

چکیده

با گسترش سیستم‌های پایگاهی و حجم بالای داده‌های ذخیره شده در آن‌ها، به ابزاری نیاز است تا بتوان این داده‌ها را پردازش کرد و اطلاعات حاصل از آن را در اختیار کاربران قرار داد. در واقع امروزه به روش‌هایی نیاز داریم که با کمترین دخالت کاربر و به صورت خودکار، الگوها و رابطه‌های منطقی را بیان نمایند. یکی از روش‌های بسیار مهمی که با آن می‌توان الگوهای مفیدی را در میان داده‌ها تشخیص داد، داده‌کاوی است. خوشه‌بندی یکی از مهم‌ترین روش‌های داده‌کاوی است. خوشه‌بندی در مسائل گوناگونی در زمینه‌های مختلف از جمله تشخیص الگو، مهندسی، پردازش تصویر و غیره مورد استفاده قرار می‌گیرد. تاکنون الگوریتم‌های مختلفی برای انجام خوشه‌بندی ارائه شده است. در این پژوهش به معرفی چند روش کارآمد در زمینه مسائل خوشه‌بندی می‌پردازیم. این روش‌ها عبارتند از :

الف) الگوریتم ژنتیک گروه‌بند

ب) الگوریتم ژنتیک با بازآرایی ژن خوشه‌بندی

الگوریتم ژنتیک گروه‌بند برگرفته از الگوریتم ژنتیک کلاسیک است که برای مسائل خوشه‌بندی توصیف می‌شود و مدل قبلی را بهبود می‌بخشد. در ادامه اجرای الگوریتم مطالعه می‌شود و تجارب به دست آمده از آزمون روی یکسری داده و نحوی مطالعه بر روی الگوریتم گروه‌بندی برای مسائل خوشه‌بندی را کامل می‌کند. در پایان یک الگوریتم پیشنهادی جهت خوشه‌بندی که کارایی بهتری نسبت به بعضی الگوریتم‌های معرفی شده دارد، ارائه می‌گردد.

فهرست مطالب

ج	مقدمه کلی
۳	فصل اول
۴	۱-۱-مقدمه
۵	۱-۲- خوشه‌بندی و طبقه‌بندی
۱۰	۱-۳-اندازه‌گیری فاصله در مسائل خوشه‌بندی
۱۳	۱-۴-الگوریتم ژنتیک
۲۷	فصل دوم
۲۸	۲-۱-مقدمه
۲۸	۲-۲-الگوریتم ژنتیک کلاسیک برای مسائل گروه‌بندی
۲۹	۲-۳-ضعف‌های الگوریتم ژنتیک کلاسیک برای مسائل گروه‌بندی
۳۰	۲-۵-الگوریتم ژنتیک گروه‌بند
۳۲	۲-۶-الگوریتم ژنتیک گروه‌بند برای مسائل خوشه‌بندی
۵۰	۲-۸- نتیجه‌گیری
۵۱	فصل سوم
۵۲	۳-۱- مقدم:
۵۳	۳-۲-تعاریف اولیه
۵۵	۳-۳- الگوریتم ژنتیک با بازآرایی ژن
۶۱	۳-۴- بررسی کارایی
۶۳	۳-۵-نتایج آزمایشات
۷۳	۳-۶-آزمایش بر روی کنترل از راه دور سنجش تصویر خوشه‌بندی
۷۷	۳-۷- روشی جدید برای بهبودالگوریتم ژنتیک گروه‌بند
۷۸	۳-۸-آزمایشات
۸۰	۳-۹- بحث و نتیجه‌گیری

منابع و مراجع.....	۸۱
واژه نامه فارسی به انگلیسی	۸۵
واژه نامه انگلیسی به فارسی	۸۶

