) - \sum_{\vec{x}_j \in C_i} \log(\vec{q}_j) + m_i - \frac{\sum_{\vec{x}_j \in C_i} \vec{q}_j}{\vec{p}_i} = 0$$

$$\Rightarrow m_i \log(\vec{p}_i) - \frac{\sum_{\vec{x}_j \in C_i} \vec{q}_j}{\vec{p}_i} \quad \text{رابطه (۱۲)}$$

$$= \sum_{\vec{x}_j \in C_i} \log(\vec{q}_j) - m_i$$

رابطه بالا را می‌توان از طریق تابع Wright Omega که توسط گرلس و جفایری جهت حل معادله $Y + \log(Y) = X$ ارائه شده، حل نمود (Corless and Jeffrey, 2002). با حل رابطه (۱۲) داریم:

$$\vec{p}_i = \frac{\sum_{\vec{x}_j \in C_i} \vec{q}_j}{m_i \times \text{WrightOmega} \left(\frac{m_i - \sum_{\vec{x}_j \in C_i} \log(\vec{q}_j)}{m_i} - \log \left(\frac{m_i}{\sum_{\vec{x}_j \in C_i} \vec{q}_j} \right) \right)} \quad \forall i = 1, 2, \dots, K$$

در اینجا نیز، کمینه‌سازی تابع هدف J به صورت تکراری، انجام می‌شود و توزیع احتمال مراکز خوشه‌ها (\vec{p}_i) ، از طریق رابطه (۱۳) در هر تکرار، بهینه می‌شوند.

۳- پیاده‌سازی و نتایج

همان‌طور که پیش‌تر بیان شد، الگوریتم‌های خوشه‌بندی K-Means و SID-Based، به صورت تکراری اجرا می‌شوند. روند اجرای این الگوریتم‌ها در شکل (۱) نشان داده شده است.

$$p_{ik} = \frac{v_{ik}}{\sum_{h=1}^n v_{ih}} \quad \forall k = 1, 2, \dots, n \quad \text{رابطه (۶)}$$

$$q_{jk} = \frac{x_{jk}}{\sum_{h=1}^n x_{jh}} \quad \forall k = 1, 2, \dots, n \quad \text{رابطه (۷)}$$

در روابط (۴) و (۵)، $I_l(\vec{x}_j)$ و $I_l(\vec{v}_i)$ به عنوان self-information \vec{v}_i و \vec{x}_j برای باند l تعریف و به صورت زیر محاسبه می‌شوند:

$$I_l(\vec{v}_i) = -\log p_{il} \quad \text{رابطه (۸)}$$

$$I_l(\vec{x}_j) = -\log q_{jl} \quad \text{رابطه (۹)}$$

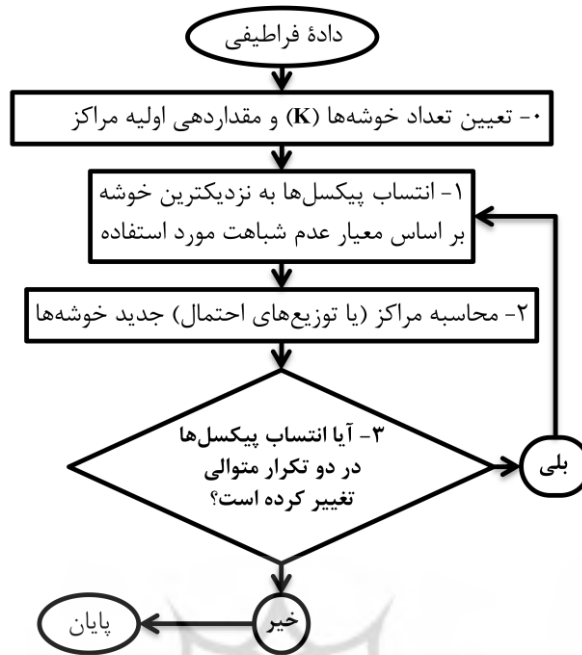
با جای‌گذاری روابط (۴) و (۵) در رابطه (۳)، رابطه محاسبه دیورژانس اطلاعات طیفی به صورت زیر، ساده می‌شود:

$$SID(\vec{v}_i, \vec{x}_j) = \sum_{l=1}^n (p_{il} - q_{jl})(\log(p_{il}) - \log(q_{jl})) \quad \text{رابطه (۱۰)}$$

حال، جهت پیاده‌سازی الگوریتم خوشه‌بندی SID-Based، دیورژانس اطلاعات طیفی جایگزین مربع فاصله اقلیدسی در تابع هدف الگوریتم خوشه‌بندی K-Means می‌شود و تابع هدف زیر به دست می‌آید:

$$J(P) = \sum_{i=1}^K \sum_{\vec{x}_j \in C_i} SID(\vec{v}_i, \vec{x}_j) = \sum_{i=1}^K \sum_{\vec{x}_j \in C_i} \left(\sum_{l=1}^n (p_{il} - q_{jl})(\log(p_{il}) - \log(q_{jl})) \right) \quad \text{رابطه (۱۱)}$$

شرط لازم برای کمینه‌سازی تابع هدف بالا، از طریق برابر با صفر قرار دادن گرادیان آن نسبت به \vec{p}_i که می‌بایست بهینه شود، حاصل می‌شود:



شکل ۱. نمودار گردش کار الگوریتم‌های K-Means و SID-Based

- همان‌طور که در شکل (۱) نشان داده شده است، الگوریتم‌های خوشه‌بندی K-Means و SID-Based، نیازمند تعیین مقادیر اولیه برای مراکز خوشه‌ها هستند. طبیعت تکراری این الگوریتم‌ها، آن‌ها را به مقداردهی اولیه حساس کرده و باعث می‌شود به کمینه‌های محلی، منجر شوند. لذا، تعیین مقادیر اولیه مناسب برای مراکز خوشه‌ها دارای اهمیت ویژه‌ایست. یک روش مقداردهی اولیه، توسط الداود برای مجموعه داده‌های مختلف ۲، ۴ و ۸ بُعدی به کار برده شده است. نتایج، نشان داده که این روش در مقایسه با روش‌های تصادفی به نتایج بهتری منجر می‌شود (Al-Daoud, 2007). این روش به صورت زیر، مقادیر اولیه مراکز خوشه‌ها را به دست می‌آورد:
- ۱- محاسبه واریانس باندهای تصویر
 - ۲- پیدا کردن باندهای که بیشترین واریانس را دارد ($b_v max$)
- ۳- منظم کردن پیکسل‌ها بر اساس مقادیر درجات خاکستری باند $b_v max$
- ۴- تقسیم پیکسل‌های منظم شده به K زیرمجموعه (با تعداد مساوی)
- ۵- یافتن پیکسل میانه هر یک از زیرمجموعه‌ها
- ۶- استخراج طیف پیکسل‌های میانه، به عنوان مقادیر اولیه مراکز خوشه‌ها
- در این مقاله نیز، از این روش، به منظور مقداردهی اولیه مراکز خوشه‌ها، استفاده شده است. با این تفاوت که به جای $b_v max$ ، از اولین مولفه PCA^۱ داده تصویری، استفاده شده است. اعمال این تغییر در این روش مقداردهی اولیه، باعث می‌شود مراکز اولیه خوشه‌ها با تمایز بیشتری نسبت به هم انتخاب شوند، زیرا به دلیل همبستگی زیاد، میان داده‌های تصویری فراطیفی در فضای طیفی، واریانس $b_v max$ خیلی بیشتر از واریانس

^۱Principal Component Analysis

تصویری فراطیفی، به عنوان مقادیر اولیه مراکز خوشه‌ها در این مقاله، از سه دسته داده تصویری فراطیفی، به منظور ارزیابی و مقایسه دقت و کارایی الگوریتم‌های خوشه‌بندی استفاده شده است. با توجه به اینکه، معیار عدم شباهت SID، برای طیف‌های پیکسلی طراحی شده است، الگوریتم‌های خوشه‌بندی K-Means و SID-Based، هر دو به طور یکسان، روی اصل داده‌های تصویری فراطیفی که شامل طیف‌های پیکسلی هستند، اجرا شده‌اند.

۳-۱- داده‌های مورد استفاده

داده‌های تصویری مورد استفاده در این مقاله، شامل تصاویر فراطیفی Berlin، Urban و Botswana هستند که مشخصات آن‌ها در شکل‌های (۲) تا (۴) نشان داده شده است.

سایر باندها (مخصوصاً باندهای مجاور) نیست. این درحالیست که محور اولین مولفه PCA در فضای تبدیل شده PCA، در جهتی قرار می‌گیرد که در آن، داده‌ها بیشترین پراکندگی را دارند و به همین دلیل، واریانس اولین مولفه PCA نسبت به سایر مولفه‌ها به طور قابل توجهی بیشتر است.

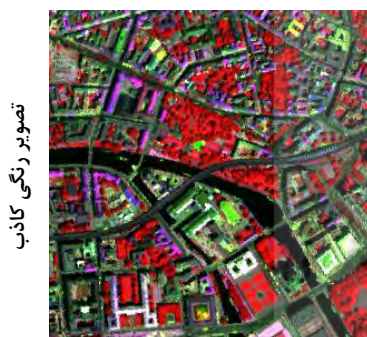
با توجه به مطالب بیان شده بالا، مقداردهی اولیه مراکز خوشه، در این مقاله به صورت زیر انجام شده است:

- ۱- اعمال تبدیل PCA بر روی داده تصویری فراطیفی
- ۲- منظم کردن پیکسل‌ها بر اساس مقادیر درجات خاکستری اولین مولفه PCA

۳- تقسیم پیکسل‌های منظم شده به K زیرمجموعه (با تعداد مساوی)

۴- یافتن پیکسل میانه هر یک از زیرمجموعه‌ها

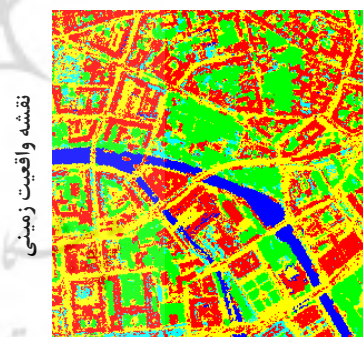
۵- استخراج طیف پیکسل‌های میانه از اصل داده



تصویر رنگی کاذب

سنجنده: HyMap

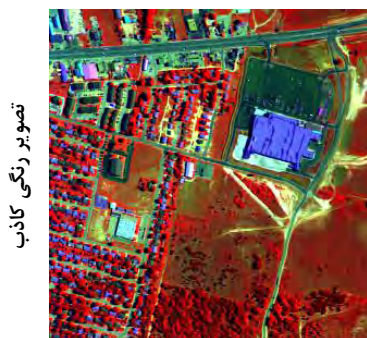
ابعاد پیکسل: ۳.۵ متر
محدوده طیفی: ۲.۵-۴۵.۰ میکرومتر
تعداد باندهای طیفی: ۱۱۴
ابعاد: ۳۰۰×۳۰۰ پیکسل