فهرست شکل ها

- شکل ۱-۱ : خوشه‌بندی نمونه‌های ورودی ۱۱
- شکل ۲-۱ مراحل فرآیند خوشه‌بندی ۱۲
- شکل ۳-۱: کدگذاری درختی ۲۲
- شکل ۱-۲ : داده ها برای اولین مثال خوشه بندی ساختگی، داده های کروی ۴۵
- شکل ۲-۲: بهترین خوشه بندی انجام شده از طریق الگوریتم GGA با شاخص DB برای داده های کروی ۴۶
- شکل ۳-۲: داده ها برای دومین مثال خوشه بندی ساختگی، داده های ساختار یافته ۴۸
- شکل ۴-۲: بهترین خوشه بندی انجام شده از طریق الگوریتم GGA با شاخص S برای داده های ساختار یافته ۴۹
- شکل ۵-۲: داده ها برای سومین مثال خوشه بندی ساختگی، داده های نامتعادل ۵۰
- شکل ۶-۲: بهترین خوشه بندی انجام شده از طریق الگوریتم GGA با شاخص DB برای داده های نامتعادل ۵۱
- شکل ۱-۳: مثالی از انحطاط در طول تقاطع ۷۲
- شکل ۲-۳ نتایج خوشه‌بندی برای شش مجموعه داده ۷۶
- شکل ۳-۳: داده زنبق دو بعدی بعد از کاهش بعد ۷۸
- شکل ۴-۳: داده نوشیدنی دو بعدی بعد از کاهش بعد ۷۹
- شکل ۵-۳. نتایج خوشه ای از داده های زنبق با استفاده از: (الف) K- میانگین (ب) خوشه‌بندی GA (ج) خوشه‌بندی KGA (د) خوشه‌بندی GAGR ۸۰
- شکل ۶-۳. نتایج خوشه ای از داده های نوشیدنی با استفاده از: (الف) K- میانگین (ب) خوشه‌بندی GA (ج) خوشه‌بندی KGA (د) خوشه‌بندی GAGR ۸۱
- شکل ۷-۳: تصویر شبه رنگ بخش هایی از میون به دست آمده از ماهواره لندست ۷ ۸۶
- شکل ۸-۳: تصویر شبه رنگ دوم بخش هایی از میون به دست آمده از ماهواره لندست ۷ ۸۶
- شکل ۹-۳: نتایج خوشه بندی از تصویر سنجش از راه دور نشان داده شده در شکل ۷-۳ با استفاده

از: (الف) K- میانگین (ب) خوشه‌بندی GA (ج) خوشه‌بندی KGA (د) خوشه‌بندی ۸۷ GAGR

شکل ۳-۱۰: نتایج خوشه‌بندی از تصویر سنجش از راه دور نشان داده شده در شکل ۳-۸ با استفاده

از: (الف) K- میانگین (ب) خوشه‌بندی GA (ج) خوشه‌بندی KGA (د) خوشه‌بندی ۸۸ GAGR

فهرست جداول

- جدول ۱-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های کروی ۴۹
- جدول ۲-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های ساختاریافته ۵۰
- جدول ۳-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های نامتعادل ۵۱
- جدول ۴-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با شاخص های DB و S، با الگوریتم های K-mean و DBSCAN برای مجموعه داده گیاه ۵۶
- جدول ۵-۲: مقایسه نتایج به دست آمده از طریق الگوریتم های بررسی شده با الگوریتم های k-mean و DBSCAN برای مجموعه داده نوشیدنی ۵۶
- جدول ۱-۳: مثالی از توابع ۵۷
- جدول ۱-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های کروی ۶۴
- جدول ۲-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های ساختاریافته ۶۶
- جدول ۳-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با الگوریتم های K-mean و DBSCAN برای داده های نامتعادل ۶۷
- جدول ۴-۲: مقایسه نتایج به دست آمده از طریق الگوریتم GGA با شاخص های DB و S، با الگوریتم های K-mean و DBSCAN برای مجموعه داده گیاه ۶۸
- جدول ۵-۲: مقایسه نتایج به دست آمده از طریق الگوریتم های بررسی شده با الگوریتم های k-mean و DBSCAN برای مجموعه داده نوشیدنی ۷۰

مقدمه کلی

داده‌ها و الگو یکی از شاخص‌های مهم در عصر اطلاعات می‌باشند و همگام با پیشرفت در زمینه طبقه‌بندی داده‌ها و استخراج آن‌ها، روش‌های دستی تحلیل داده‌ها از رده خارج شده و استفاده از روش‌های کامپیوتری تحلیل معمول و غیرقابل انکار گردیده‌اند.

بشر به طرز ذاتی با مواجه شدن با مسائل پیچیده شروع به طبقه‌بندی مشاهدات با استفاده از ویژگی‌های آن‌ها که در علوم مختلف کاربرد دارد می‌کند. برای رسیدن به این منظور، اولین کار مورد نظر انسان طبقه‌بندی داده‌های جمع‌آوری شده و استخراج گروه‌های نمایانگر با خصوصیتی از سیستم مورد مطالعه می‌باشد. همچنین انسان‌ها بیشتر مواقع به استفاده از روش‌هایی برای نمایش روابط موجود در سیستم‌های پیچیده می‌پردازند. نمایش روابط موجود در سیستم‌های پیچیده با داشتن چندین ویژگی از آن‌ها بسیار مشکل می‌باشد. خوشه‌بندی از جمله روش‌های پرکاربرد در تجزیه و تحلیل داده‌ها است.

خوشه‌بندی شامل روش‌های بسیار متنوعی است که علیرغم کاربرد در بسیاری از علوم، روش‌ها از نظر هدف، الگوریتم و نتایج با یکدیگر یکسان نیستند. مقالات و کتاب‌های زیادی در ارتباط با روش‌های خوشه‌بندی نوشته شده است، اما هنوز در انتخاب صحیح روش خوشه‌بندی ملاک روشنی وجود ندارد و عدم انتخاب روش صحیح از بین روش‌های موجود، استخراج نادرست الگوها را موجب می‌گردد. خوشه‌بندی، تقسیم‌بندی الگوریتم‌های خوشه‌بندی، طبقه‌بندی، مراحل خوشه‌بندی و معیارهای ارزیابی خوشه‌بندی و الگوریتم ژنتیک در فصل ۱ به طور مفصل تشریح شده است. یکی دیگر از مسائلی که در این پایان‌نامه مورد بررسی می‌گیرد، مسائل گروه‌بندی است. الگوریتم ژنتیک برای مسائل گروه‌بندی مشابه الگوریتم ژنتیک معمولی است و تنها نحوه کدگذاری آن به شکلی است که برای مسائل گروه‌بندی مناسب می‌باشد. با استفاده از الگوریتم ژنتیک برای مسائل گروه‌بندی، امانوئل فالکنایر^۱ متوجه شد که یکسری از اشکالات قابل توجه در الگوریتم ژنتیک برای مسائل گروه‌بندی وجود دارد.

^۱Emanuel Falkenauer

فالکنایر در مقاله‌اش با عنوان مسائل گروه‌بندی و الگوریتم‌های ژنتیک [۱] وارد جزئیاتی از نقص‌های ناشی از اعمال الگوریتم ژنتیک کلاسیک بر روی مسائل گروه‌بندی، شده است. به بیان او کدگذاری بدون واسطه، اثر بسیار نامطلوبی بر روی اطلاعات دارد. همچنین فرآیند ترکیب یا تقاطع یک کروموزوم می‌تواند به راحتی به راه‌حلی برسد که از هیچ یک از صفات والدین استفاده نمی‌کند. به طور مشابه مسائل در رابطه با به کار بردن طرح کدگذاری استاندارد و عملگر جهش وجود دارد. در ادامه الگوریتم ژنتیک گروه‌بند برای مسائل خوشه‌بندی معرفی و سپس به آن پرداخته می‌شود. در ادامه رفتار دو الگوریتم k - میانگین و خوشه‌بندی مبتنی بر چگالی را با الگوریتم GGA مقایسه می‌کنیم. نتایج و آزمایشات در انتهای فصل ۲ آمده است.