نقشه واقعیت زمینی

Vegetation
Built-up
Impervious
Soil
Water

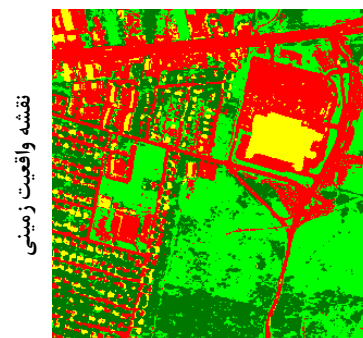
شکل ۲. مشخصات داده تصویری Berlin



تصویر رنگی کاذب

سنجنده: HYDICE

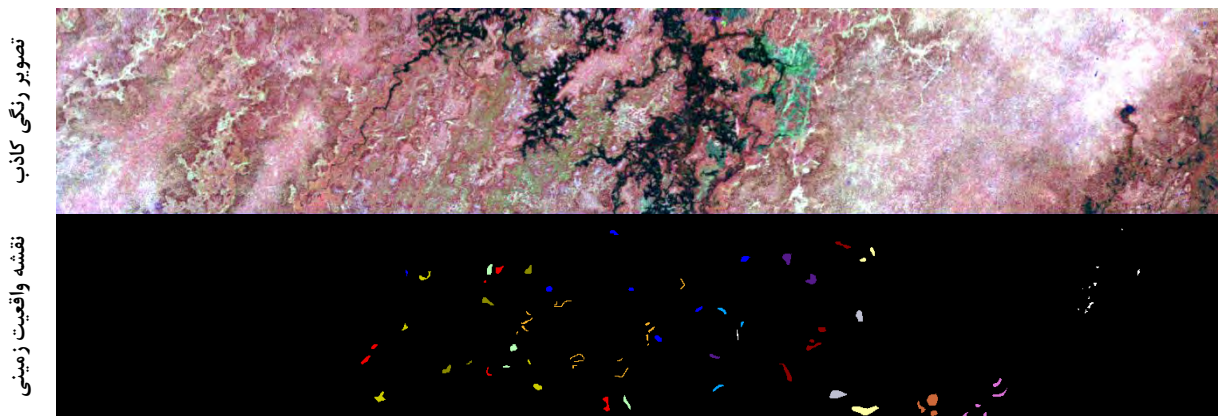
ابعاد پیکسل: ۲ متر
محدوده طیفی: ۲.۵-۴۰.۰ میکرومتر
تعداد باندهای طیفی: ۱۶۲
ابعاد: ۳۰۷×۳۰۷ پیکسل



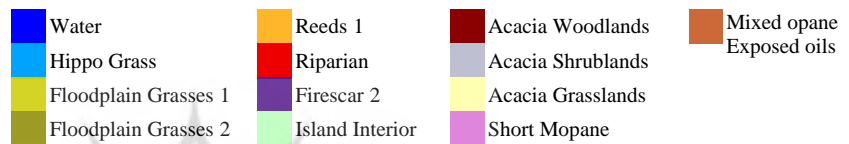
نقشه واقعیت زمینی

Grass
Tree
Asphalt
Roof

شکل ۳. مشخصات داده تصویری Urban



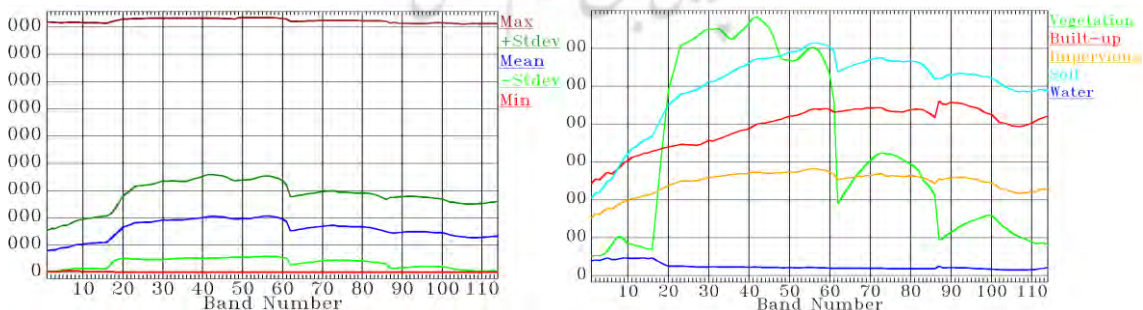
سنجده: Hyperion
 ابعاد پیکسل: ۳۰ متر
 محدوده طیفی: ۲,۵-۴,۰ میکرومتر
 تعداد باندهای طیفی: ۱۴۵
 ابعاد: ۱۴۷۶×۲۵۶ پیکسل