در فصل ۳ ضمن ارائه تعاریف مورد نیاز، الگوریتم ژنتیک برای خوشه‌بندی با روش k - میانگین با استفاده از بازآرایی ژن معرفی می‌گردد. سپس به منظور آزمون عملکرد الگوریتم خوشه‌بندی با استفاده از بازآرایی ژن، آزمایش‌ها بر روی داده‌ها در دنیای واقعی از UCI تصاویر سنجش از راه دور، انجام و نتایج نشان داده خواهد شد. نهایتاً در این بخش، اجرای خوشه‌بندی با GAGR، خوشه‌بندی با GA، خوشه‌بندی با KGA و خوشه‌بندی با الگوریتم k - میانگین با آزمایشاتی که بر مبنای شش مجموعه داده‌ها در دنیای واقعی است، مقایسه می‌گردند. سپس روشی جدید برای بهبود الگوریتم ژنتیک گروه‌بند ارائه و نتایج و آزمایشات که بیانگر بهبود الگوریتم هستند، ارائه می‌گردد.

فصل اول

مقدمات و پیش‌نیازها

۱-۱-۱- مقدمه

جامعه‌ی مبتنی بر اطلاعات را می‌توان به عنوان جامعه‌ای تعریف نمود که بخش غالب اجتماع به جای کارهای فیزیکی در گیرکارهای فکری هستند. در چنین جامعه‌ای بیشترین توجه به فعالیت‌های اطلاعاتی از قبیل: فراهم‌آوری، پردازش، تولید، ثبت، انتقال، اشاعه و مدیریت اطلاعات مبذول می‌گردد و بیشترین هزینه‌ها صرف فرایندهای اطلاعاتی می‌شود.

با گسترش سیستم‌های پایگاهی و حجم بالای داده‌های ذخیره شده در این سیستم‌ها، به ابزاری نیازاست تا بتوان این داده‌ها را پردازش کرد و اطلاعات حاصل از آن را در اختیار کاربران قرار داد. در واقع امروزه به روش‌هایی نیاز داریم که به اصطلاح به کشف دانش^۱ بپردازند. یعنی روش‌هایی که با کمترین دخالت کاربر و به صورت خودکار، الگوها و رابطه‌های منطقی را بیان نمایند. یکی از روش‌های بسیار مهمی که با آن می‌توان الگوهای مفیدی را در میان داده‌ها تشخیص داد، داده‌کاوی است [۲].

۱-۱-۱-۱ داده‌کاوی

از روش‌های بسیار مهم تشخیص الگوهای مفید در میان داده‌ها، داده‌کاوی است. این روش که با حداقل دخالت کاربران همراه است اطلاعاتی را در اختیار آن‌ها و تحلیل‌گران قرار می‌دهد تا براساس آن‌ها تصمیمات مهم و حیاتی در سازمانشان اتخاذ نمایند.

باید توجه داشت که اصطلاح داده‌کاوی زمانی به کار برده می‌شود که با حجم بزرگی از داده‌ها، در حد گیگا یا ترابایت، مواجه باشیم. در تمامی منابع داده‌کاوی بر این مطلب تاکید شده است. هرچه حجم داده‌ها بیشتر و روابط میان آن‌ها پیچیده‌تر باشد، دسترسی به اطلاعات نهفته در میان داده‌ها مشکل‌تر می‌شود و نقش داده‌کاوی به عنوان یکی از روش‌های کشف دانش، آشکارتر می‌گردد [۳].