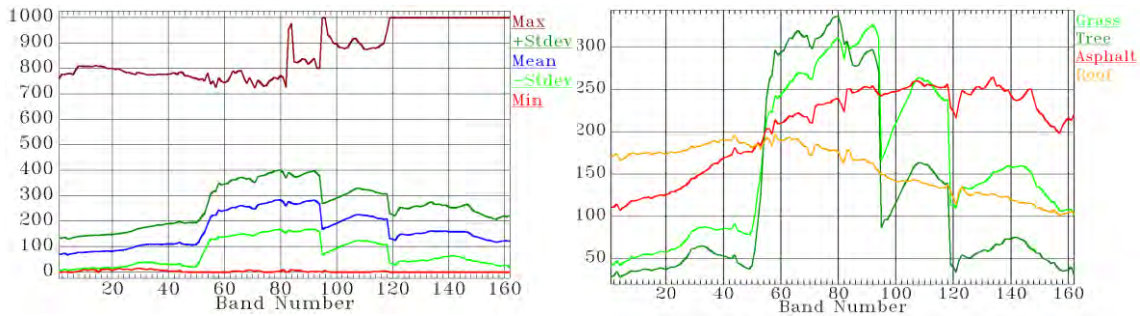
شکل ۴. مشخصات داده تصویری Botswana

باند‌های کالیبره نشده و نویزی که خصوصیات جذبی آب را پوشش می‌دهند، در این سه دسته داده تصویری فراطیفی پیش‌تر حذف شده‌اند. باند‌های حذف‌شده، شامل [تصویر Berlin: ۱، ۶۶-۶۲، ۹۶-۹۲ و ۱۲۶]، [تصویر Urban: ۴-۱، ۷۶، ۸۷، ۱۱۱-۱۰۱، ۱۵۳-۱۳۶] و [تصویر Botswana: ۱-۹، ۸۱-۵۶] و [۱۹۸-۲۱۰] هستند. به‌منظور معرفی و تشریح بیشتر این داده‌ها، اطلاعات آماری (شامل میانگین، انحراف معیار، مینیمم و ماکزیمم باند‌های طیفی) و امضای طیفی میانگین کلاس‌های نقشه واقعیت زمینی آن‌ها در شکل‌های (۵) الی (۷) نشان داده شده است.

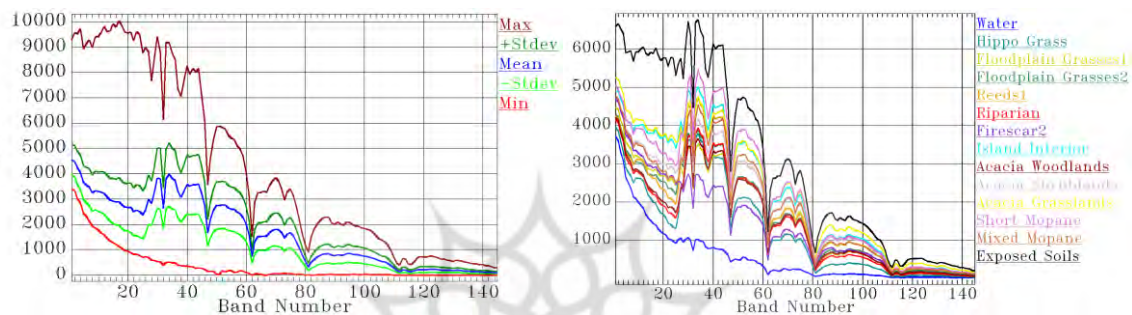
باند‌های کالیبره نشده و نویزی که خصوصیات جذبی آب را پوشش می‌دهند، در این سه دسته داده تصویری فراطیفی پیش‌تر حذف شده‌اند. باند‌های حذف‌شده، شامل [تصویر Berlin: ۱، ۶۶-۶۲، ۹۶-۹۲ و ۱۲۶]، [تصویر Urban: ۴-۱، ۷۶، ۸۷، ۱۱۱-۱۰۱، ۱۵۳-۱۳۶] و [تصویر Botswana: ۱-۹، ۸۱-۵۶] و [۱۹۸-۲۱۰] هستند. به‌منظور معرفی و تشریح بیشتر این داده‌ها، اطلاعات آماری (شامل میانگین، انحراف معیار، مینیمم و ماکزیمم باند‌های طیفی) و امضای طیفی میانگین کلاس‌های نقشه واقعیت زمینی آن‌ها در شکل‌های (۵) الی (۷) نشان داده شده است.



شکل ۵. اطلاعات آماری و امضای طیفی میانگین کلاس‌های داده تصویری Berlin



شکل ۶. اطلاعات آماری و امضای طیفی میانگین کلاس‌های داده تصویری Urban

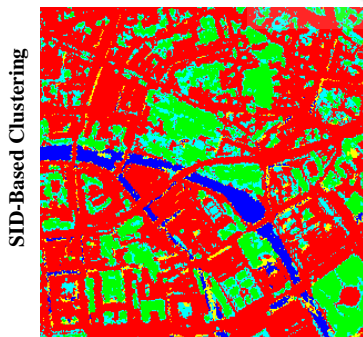
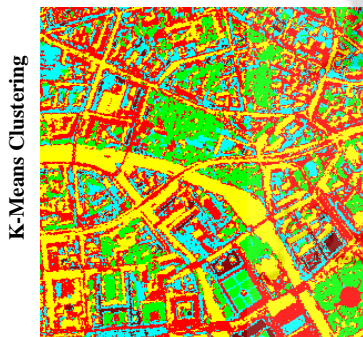


شکل ۷. اطلاعات آماری و امضای طیفی میانگین کلاس‌های داده تصویری Botswana

۲-۳- ارزیابی نتایج

با اجرای هر یک از الگوریتم‌های K-Means و SID-Based، بر روی تصاویر فراطیفی مورد استفاده، یک تصویر خوشه‌بندی تولید شد که در آن، هر پیکسل به یک خوشه اختصاص داده شده است. جهت ارزیابی کمی دقت این تصاویر خوشه‌بندی، ابتدا با شناسایی خوشه متناظر هر کلاس در نقشه واقعیت زمینی، ماتریس خطا محاسبه شده است، سپس شاخص‌های دقت طبقه‌بندی شامل دقت کلی (OA)، دقت متوسط (AA) و ضریب کاپا (KC) به دست آمده است (Jie et al., 2008; Ayday and Minz, 2014).

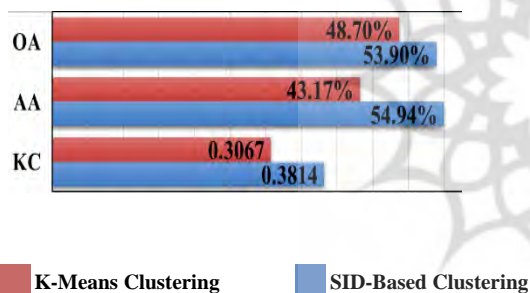
با اجرای الگوریتم‌های K-Means و SID-Based، بر روی تصویر Berlin، تصاویر خوشه‌بندی نشان داده شده در شکل (۸) به دست آمده است.