^۱ Knowledge Discovery

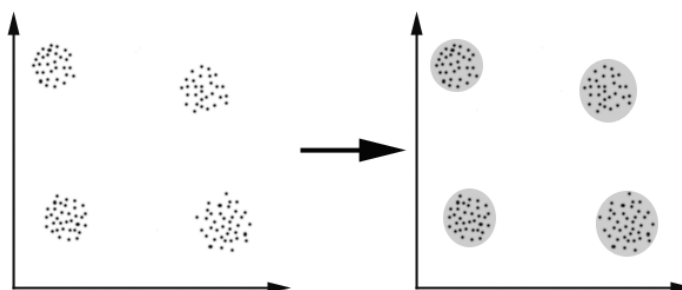
۱-۲- خوشه‌بندی و طبقه‌بندی

۱-۲-۱ خوشه‌بندی

خوشه‌بندی از مهم‌ترین روش‌های داده‌کاوی و یک زیرگروه مهم از تکنیک‌های یادگیری ماشین بدون نظارت است که شامل گروه‌بندی داده‌ها به خوشه‌های جدا از هم است. به صورت کلی به گروه‌بندی اشیاء داده‌ای به گروه‌های مجزا (خوشه)، خوشه‌بندی می‌گویند. خوشه‌بندی در مسائل گوناگونی در زمینه‌های مختلف از جمله تشخیص الگو، مهندسی، پردازش تصویر و غیره مورد استفاده قرار می‌گیرد. دلیل اصلی استفاده از خوشه‌بندی این است که بعضی مواقع اطلاعات خاصی در مورد ویژگی‌ها و کیفیت خوشه‌ها به علت پیچیدگی پایگاه داده در دسترس نمی‌باشد. در نتیجه بایستی از روشی که توانایی یافتن الگوها از فضای جستجو را دارد، استفاده کنیم [۴، ۵].

در تعریفی دیگر به سازماندهی داده‌ها در خوشه‌ها، تحلیل خوشه‌بندی می‌گویند. تحلیل خوشه‌بندی به معنی سازماندهی الگوهای جمع‌آوری شده است. سازماندهی و تحلیل بر این مبنا صورت می‌گیرد که داده‌های هر خوشه حداکثر درجه شباهت و داده‌های متعلق به خوشه‌های مختلف دارای حداکثر درجه عدم شباهت باشند. برای مشابه بودن می‌توان معیارهای مختلفی را در نظر گرفت. مثلاً می‌توان معیار فاصله را برای خوشه‌بندی مورد استفاده قرار داد و اشیائی را که به یکدیگر نزدیک‌تر هستند را به عنوان یک خوشه در نظر گرفت که به این نوع خوشه‌بندی، خوشه‌بندی مبتنی بر فاصله^۱ نیز گفته می‌شود. به عنوان مثال در شکل ۱-۱ نمونه‌های ورودی در سمت چپ به چهار خوشه مشابه شکل سمت راست تقسیم می‌شوند. در این مثال هر یک از نمونه‌های ورودی به یکی از خوشه‌ها تعلق دارد و نمونه‌ای وجود ندارد که متعلق به بیش از یک خوشه باشد.

^۱ Distance-based Clustering



شکل ۱-۱: خوشه‌بندی نمونه‌های ورودی

هر داده نشان دهنده برداری از مقادیر کمی در یک فضای چند بعدی است.

روش‌های تحلیل به دو گروه عمده سلسله مراتبی و تفکیکی تقسیم می‌گردند:

۱. خوشه‌بندی سلسله مراتبی عبارت است از تشکیل متوالی گروه‌هایی که اعضای آنها بیشترین شباهت را باهم دارند، یا جداسازی متوالی گروه‌هایی که اعضای آنها دارای بیشترین اختلاف با هم هستند. در این نوع خوشه‌بندی مشکلات ناشی از مقداردهی اولیه و کمینه‌های محلی وجود ندارد.

۲. خوشه بندی تفکیکی، شامل به دست آوردن تفکیکی از داده‌های ورودی در تعداد مشخصی از خوشه‌ها می‌باشد. چنین روش‌هایی معمولاً به دنبال جداسازی هستند که تابع شایستگی را به صورت محلی بهینه کند. برای بهبود کیفیت خوشه‌بندی الگوریتم چندین بار و در نقاط شروع مختلف اجرا می‌شود و بهترین وضعیت بدست‌آمده از کل دفعات اجرا به عنوان خروجی خوشه‌بندی انتخاب می‌گردد [۷،۶].

۲-۲-۱-۲ تقسیم‌بندی الگوریتم‌های خوشه‌بندی

روش‌های مختلفی برای تقسیم بندی الگوریتم‌های خوشه‌بندی وجود دارد، اما در اینجا با یک دیدگاه متفاوت آن‌ها را به چهار گروه تقسیم می‌کنیم:

۱. دسته اول یک مفهوم کلی از خوشه‌بندی را بر اساس این نظر که داده‌های نزدیک به هم باید در یک خوشه قرار گیرند به کار می‌برد. به عنوان نمونه از الگوریتم‌هایی که بر اساس

این اصل پیاده‌سازی می‌شوند می‌توان موارد زیر را نام برد: روش‌های بر پایه‌ی چگالی [۸، ۹، ۱۰]، روش‌های نزدیکترین هم‌سایه [۱۱]، خوشه‌بندی تراکمی اتصال تک [۱۲]. این روش‌ها به راحتی هر شکلی از خوشه‌بندی را انجام می‌دهند، اما وقتی فاصله مجزا بین خوشه‌ها کم است، در تعیین خوشه‌ها دچار مشکل می‌شوند.

۲. دسته دوم با توجه به متغیر داخل خوشه‌ها، جواب نهایی را به دست می‌آورند. این دسته شامل الگوریتم‌هایی مثل k -میانگین [۱۳، ۱۴]، خوشه‌بندی یادگیری شبکه [۱۵، ۱۶] و خوشه‌بندی بر اساس مدل [۱۷] هستند. این روش‌ها برای خوشه‌های کاملاً مجزا و کروی موثرتر هستند، اما برای خوشه‌هایی با ساختار پیچیده‌تر، ممکن است منجر به شکست شوند.

۳. دسته سوم خوشه‌بندی سطر و ستون ماتریس را به صورت همزمان انجام می‌دهند. از این نوع می‌توان الگوریتم‌های خوشه‌بندی دوگانه [۱۸] را نام برد. هدف این روش این است که با اجرای خوشه‌بندی همزمان سطر و ستون ماتریس‌ها، به جای خوشه‌بندی مجزا برای سطر و ستون، زیرگروه‌هایی از سطرها و زیرگروه‌هایی از ستون‌ها را بشناسند. پس خوشه‌بندی دوگانه مدل‌های محلی را تولید می‌کند، در حالی که این روش خوشه‌بندی مدل‌های کلی را محاسبه می‌کند. خوشه‌های معرفی شده توسط این الگوریتم‌ها دو به دو ناسازگار یا جامع نیستند. یک داده ممکن است متعلق به هیچ خوشه‌ای نباشد یا در بیشتر از یک خوشه باشد.

۴. دسته چهارم ویژگی‌های مختلف مجموعه داده را بهینه می‌کند. این الگوریتم‌ها می‌توانند به خوشه‌بندی دسته جمعی تقسیم شوند [۱۹] که نتیجه را به یک موردی با کیفیت بهتر ترکیب می‌کنند و روش‌های خوشه‌بندی چند منظوره [۱۸] که یک تخمینی از کیفیت همه حل‌های خوشه‌بندی فراهم می‌کنند. اگرچه اغلب الگوریتم‌ها قوی هستند اما بدون اشکال هم نیستند.