شکل ۸. خوشه‌بندی تصویر Berlin توسط دو الگوریتم SID-Based و K-Means

فاصله اقلیدسی بین آن‌ها زیاد است. در مقابل، دو امضای طیفی می‌توانند شکل متفاوت داشته ولی فاصله اقلیدسی بین آن‌ها چندان زیاد نباشد. این در حالیست که الگوریتم SID-Based، شکل امضای طیفی را مورد توجه قرار می‌دهد، زیرا دیورژانس اطلاعات طیفی، شکل امضای طیفی را مبنای اندازه‌گیری عدم شباهت قرار می‌دهد (Chen et al., 2009).

با مقایسه دیداری نتایج خوشه‌بندی در شکل (۸) و همچنین با توجه به مطالب بیان شده بالا می‌توان نتیجه گرفت که در مجموع، الگوریتم SID-Based عملکرد بهتری داشته است. شاخص‌های دقت طبقه‌بندی این نتایج که در شکل (۹) نشان داده شده است نیز، صحت این موضوع را تایید می‌کند.



شکل ۹. شاخص‌های دقت طبقه‌بندی دو الگوریتم

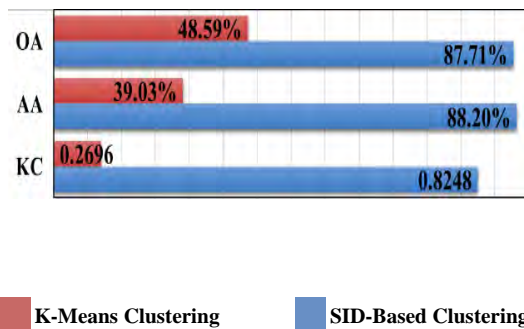
خوشه‌بندی K-Means و SID-Based برای تصویر Berlin

مقایسه شاخص‌های دقت طبقه‌بندی دو الگوریتم K-Means و SID-Based برای تصویر Berlin در شکل (۹) نشان می‌دهد که دقت کلی، دقت متوسط و ضریب کاپای الگوریتم SID-Based به ترتیب ۰/۵۲، ۰/۱۱، ۰/۷۷ و ۰/۷۴ بزرگ‌تر از الگوریتم K-Means است. الگوریتم‌های K-Means و SID-Based بر روی تصویر Urban نیز اجرا شده‌اند و تصاویر خوشه‌بندی نشان داده شده در شکل (۱۰) به دست آمده است.

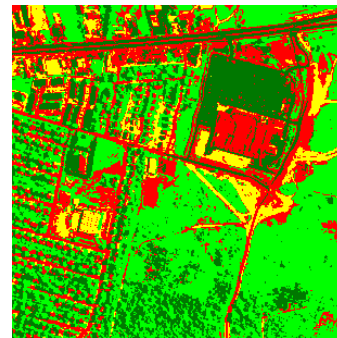
تصویر خوشه‌بندی به‌دست آمده از الگوریتم K-Means در شکل (۸) نشان می‌دهد که این الگوریتم، کلاس Vegetation را نسبتاً خوب استخراج نموده است، اما در استخراج کلاس Water ناموفق بوده است. در ادامه، به تشریح دلیل این امر و نحوه خوشه‌بندی پیکسل‌های سایر کلاس‌ها پرداخته می‌شود. بررسی انجام شده روی داده تصویر Berlin نشان می‌دهد، فاصله اقلیدسی بین امضاهای طیفی پیکسل‌های کلاس Built-up زیاد بوده است و به همین دلیل الگوریتم K-Means پیکسل‌های این کلاس در دو خوشه (نواحی قرمز و قهوه‌ای رنگ) قرار داده است. این امر، باعث شده با توجه به نزدیکی (مبتنی بر فاصله اقلیدسی) امضاهای طیفی کلاس‌ها در شکل (۵)، بخشی از پیکسل‌های کلاس Soil و حدود نیمی از پیکسل‌های کلاس Impervious به خوشه دوم کلاس Built-up (نواحی قرمز رنگ) اختصاص داده شوند و پیکسل‌های کلاس Water نیز با پیکسل‌های باقی‌مانده کلاس Impervious ادغام شده و در یک خوشه (نواحی زرد رنگ) قرار گیرند.

بررسی تصویر خوشه‌بندی حاصل از الگوریتم SID-Based در شکل (۸) نشان می‌دهد که این الگوریتم هر دو کلاس Vegetation و Water را به خوبی استخراج کرده است، اما بخش عمده‌ای از پیکسل‌های کلاس Impervious و Soil را با پیکسل‌های کلاس Built-up ادغام کرده و در یک خوشه (نواحی قرمز رنگ) قرار داده است. چون این سه کلاس از نظر طیفی به هم شبیه بوده و امضای طیفی آن‌ها به هم نزدیک است (شکل (۵)).

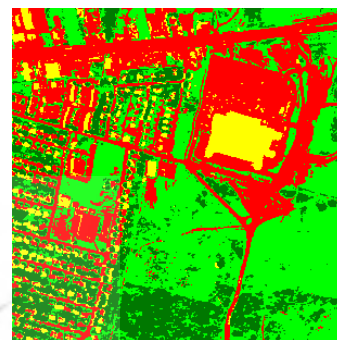
در الگوریتم خوشه‌بندی K-Means، عدم شباهت از طریق فاصله اقلیدسی اندازه‌گیری می‌شود و به شکل امضای طیفی توجهی نمی‌شود. دو امضای طیفی، می‌توانند شکل شبیه به هم داشته باشند، در حالیکه



K-Means Clustering



SID-Based Clustering



شکل ۱۱. شاخص‌های دقت طبقه‌بندی دو الگوریتم خوشه‌بندی K-Means و SID-Based برای تصویر Urban

مقایسه شاخص‌های دقت طبقه‌بندی دو الگوریتم K-Means و SID-Based در شکل (۱۱) نشان می‌دهد که بکارگیری الگوریتم SID-Based، دقت کلی، دقت متوسط و ضریب کاپا را به ترتیب ۳۹،۱۲٪، ۴۹،۱۷٪ و ۵۵،۵۲٪ برای تصویر Urban افزایش داده است.

تصویر Botswana نیز توسط الگوریتم‌های K-Means و SID-Based خوشه‌بندی شده است. پیکسل‌های پس‌زمینه داده‌های واقعیت زمینی این داده تصویری، در فرآیند خوشه‌بندی شرکت نکرده‌اند، زیرا اطلاعاتی از آن‌ها جهت ارزیابی دقت طبقه‌بندی در دسترس نبوده است. با توجه به اینکه، کلاس‌های موجود در این داده تصویری زیاد و دارای خصوصیات طیفی شبیه به هم هستند (شکل (۷))، بررسی و مقایسه دیداری نتایج حاصل از خوشه‌بندی این داده تصویری پیچیده و دشوار بوده است. به همین دلیل، نتایج حاصل از خوشه‌بندی این داده تصویری فقط به صورت کمی مورد ارزیابی گرفته است. در شکل (۱۲)، شاخص‌های دقت طبقه‌بندی نتایج حاصل از خوشه‌بندی تصویر Botswana، توسط دو الگوریتم K-Means و SID-Based نشان داده شده است.

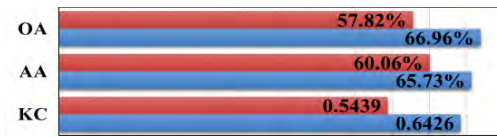
شکل ۱۰. خوشه‌بندی تصویر Urban توسط دو الگوریتم SID-Based و K-Means

تصاویر خوشه‌بندی در شکل (۱۰)، نشان می‌دهد که الگوریتم SID-Based همه کلاس‌ها را خیلی خوب استخراج کرده است اما الگوریتم K-Means فقط کلاس Grass را نسبتاً خوب استخراج کرده و در طبقه‌بندی سایر کلاس‌ها، ناموفق بوده است. با توجه به مطالبی که در ارزیابی خوشه‌بندی تصویر Berlin بیان شد، می‌توان گفت الگوریتم SID-Based در استخراج همه کلاس‌های داده تصویری Urban موفق بوده است، زیرا امضای طیفی کلاس‌های موجود در این داده تصویری دارای شکل‌های متمایزی نسبت به هم هستند (شکل (۶)). الگوریتم K-Means، به نتایج خوبی منجر نشده است، زیرا فاصله اقلیدسی بین امضای طیفی کلاس‌های موجود در داده تصویری Urban زیاد نیست. به‌منظور ارزیابی کمی این نتایج، شاخص‌های دقت طبقه‌بندی دو الگوریتم K-Means و SID-Based محاسبه شده و در شکل (۱۱) نشان داده شده است.

باید توجه داشت که هزینه زمانی در الگوریتم‌های K-Means و SID-Based، به تعداد باندهای تصویر، تعداد پیکسل‌های تصویر و تعداد خوشه‌ها وابسته است. نتایج جدول (۱) نشان می‌دهد، الگوریتم خوشه‌بندی SID-Based، به‌ویژه وقتی که تعداد خوشه‌ها زیاد است، زمان محاسباتی بیشتری را صرف می‌کند. این امر به دلیل انجام محاسبات وقت‌گیر در تعیین مقدار SID (رابطه (۱۰) و تابع Wright Omega (رابطه (۱۳) است.

۴- نتیجه‌گیری

در این مقاله، یک الگوریتم خوشه‌بندی مبتنی بر معیار SID (SID-Based)، به‌منظور خوشه‌بندی تصاویر فراطیفی پیشنهاد شد. این الگوریتم در فرآیند خوشه‌بندی، به جای فاصله اقلیدسی از معیار آماری SID، جهت اندازه‌گیری عدم شباهت استفاده می‌کند. فاصله اقلیدسی یک معیار قطعی است و امضاهای طیفی پیکسل‌ها و مراکز خوشه‌ها را به صورت نقاطی، در یک فضای n -بعدی در نظر می‌گیرد و فاصله فضایی بین آن‌ها را در آن فضا اندازه‌گیری می‌کند. معیار SID، در دسته معیارهای آماری قرار می‌گیرد و فاصله بین توزیع‌های احتمال دو طیف را اندازه‌گیری می‌کند. از آنجاکه داده‌های فراطیفی، به دلیل عواملی همچون خطا در عملکرد سنجنده، تغییرات توپوگرافی سطح زمین، سایه‌ها و شرایط گوناگون محیطی و جوی، همواره با عدم قطعیت همراهند، عدم شباهت در آن‌ها با یک معیار آماری (غیرقطعی) بهتر برآورد می‌شود. بنابراین می‌توان نتیجه گرفت اندازه‌گیری عدم شباهت، بین امضای طیفی یک پیکسل تصویر فراطیفی و امضای طیفی مرکز هر خوشه در الگوریتم خوشه‌بندی پیشنهادی SID-Based، به طور مؤثرتری انجام می‌شود.



شکل ۱۲. شاخص‌های دقت طبقه‌بندی دو الگوریتم خوشه‌بندی K-Means و SID-Based برای تصویر Botswana

با مقایسه شاخص‌های دقت طبقه‌بندی به دست آمده در شکل (۱۲) می‌توان گفت به‌کارگیری الگوریتم SID-Based دقت کلی، دقت متوسط و ضریب کاپا را به ترتیب ۹.۱۴٪، ۵.۶۷٪ و ۹.۸۷٪ برای تصویر Botswana افزایش داده است.

در ارزیابی پایانی، هزینه زمانی الگوریتم‌های خوشه‌بندی K-Means و SID-Based مورد بررسی قرار گرفته است. تمامی محاسبات و پردازش‌های این مقاله، توسط یک لپ-تاپ با پردازنده Core i7-6700HQ و رم ۱۲ گیگابایت در محیط برنامه‌نویسی MATLAB انجام شده است. همان‌طور که پیش‌تر بیان شد، الگوریتم‌های خوشه‌بندی K-Means و SID-Based به صورت تکراری اجرا می‌شوند. در جدول (۱)، هزینه زمانی هر تکرار الگوریتم‌های K-Means و SID-Based نشان داده شده است.

جدول ۱. هزینه زمانی هر تکرار الگوریتم‌های خوشه‌بندی K-Means و SID-Based

Data	Computation Time of each Iteration				
	Name	Sample Size	Bands	Clusters	K-Means Clustering
Berlin	90000	114	5	0.12 s	0.59 s
Urban	94249	162	4	0.17 s	0.73 s
Botswana	3248	145	14	0.01 s	0.98 s

- Fuzzy C-Means with Spatial Information for Clustering of Hyperspectral Images**, Journal of Basic and Applied Engineering Research, 1(7), 38-42.
- Brereton, R.G., 1992, **Multivariate Pattern Recognition in Chemometrics**, Illustrated by Case Studies (Data Handling in Science and Technology, Vol. 9), Elsevier Science.
- Chang, C.-I., 2000, **An information-theoretic approach to spectral variability, similarity, and discrimination for hyperspectral image analysis**, IEEE Transactions on Information Theory, 46(5), 1927-1932.
- Chang, C.-I., 2003, **Hyperspectral Imaging: Techniques for spectral Detection and Classification**, Springer US, New York.
- Chen, J., Jia, X., Yang, W. & Matsushita, B., 2009, **Generalization of Subpixel Analysis for Hyperspectral Data with Flexibility in Spectral Similarity Measures**, IEEE Transactions on Geoscience and Remote Sensing, 47(7), 2165-2171.
- Corless, R.M. & Jefirey, D.J., 2002, **The Wright omega Function**, Paper presented at the Artificial Intelligence, Automated Reasoning, and Symbolic Computation (Ed. J. Calmet, B. Benhamou, O. Caprotti, L. Henocque and V. Sorge), Berlin.
- Du, Y., Chang, C.-I., Ren, H., Chang, C.-C., Jensen, J.O., & D'Amico, F.M., 2004, **New** آزمون‌های انجام‌شده بر روی سه داده تصویری فراطیفی واقعی، نشان داد بکارگیری الگوریتم خوشه‌بندی SID-Based، نتایج خوشه‌بندی را به خوبی بهبود می‌بخشد. به طوری که با بررسی‌های انجام شده در بخش ارزیابی نتایج، مشخص شد الگوریتم SID-Based در مقایسه با الگوریتم K-Means، ضریب کاپای نتایج حاصل از خوشه‌بندی تصاویر فراطیفی Berlin، Urban و Botswana را به ترتیب حدود ۷٪، ۵۶٪ و ۱۰٪ افزایش می‌دهد. به عنوان تحقیقات آتی، می‌توان کارایی الگوریتم پیشنهادی SID-Based را در خوشه‌بندی سایر داده‌های چندبعدی در حوزه سنجش از دور و سیستم اطلاعات مکانی مورد ارزیابی قرار داد.
- ۵- سپاس‌گزاری**
- نویسنده بر خود لازم می‌داند از نظرها و پیشنهادهای داوران محترم نشریه علمی - پژوهشی «سنجش از دور و GIS ایران» که موجب بالابردن سطح علمی و رفع کاستی‌های این پژوهش شد، تشکر و قدرانی نماید.
- ۶- منابع**
- Adep, R.N., Vijayan, A.P., Shetty, A. & Ramesh, H., 2016, **Performance evaluation of hyperspectral classification algorithms on AVIRIS mineral data**, Perspectives in Science, 8, 722-726.
- Al-Daoud, M.B., 2007, **A New Algorithm for Cluster Initialization**, International Journal of Computer, Electrical, Automation, Control and Information Engineering, 1(4), 1031-1033.
- Aydav, P.S.S. & Minz, S., 2014, **Soft Subspace**

- Matching**, XX ISPRS Congress, Istanbul.
- Jain, A.K., 2010, **Data clustering: 50 years beyond K-means**, Pattern Recognition Letters, 31(8), 651-666.
- Jain, A.K. & Dubes, R.C., 1988, **Algorithms for clustering data**, Prentice Hall, Englewood Cliffs, New Jersey.
- Jensen, J.R., 1996, **Introductory Digital Image Processing: A Remote Sensing Perspective**, Prentice Hall, Upper Saddle River, New Jersey.
- Jie, Y., Peihuang, G., Pinxiang, C., Zhongshan, Z. & Wenbin, R., 2008, **Remote Sensing Image Classification Based on Improved Fuzzy c-Means**, Geo-spatial Information Science, 11(2), 90-94.
- MacQueen, J.B., 1967, **Some Methods for classification and Analysis of Multivariate Observations**, 5th Berkeley Symposium on Mathematical Statistics and Probability, Berkeley.
- Palsson, F., Sigurdsson, J., Sveinsson, J.R. & Ulfarsson, M. O., 2017, **Neural network hyperspectral unmixing with spectral information divergence objective**, Paper presented at the 2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS), Fort Worth.
- Shi, W., 2009, **Principles of Modeling Uncertainties in Spatial Data and Spatial Analyses**, CRC Press, Boca Raton.
- hyperspectral discrimination measure for spectral characterization**, Society of Photo-Optical Instrumentation Engineers, 43(8), 1777-1786.
- Duda, R.O., Hart, P.E. & Stork, D.G., 2001, **Pattern Classification**, Wiley, New York.
- Erudel, T., Fabre, S., Houet, T., Mazier, F., & Briottet, X., 2017, **Criteria Comparison for Classifying Peatland Vegetation Types Using In Situ Hyperspectral Measurements**, Remote Sensing, 9(7), 748-806.
- Galal, A., Hassan, H. & Imam, I.F., 2012, **A novel approach for measuring hyperspectral similarity**, Applied Soft Computing, 12(10), 3115-3123.
- Gholizadeh, H., Gamon, J.A., Zygielbaum, A.I., Wang, R., Schweiger, A.K. & Cavender-Bares, J., 2018, **Remote sensing of biodiversity: Soil correction and data dimension reduction methods improve assessment of α -diversity (species richness) in prairie ecosystems**, Remote Sensing of Environment, 206, 240-253.
- Guha, S., Rastogi, R. & Shim, K., 2001, **Cure: an efficient clustering algorithm for large databases**, Information Systems, 26(1), 35-58.
- Homayouni, S., & Roux, M., 2004, **Hyperspectral image Analysis for Material Mapping Using Spectral**

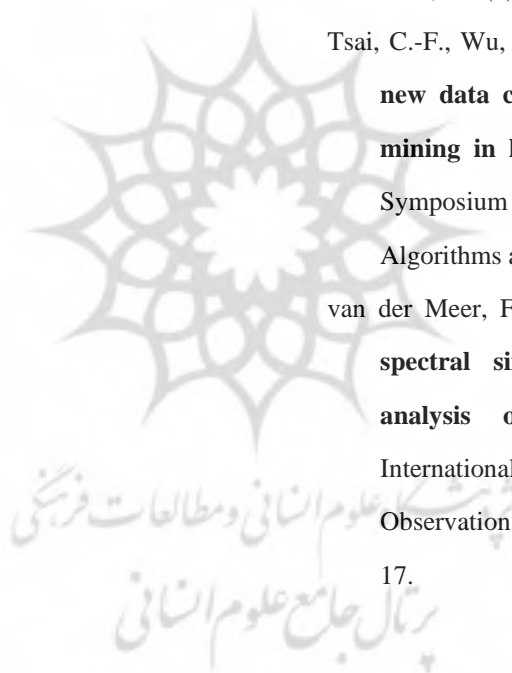
Timm, H., Borgelt, C., Döring, C., & Kruse, R.,
2004, **An extension to possibilistic fuzzy
cluster analysis**, Fuzzy Sets and Systems,
147(1), 3-16.

Tou, J.T. & Gonzalez, R.C., 1974, **Pattern
Recognition Principles**, Addison-Wesley,
Massachusetts.

Tran, T.N., Wehrens, R. & Buydens, L.M.C.,
2003, **SpaRef: a clustering algorithm for
multispectral images**, Analytica Chimica
Acta, 490(1), 303-312.

Tsai, C.-F., Wu, H.-C. & Tsai, C.-W., 2002, **A
new data clustering approach for data
mining in large databases**, International
Symposium on Parallel Architectures,
Algorithms and Networks, Makati.

van der Meer, F., 2006, **The effectiveness of
spectral similarity measures for the
analysis of hyperspectral imagery**,
International Journal of Applied Earth
Observation and Geoinformation, 8(1), 3-
17.



بهبود خوشه‌بندی تصاویر فراطیفی با به‌کارگیری دیورژانس اطلاعات طیفی





نخستین مجله از دور
,
GIS ایران



سنجش از دور و GIS ایران سال دهم، شماره سوم، پاییز ۱۳۹۷
Iranian Remote Sensing & GIS Vol.10, No.3, Autumn 2018

17-32

Improvement of Clustering for Hyperspectral Images using Spectral Information Divergence

Ezzatabadi Pour H.^{1*}

Instructor, Department of Civil Engineering, Sirjan University of Technology, Sirjan, Iran

Abstract

K-Means is one of the most frequently used unsupervised classification approaches for remotely sensed image analysis. In standard K-Means version, the Euclidean distance (ED) has used to estimate the dissimilarity between an unknown vector data and the cluster center. Since, this measure is very sensitive to topographic and environmental effects on spectral observations, we have proposed to replace it with a new one for goal of hyperspectral image clustering. The Spectral Information Divergence (SID) is a stochastic measure that is a more reliable dissimilarity measure when compared to ED as a deterministic measure. Where the ED measure the spectral distance between vector data and the clusters, SID models the probability distributions for vector data and clusters by normalizing their spectral signatures and measures the distances between them. This idea has applied to develop an enhanced clustering framework. The experimental results on three real hyperspectral images collected by HyMap, HYDICE and Hyperion sensors show that the proposed method improves classification results. In the manner that the Kappa coefficient of the classification results of three hyperspectral imagery datasets increased by about 7%, 56% and 10%, respectively.

Keywords: Clustering, Dissimilarity Measure, Spectral Information Divergence, Hyperspectral Images