یک خوشه‌بندی خوب باید هم توسط الگوریتم تولید شده و هم یک معیار خارجی تایید شود. از بین الگوریتم‌های دسته دوم، الگوریتم k -میانگین کاربرد بیشتری دارد. مراکز اولیه که برای

خوشه‌ها تعیین می‌شود، بر عملکرد این الگوریتم تاثیر می‌گذارد و براحتی راه‌حل‌های بهینه محلی را ارائه می‌کند. در این الگوریتم زمانی که داده‌ها بزرگ هستند رسیدن به جواب به زمان زیادی نیاز دارد.

۱-۲-۳ طبقه‌بندی

طبقه‌بندی در واقع ارزشیابی ویژگی‌های مجموعه‌ای از داده‌ها و سپس اختصاص دادن آن‌ها به مجموعه‌ای از گروه‌های از پیش تعریف شده است. مسائل طبقه‌بندی به شناسایی خصوصیتی منجر می‌شوند که مشخص می‌نمایند هر مورد به کدام گروه تعلق دارد. این الگو هم می‌تواند برای فهم داده موجود و هم برای پیش‌بینی این‌که هر نمونه جدید چگونه کار می‌کند استفاده شود.

۱-۲-۴ طبقه‌بندی در مقابل خوشه‌بندی

درک تفاوت بین خوشه‌بندی و طبقه‌بندی بسیار مهم می‌باشد. در طبقه‌بندی، همه دسته‌ها و ویژگی‌های آنها از همان ابتدا مشخص است و برای هر شیئی تنها کافیست بررسی کنیم که دارای ویژگی‌های کدام طبقه می‌باشد، اما در خوشه‌بندی هیچ اطلاعاتی در مورد خوشه‌ها نداریم و هنگام بررسی اشیاء، بایستی مدلی برای شباهت بین داده‌ها استخراج گردد. در واقع در اغلب الگوریتم‌های خوشه‌بندی در ابتدا بایستی راهی برای درک الگوی مابین داده‌ها به دست آورد.

۱-۲-۵ هدف از خوشه‌بندی چیست؟

هدف از خوشه‌بندی یافتن خوشه‌های مشابه از اشیاء در بین نمونه‌های ورودی می‌باشد، اما چگونه می‌توان گفت که یک خوشه‌بندی مناسب است و دیگری مناسب نیست؟ می‌توان نشان داد که هیچ معیار مطلقاً برای بهترین خوشه‌بندی وجود ندارد، بلکه این بستگی به مساله و نظر کاربر دارد که باید تصمیم بگیرد که آیا نمونه‌ها به درستی خوشه‌بندی شده‌اند یا خیر. با این حال معیارهای مختلفی برای خوب بودن یک خوشه‌بندی ارائه شده است که می‌تواند کاربر را برای رسیدن به یک خوشه‌بندی مناسب راهنمایی کند که در بخش‌های بعدی چند نمونه از این معیارها

آورده شده است. یکی از مسایل مهم در خوشه‌بندی انتخاب تعداد خوشه‌ها می باشد. در بعضی از الگوریتم‌ها تعداد خوشه‌ها از قبل مشخص شده است و در بعضی دیگر خود الگوریتم تصمیم می گیرد که داده‌ها به چند خوشه تقسیم شوند.

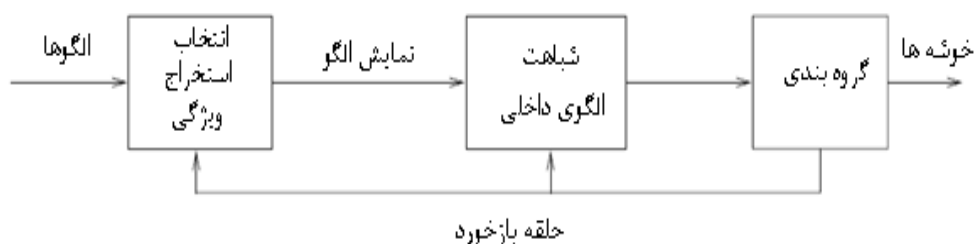
۱-۲-۶ مراحل خوشه‌بندی

مراحل خوشه‌بندی به صورت زیر می باشد:

- ۱- نمایش الگو که شامل استخراج و انتخاب ویژگی است.
- ۲- تعریف از اندازه‌گیری شباهت الگوها برای دامنه داده‌ها
- ۳- خوشه‌بندی یا گروه‌بندی
- ۴- انتزاع داده (در صورت نیاز)
- ۵- ارزیابی خروجی (در صورت نیاز)

شکل زیر ترتیب سه مرحله اولیه خوشه‌بندی را نشان می‌دهد که شامل بازخوردی از مرحله

گروه‌بندی به مراحل قبلی می‌باشد.



شکل ۱-۲: مراحل فرآیند خوشه‌بندی

۱. نمایش الگو: شامل مشخص کردن تعداد کلاس‌ها و الگوهای در دسترس، انواع و میزان اهمیت ویژگی‌ها برای خوشه‌بندی است.
۲. انتخاب ویژگی: فرآیندی است که تاثیرگذارترین زیر مجموعه‌ها از ویژگی‌های اصلی را برای خوشه‌بندی استفاده می‌نماید.
۳. شباهت الگو: معمولاً به وسیله تابع فاصله برای هر زوج الگو، میزان شباهت، اندازه‌گیری می‌شود. معیارهای مختلفی برای اندازه‌گیری فاصله میان الگوها استفاده می‌شود که

معروف‌ترین آن فاصله اقلیدسی است.

۴. خوشه‌بندی: بر اساس روش‌های مختلفی انجام می‌گیرد. مرزبندی میان خوشه‌ها می‌تواند سخت و یا به صورت فازی به شکلی که هر الگو درجه عضویت متفاوتی نسبت به هر گروه داشته باشد، بررسی گردد.
۵. انتزاع داده: فرآیندی است که یک نمایش ساده و فشرده از مجموعه داده‌ها، را استخراج می‌نماید. در واقع در خوشه‌بندی، انتزاع داده توصیف فشرده‌ای از هر خوشه است.
۶. تحلیل خوشه‌بندی: به دنبال سازمان‌دهی مجموعه‌ای از داده‌ها در یک سری خوشه اتفاق می‌افتد [۲۰].

۱-۲-۷ معیارهای ارزیابی خوشه‌بندی

برای استفاده از خوشه‌بندی در حل مسائل، بایستی آن را در دنیای ریاضی بررسی نمائیم. به همین دلیل در این بخش، نحوه تعریف خوشه‌بندی در ریاضیات مورد توجه قرار گرفته است. در ریاضیات، یک مساله خوشه‌بندی به صورت ترکیبی از عناصر $X = \{x_1, \dots, x_n\}$ که نشان دهنده مجموعه‌ای از n بردار در فضای ویژگی S است، در نظر گرفته می‌شود. هدف یک مساله خوشه‌بندی، یافتن یک افراز بهینه $U^* = \{C_1^*, \dots, C_K^*\}$ ، $C_i \cap C_j = \emptyset$ برای مجموعه X است، که C_i^* نمادی برای i امین خوشه از افراز U^* می‌باشد. برای هر مساله تابع هدفی تعریف و $m(u)$ نامیده می‌شود که برحسب آن، تا حد امکان الگوهای مربوط به یک خوشه، یکسان و الگوهای مربوط به خوشه‌های مختلف، متفاوت می‌باشد.

۱-۳-اندازه‌گیری فاصله در مسائل خوشه‌بندی

در خوشه‌بندی، اندازه‌گیری فواصل یکی از مسائل کلیدی است، زیرا شباهت بین دو بردار مختلف x_i و x_j از طریق میزان فاصله‌شان از یکدیگر در فضای ویژگی S سنجیده می‌شود. اغلب رایج است که برای محاسبه فاصله از نرم اقلیدسی، رابطه زیر استفاده می‌گردد: