

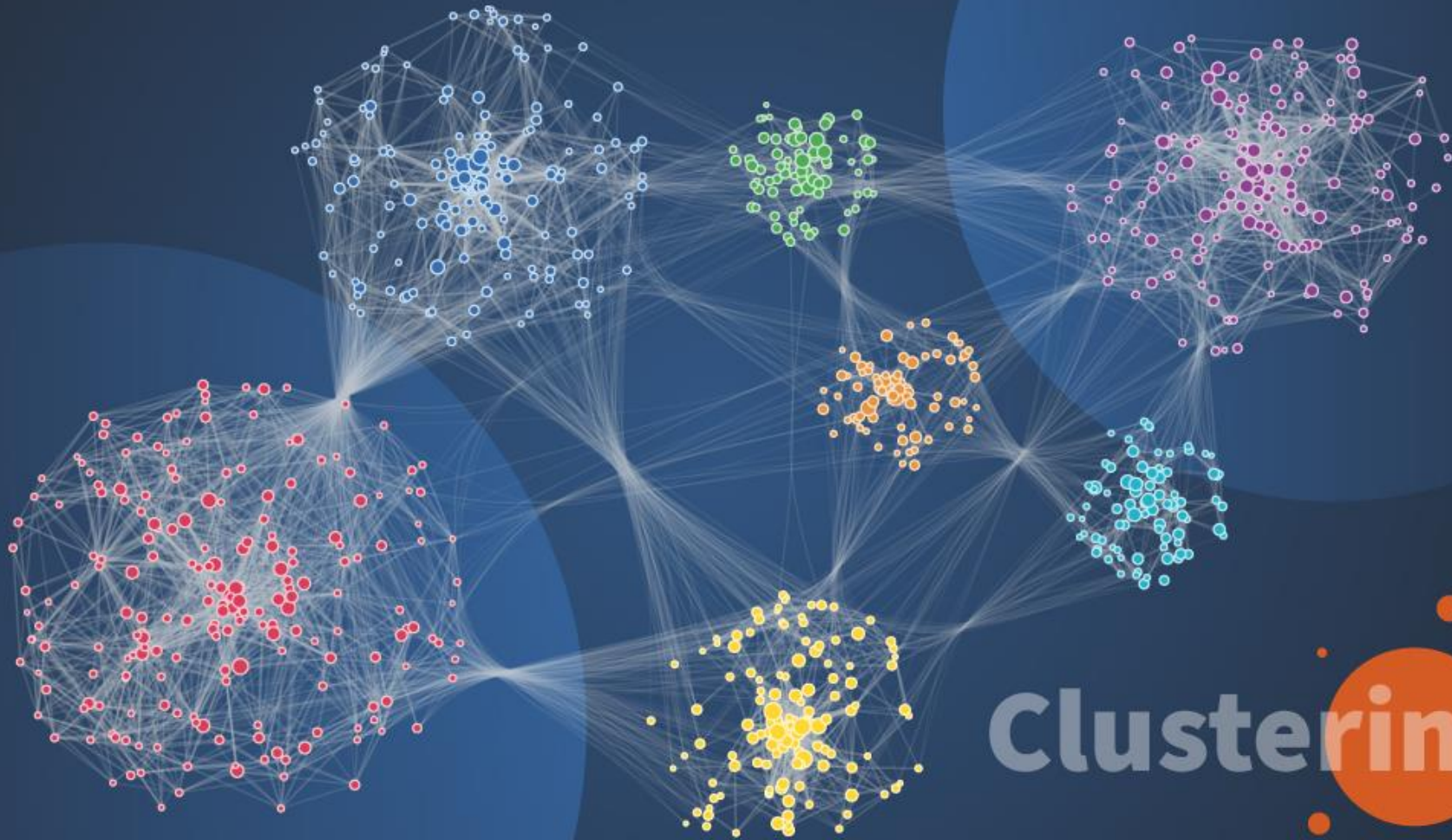
یادگیری بدون ناظر

تحلیل خوشه‌ای

یادگیری بدون ناظر:

- اگر در پایگاه داده ها، متغیر (یا فیلد) بر چسب دار وجود نداشته باشد برای مدلسازی این نوع از پایگاه داده ها فقط می توان از روش های یادگیری بدون ناظر استفاده کرد.
- از جمله این روش ها می توان به **تحلیل خوشه ای** و روش های پیشرفته آن اشاره کرد.

Data Mining



Clustering

آنالیز خوشه‌بندی

- آنالیز خوشه‌بندی یک روش قدیمی برای دسته‌بندی داده‌ها است که اولین بار توسط تریون در سال 1935 پیشنهاد شد.
- امروزه به دلیل استفاده از ابر کامپیوترها برای محاسبات سنگین، پیشرفت‌های زیاد و قابل توجهی برای آن صورت گرفته است.
- تفاوت روش خوشه‌بندی با دسته‌بندی در آن است که در این روش متغیر هدف و یا متغیر برچسب‌دار نداریم، لذا این روش تحلیل داده در دسته روش‌های بدون ناظر قرار می‌گیرد.
- به عبارت دیگر متغیرهای موجود به دو دسته مستقل و وابسته تقسیم نمی‌شوند.
- **هدف اصلی در خوشه‌بندی** تقسیم بندی اشیاء به گونه ای است که بیشترین شباهت در یک گروه و بیشترین تفاوت با اشیاء گروه های دیگر را دارا باشد.

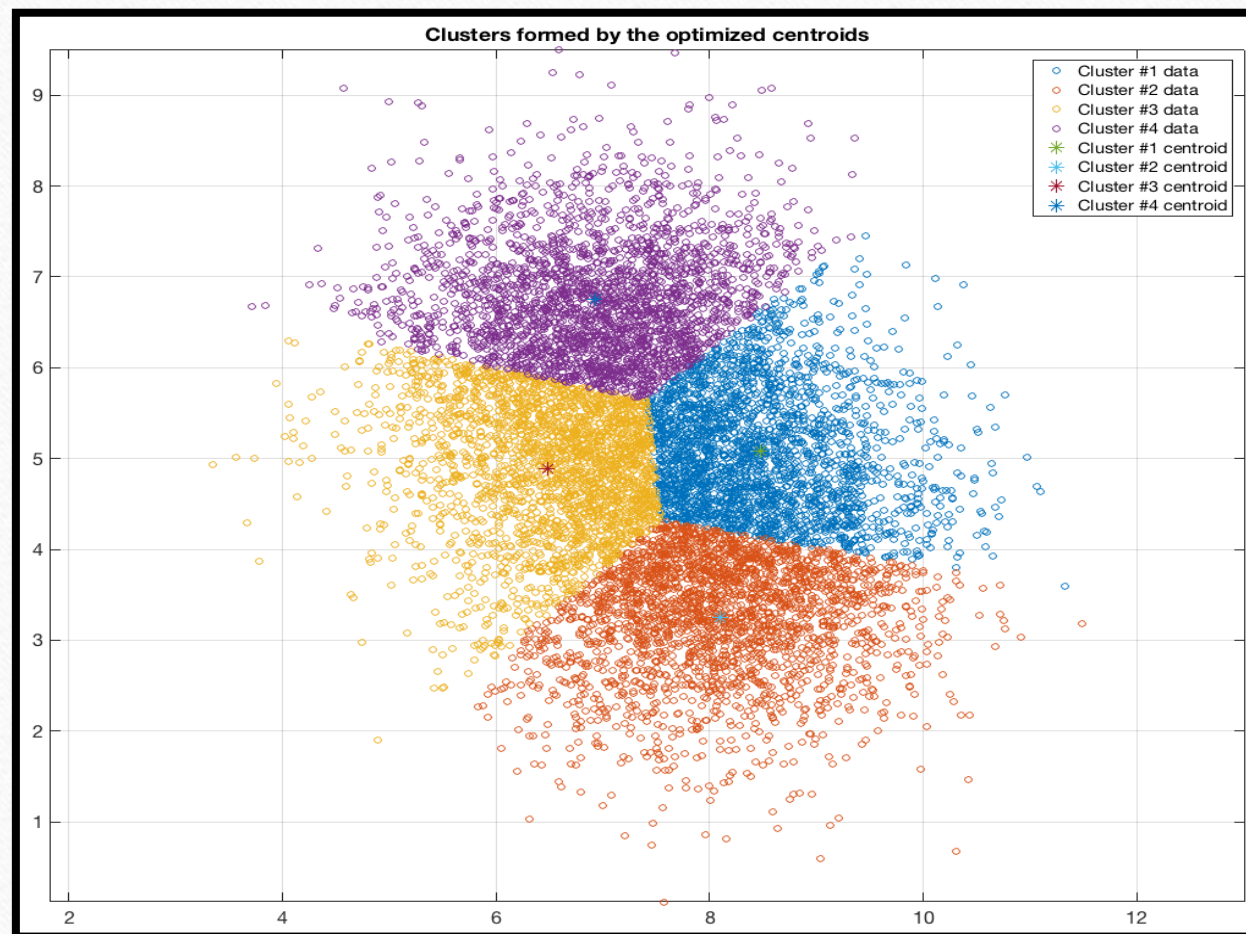
- تحلیل خوشه‌ای فرآیند تقسیم مجموعه‌ای از داده‌ها یا مشاهدات در زیر مجموعه‌ها می‌باشد، **هر زیر مجموعه یک خوشه است**. همه اعضای هر خوشه **مشابه** یکدیگرند و با اعضای دیگر خوشه‌ها **تفاوت** دارند. امکان این تقسیم‌بندی توسط انسان‌ها وجود ندارد و باید توسط الگوریتم‌های خوشه‌بندی صورت گیرد.
- از روش تحلیل خوشه‌ای برای شناخت گروه‌های ناشناخته در داده‌ها استفاده می‌شود. از روش تحلیل خوشه‌ای در زمینه‌های مختلف از جمله هوش تجاری، روش‌های جستجو اینترنتی، زیست‌شناسی، ... استفاده می‌شود.
- **مثال:** در یک شرکت مشاوره‌ای با خوشه‌بندی پروژه‌ها می‌توان آن دسته از پروژه‌هایی که از نظر ویژگی‌های تعریف شده برای پروژه‌ها، شباهت دارند در یک خوشه قرار گیرند.

-
- تحلیل خوشه‌ای یکی از روش‌های آماری است که برای کاهش داده‌ها و پیدا کردن گروه‌های مفید بسیار کاربرد دارد.
 - توجه شود در خوشه‌بندی، داده‌ها در گروه‌هایی قرار می‌گیرند که از نظر محقق در آن زمینه مشابه باشند. اگر زاویه دید محقق تغییر کند، گروه‌ها در خوشه‌بندی تغییر می‌کنند.
 - یکی از دلایل استفاده از خوشه‌بندی کشف ساختارهای جدید یا الگوهای جدید می‌باشد که به صورت طبیعی در داده‌ها وجود دارد اما بایستی توسط تحلیل خوشه‌ای کشف شود.

(Clustering) یا خوشه بندی از جمله الگوریتم های یاد گیری مدل به حساب می آید. الگوریتم خوشه بندی رکوردهایی را که ویژگی های نزدیک به هم و مشابه دارند را در یک خوشه قرار می دهد. وقتی با یک مجموعه کوچک از صفات روبرو باشیم خوشه بندی به سادگی قابل اجرا است

✓ برای مثال در یک مجموعه از خودکارهای آبی، مشکی، قرمز و سبز به راحتی می توانیم آنها را در 4 دسته قرار دهیم اما اگر در همین مجموعه ویژگی های دیگری مثل سایز، شرکت سازنده، وزن، قیمت و... مطرح باشد کار کمی پیچیده می شود. حال فرض کنید در یک مجموعه متشکل از هزاران رکورد و صدها ویژگی قصد دسته بندی دارید، چگونه باید این کار را انجام دهید؟

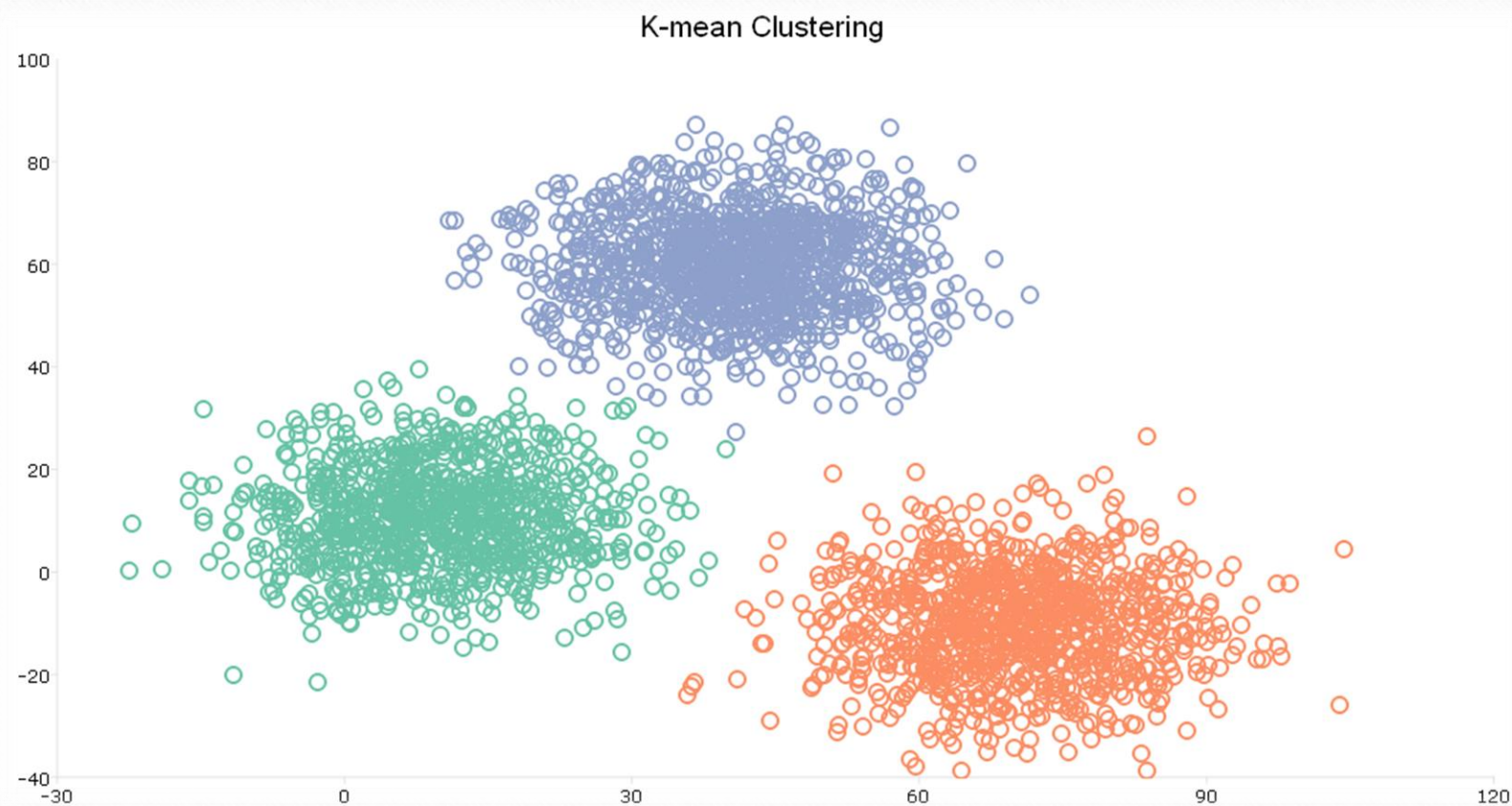
خوشه بندي در واقع تقسيم بندي يك جمعيت ناهمگون به تعدادي زیر مجموعه همگون است.



- تقسیم داده‌ها به گروه‌ها یا خوشه‌های معنادار به طوری که محتویات هر خوشه ویژگی‌های مشابه و در عین حال نسبت به اشیاء دیگر در سایر خوشه‌ها غیر مشابه باشند را خوشه بندی می‌گویند. از این الگوریتم در مجموعه داده‌های بزرگ و در مواردی که تعداد ویژگی‌های داده زیاد باشد استفاده می‌شود.
- زمانی که با یک مجموعه کوچک از خصوصیات سروکار داریم، خوشه‌بندی یک عمل ساده‌ای است که می‌توانیم آن را انجام دهیم. اما زمانی که خصوصیات رشد می‌کنند مشکلات خوشه بندی افزایش پیدا کرده و حتی ممکن است از طریق ذهن آدمی غیر ممکن باشد. عمل خوشه بندی فقط در 5 یا 6 ویژگی برای افراد نظریه پرداز که فهم عمیقی از داده ها دارند امکان پذیر است. اما مجموعه داده‌های مدرن، عموماً شامل ده ها (اگر نگوییم صدها) ویژگی هستند.

-
- خوشه بندی، اشیاء را بر اساس ویژگی‌هایی که با هم دارند گروه‌بندی می‌کند.
 - هدف اصلی در خوشه بندی تقسیم بندی اشیاء به گونه‌ای است که بیشترین شباهت در یک گروه و بیشترین تفاوت با اشیاء گروه‌های دیگر را دارا باشد.
 - به عنوان تعریف ساده‌تر می‌توان گفت که اشیاء در خوشه مخصوص خود دارای بیشترین شباهت و در برابر اشیای متعلق به خوشه‌های دیگر دارای بیشترین تفاوت هستند.
 - در روش خوشه‌بندی هیچ خوشه‌ای از ابتدا مشخص نیست. به عبارتی متغیرهای مورد نظر در ابتدا کاملاً از هم مستقل فرض می‌شوند.
 - یکی از کاربردهای خوشه‌بندی **تشخیص نقاط پرت** بعد از تشکیل خوشه‌هاست.

در خوشه‌بندی، هدف تقسیم داده‌ها به گروه‌های مختلف است که با رنگ‌های مختلف در اینجا نشان داده شده‌اند.



- در این قسمت تفاوت تجزیه و تحلیل خوشه ای با دیگر تکنیک های گروه بندی را بیان می کنیم. در واقع گروه بندی اشیاء یا داده ها را به کلاس های نام گذاری شده تخصیص می دهد، در دسته بندی هر شیء دارای یک سرپرست و یا ناظر می باشد ولی می توان گفت خوشه بندی یک نوع گروه بندی بدون سرپرست یا ناظر است یعنی اشیاء بر اساس شباهت هایی که باهم دارند تقسیم می شوند و نه بر اساس معیارهای از پیش تعیین شده.
- به همین دلیل به خوشه بندی گاهی اوقات دسته بندی بدون ناظر نیز می گویند.
- در داده کاوی هنگامی که از اصطلاح دسته بندی استفاده می شود منظور همان دسته بندی به صورت با ناظر می باشد. در واقع استفاده از دسته بندی، استفاده از تکنیک های ساده ای برای گروه بندی داده ها می باشد.



الف) خوشه بندی با دو خوشه



ج) خوشه بندی با شش خوشه



ب) خوشه بندی با چهار خوشه

کاربرد های خوشه بندی:

- تشخیص داده های دور افتاده (مقادیری که با خوشه ها فاصله زیاد دارند).
- گروه بندی داده ها برای پیش بینی

خوشه بندی شاخه ای از علم آمار است، که یکی از موضوعات بسیار پرکاربرد در داده کاوی می باشد.

مراحل خوشه‌بندی

- انتخاب معیار مناسب برای تشخیص شباهت‌ها یا نزدیکی داده‌ها
- انتخاب روش تجزیه و تحلیل خوشه‌ها
- تصمیم‌گیری در مورد تعداد خوشه‌ها
- تفسیر نتایج خوشه‌ها و یا گروه‌ها

• از جمله روش‌های خوشه‌بندی می‌توان به موارد زیر اشاره کرد:

1. روش افزایشی
2. روش سلسله مراتبی
3. روش‌های مبتنی بر تابع چگالی
4. روش‌های مبتنی بر مدل
5. روش Grid-based

روش افراز بندی

تحلیل خوشه‌ای

روش افرازبندی

- در این روش k افراز (k خوشه) مجزا درست می‌شود.
- هر افراز یک خوشه را درست میکند یعنی مشاهدات (n مشاهده) در این k خوشه قرار می‌گیرند. به طوریکه:

I. هر گروه یا هر خوشه حداقل یک عضو داشته باشد.

II. هر عضو (هر واحد) فقط به یک خوشه تعلق داشته باشد.

✓ به عبارت دیگر در روش افرازبندی، خوشه‌ها کاملاً از هم مستقل بوده و با یکدیگر همپوشانی ندارند.

انواع الگوریتم‌های مربوط به روش افرایندی

- (1) روش k-means
- (2) روش k-medoids
- (3) روش k-medians
- (4) روش Fuzzy C-means
- (5) روش K-SVD و ...

الگوریتم k-means

- الگوریتم k میانگین یکی از ساده‌ترین و محبوب‌ترین الگوریتم‌هایی است که در «داده کاوی» بخصوص در حوزه «مدل یادگیری بدون ناظر» کار می‌رود.
- معمولاً در حالت چند متغیره، باید از ویژگی‌های مختلف اشیا به منظور خوشه کردن آن‌ها استفاده کرد. به این ترتیب با داده‌های چند بعدی سروکار داریم که معمولاً به هر بعد از آن، ویژگی یا خصوصیت گفته می‌شود.
- با توجه به استفاده از **توابع فاصله** مختلف که در این جا مطرح می‌شود، ممکن است بعضی از ویژگی‌های اشیا **کمی** و بعضی دیگر **کیفی** باشند. به هر حال آنچه اهمیت دارد روشی برای اندازه گیری میزان شباهت یا عدم شباهت بین اشیاء است که باید در روش‌های خوشه بندی لحاظ شود.

- در این الگوریتم داده‌ها به تعداد خوشه‌هایی از پیش تعیین شده (k معلوم) تقسیم می‌شوند. به طوریکه برای هر خوشه، مرکز خوشه به صورت مقدار اولیه در نظر گرفته شده و سپس فاصله هر نقطه تا مرکز محاسبه می‌شود و مشاهداتی که به هر مرکز نزدیک‌تر باشند، در خوشه مربوطه قرار گرفته یا تخصیص داده می‌شوند.
- پس از تخصیص هر مشاهده به یکی از خوشه‌ها و تشکیل خوشه‌ها، برای هر خوشه یک نقطه جدید به عنوان مرکز خوشه محاسبه می‌شود.
- در این الگوریتم مرکز خوشه، میانگین مشاهدات موجود در هر خوشه است.
- این فرآیند تا جایی تکرار می‌شود که تغییر در میانگین خوشه‌ها منجر به تغییر در مشاهدات درون هر خوشه نشود.
- **توجه شود که اگر مراکز اولیه خوشه‌ها به درستی انتخاب نشود، خوشه‌های حاصل شده در پایان الگوریتم کیفیت مناسبی نخواهند داشت و این یکی از ضعف‌های این الگوریتم است.**

توابع فاصله:

- هدف خوشه بندی قرار گرفتن مشاهدات یا رکوردهای نزدیک به هم در یک خوشه باشد، به طوریکه بین مشاهدات درون خوشه کمترین فاصله و بین خوشه ها بیشترین فاصله باشد. خوشه بندی را مبتنی بر فاصله می گویند.

انواع توابع فاصله

1. فاصله اقلیدوسی:

- با استفاده از «فاصله اقلیدسی» Euclidean Distance کوتاهترین فاصله بین دو نقطه بر اساس رابطه فیثاغورث، محاسبه می‌شود. اگر x و y دو نقطه با p مؤلفه باشند، فاصله اقلیدسی بین این دو به صورت زیر قابل محاسبه است:

$$D_{euc} = \left(\sum_{i=1}^p (x_i - y_i)^2 \right)^{\frac{1}{2}}$$

✓ فاصله اقلیدسی نامنفی است.

✓ وجود رابطه انعکاسی برای فاصله اقلیدسی اثبات می‌شود.

✓ رابطه مثلثی برای فاصله اقلیدسی برقرار است.

2. فاصله منهتن:

- اگر به جای مربع فاصله بین مولفه‌ها، از قدر مطلق فاصله بین مولفه‌های نقاط استفاده شود، تابع فاصله را «منهتن» Manhattan می‌نامند. این نام به علت تقاطع منظم خیابان‌ها در محله منهتن نیویورک انتخاب شده است.
- اگر x و y نقطه با p مؤلفه (p بعدی) باشند، شیوه محاسبه فاصله منهتن به صورت زیر خواهد بود:

$$D_{man} = \sum_{i=1}^p |x_i - y_i|$$

✓ نکته: از آنجایی که میانه مولفه‌ای، نقطه‌ای است که فاصله منهتن را کمینه می‌کند، در بعضی از الگوریتم‌های خوشه بندی از میانه به عنوان مرکز هر خوشه استفاده می‌شود.

3. فاصله مینکوفسکی:

- **فاصله مینکوفسکی** حالت کلی تری از فاصله اقلیدسی و منهتن است. اگر A و B دو نقطه در فضای p بعدی باشند، فاصله مینکوفسکی برای آن‌ها به صورت زیر محاسبه می‌شود. پارامتر فاصله مینکوفسکی در اینجا d در نظر گرفته شده است، این فاصله به صورت زیر نشان داده می‌شود.

$$D_{mink}(A, B; d) = \left(\sum_{i=1}^p |x_i - y_i|^d \right)^{\frac{1}{d}}$$

مثال از الگوریتم k-means

- مشاهدات زیر مفروض هستند.
- $\{2, 4, 10, 12, 3, 20, 30, 11, 25\}$
مطلوبست الگوریتم k-means برای $k=2$ و مقادیر اولیه $m_{01} = 3$ و $m_{02} = 4$.

	خوشه اول	خوشه دوم
مرحله اول	$m_{01} = 3$ $\{2,3\}$	$m_{02} = 4$ $\{4,10,12,3,20,30,11,25\}$
مرحله دوم	$m_{11} = 2.5$ $\{2,3,4\}$	$m_{12} = 16$ $\{10,12,20,30,11,25\}$
مرحله سوم	$m_{21} = 3$ $\{2,3,4,10\}$	$m_{22} = 18$ $\{12,20,30,11,25\}$
مرحله چهارم	$m_{31} = 4.75$ $\{2,3,4,10,11,12\}$	$m_{32} = 19.6$ $\{20,30,25\}$
مرحله پنجم	$m_{41} = 7$ $\{2,3,4,10,11,12\}$	$m_{42} = 25$ $\{20,30,25\}$

نقاط ضعف الگوریتم k-means

1. این روش زمانی کاربرد دارد که بتوان مرکز خوشه را حدس زد. چرا که جواب نهایی به انتخاب مراکز اولیه خوشه وابسته است.
2. روال خاصی برای تعیین مقدار k و مراکز اولیه خوشه وجود ندارد.
3. این روش برای کشف خوشه‌های پیچیده مناسب نیست.
4. این روش نسبت به داده‌های پرت (داده‌های دور از مرکز) حساس است چرا که داده‌های پرت به راحتی می‌توانند مرکز خوشه را تغییر دهند و نتایج را نامناسب نشان دهند.
5. امکان تولید خوشه‌های خالی توسط الگوریتم وجود دارد.
6. تفسیر نتایج حاصل از خوشه‌بندی نیاز به تجربه و افراد خبره در این زمینه دارد.
7. انتخاب روش‌های صحیح برای تعیین فواصل رکوردها تا مراکز خوشه کار آسانی نیست.

راه حل‌های پیشنهادی برای رفع ضعف‌های الگوریتم k-means

- راه حل‌های رفع **مشکل انتخاب مراکز اولیه** در الگوریتم k-means:
 1. اجرای چندین باره این الگوریتم به ازای مقادیر اولیه متفاوت می‌باشد.
 2. راه حل دیگر انتخاب تعداد مراکز اولیه بیش از تعداد خوشه‌ها و سپس ادغام کردن آن‌ها باتوجه به مقدار k است.
 3. سومین راه حل انتخاب مقادیری است که بیشترین جداکنندگی را داشته باشند.

- برای رفع **ضعف الگوریتم k-means در تولید خوشه‌های خالی** راه حل‌های زیر می‌توانند کارآمد باشند:

1. مشاهده‌ای که دارای دورترین فاصله از مرکز خوشه است را به عنوان یک خوشه جدید انتخاب کنیم و خوشه‌های خالی را پر کنیم.

2. در صورتیکه مرکز خوشه خالی یکی از رکوردهای ما باشد، می‌توان آن را به عنوان داده پرت تشخیص داد و از میان مشاهدات خارج کرد.

3. اگر چندین خوشه خالی وجود داشته باشد، می‌توان با تکرار روش‌های فوق ضعف خوشه‌های خالی را برطرف کرد.

✓ اگر در میان داده‌ها نقاط پرت زیادی وجود داشته باشد در این صورت پیشنهاد می‌شود به جای استفاده از میانگین داده‌ها برای تعیین مرکز خوشه، از میانه خوشه‌ها استفاده شود که به این الگوریتم، الگوریتم k-median گفته می‌شود.

بهبود کیفی خوشه‌های تولیدشده:

1. شکستن خوشه‌ها به خوشه‌های کوچکتر:

در صورتی که رکوردهایی وجود داشته باشد که فواصل آن‌ها از تا مرکز از سایر مشاهدات داخل خوشه بیشتر باشد، توصیه به شکستن آن خوشه می‌شود.

2. ادغام خوشه‌های مشابه یا نزدیک:

خوشه‌هایی که مجموع فواصلشان تا مرکز آن‌ها از سایر خوشه‌ها کمتر باشد (برای مثال دو الی چند خوشه به هم بسیار نزدیک باشند)، توصیه به ادغام این خوشه‌ها می‌شود.

3. حذف خوشه‌های کوچک:

در صورتی که بعضی از نقاط پرت باشند، خوشه‌های کوچک ایجاد شده توسط این نقاط ممکن است نامناسب بوده و توصیه به حذف خوشه می‌گردد.

خوشه بندی سلسله مراتبی

تحلیل خوشه‌ای

خوشه بندی سلسله مراتبی

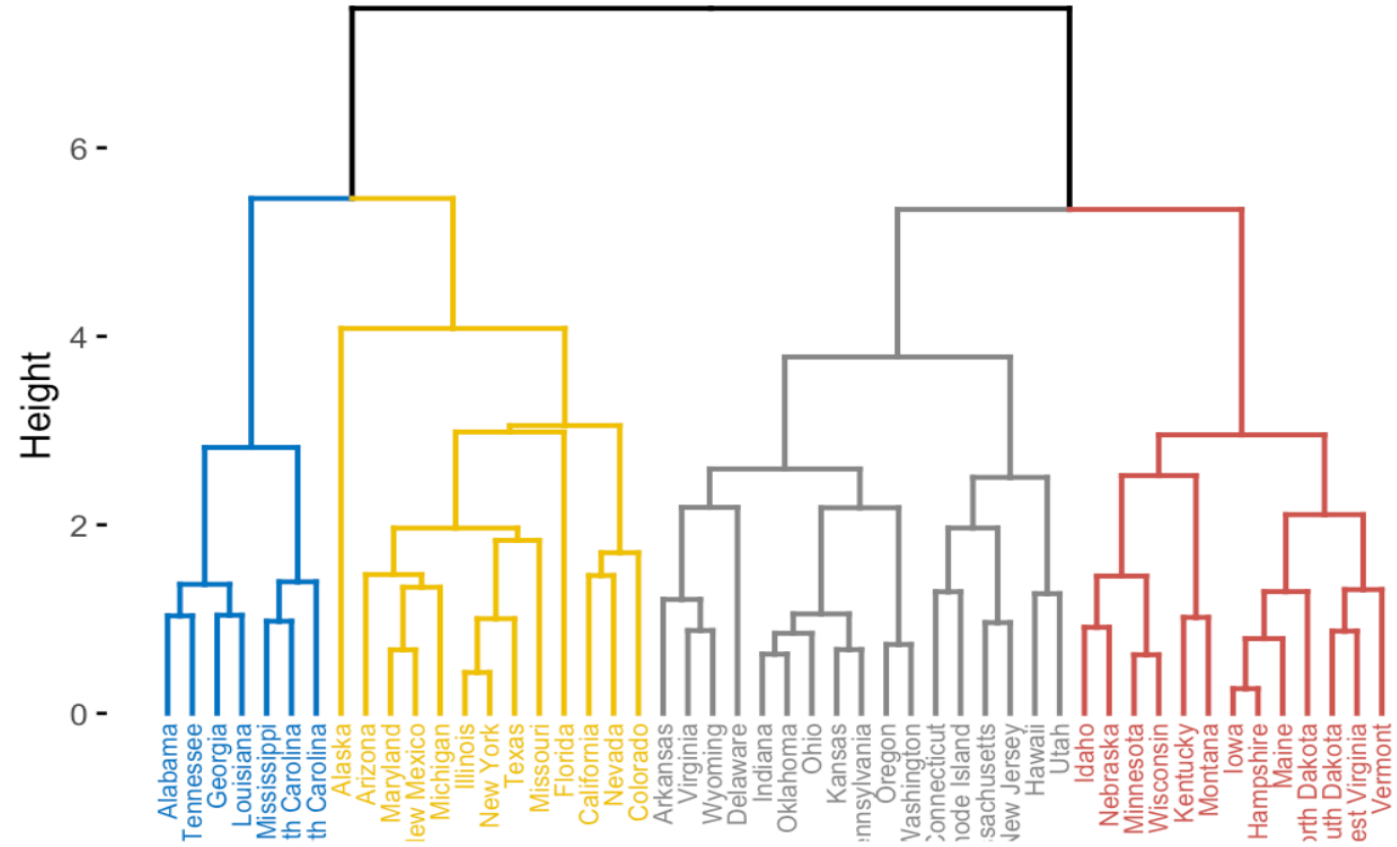
- با توجه به آن که روش‌های خوشه‌بندی مطرح شده تا حد زیادی پاسخ‌گوی نیاز ما هستند، اما در بعضی از موارد هدف تقسیم‌بندی داده‌ها در گروه‌هایی است که سطوح متفاوتی داشته و نسبت به یکدیگر در سلسله مراتب مختلفی قرار می‌گیرند.
- از این روش‌های خوشه‌بندی سلسله مراتبی از اهمیت بیشتری برخوردار هستند.
- مثال) مدیر منابع انسانی شرکت الکترونیکی، کارکنان شرکت خود را در گروه‌هایی مانند مدیران اجرایی، مدیران و کادر اداری سازماندهی می‌کند. کادر اداری را نیز به زیرگروه‌هایی مانند کارمندان ارشد، کارمندان و کارآموزان تقسیم می‌کند. با این تقسیم‌بندی کسب اطلاعات و خلاصه‌سازی و استخراج اطلاعات راحت‌تر می‌باشد.

- برعکس خوشه بندی تفکیکی که اشیاء را در گروه های مجزا تقسیم می کند، «خوشه بندی سلسله مراتبی» در هر سطح از فاصله، نتیجه خوشه بندی را نشان می دهد. این سطوح به صورت سلسله مراتبی هستند.

-
- برای نمایش نتایج خوشه بندی به صورت سلسله مراتبی از «**درختواره**» **Dendrogram** استفاده می شود. این شیوه، روشی موثر برای نمایش نتایج خوشه بندی سلسله مراتبی است.

- روش های خوشه بندی سلسله مراتبی خود نیز شامل چندین روش می باشند که از جمله آنها می توان به خوشه بندی سلسله مراتبی چندمرحله ای، خوشه بندی سلسله مراتبی با استفاده از مدل سازی پویا، خوشه بندی سلسله مراتبی احتمالی و ... اشاره کرد.

Cluster Dendrogram



بر خلاف تکنیک های خوشه بندی که بر پایه افراز که تعداد خوشه ها به عنوان پارامتر ورودی توسط کاربر مشخص می شد ، در تکنیک های خوشه بندی سلسله مراتبی مجموعه داده ها و معیاری جهت ارزیابی تشابه به عنوان ورودی معین می شوند .

بسته به اینکه تحلیل این ساختار سلسله مراتبی از پایین به بالا (Bottom-up) و یا برعکس از بالا به پایین (Top-down) انجام شود ، عملیات اصلی این الگوریتم ها را می توان در دو دسته ی ادغام و تقسیم قرار داد . فرآیند ادغام و یا تقسیم چند خوشه در این تکنیک ها نقش مهمی را ایفا می کند . چرا که در صورت اخذ تصمیمی ضعیف جهت ادغام و تقسیم خوشه ها تکنیک توانایی برگشت و اصلاح آن را ندارد و این عمل باعث کاهش کیفیت در خوشه های تولید شده نهایی خواهد داشت .

الگوریتم های خوشه بندی سلسله مراتبی که یک فرآیند پایین به بالا را طی می کنند ، کارشان با قرار دادن هر نمونه داده در یک خوشه ی مجزا شروع می شود با ادغام خوشه ها الگوریتم تا جایی پیش می رود که یا کلیه ی نمونه ها در یک خوشه قرار گیرند و یا شرط از پیش تعیین شده ای به عنوان پایان اجرا مشخص شده باشد . اغلب تکنیک های موجود سلسله مراتبی متعلق به این دسته هستند .

در طرف دیگر الگوریتم های سلسله مراتبی قرار دارند که یک استرانژی بالا به پایین را دنبال می کنند . در این روش ها بر عکس تکنیک های قبلی در ابتدا کلیه نمونه ها در یک خوشه قرار می گیرند . پس از آن با کمک از یک معیار تشابه در چند مرحله بصورت سلسله مراتبی این خوشه به خوشه های کوچکتر تقسیم می شود . در این الگوریتم ها نیز می توان کار را تا جایی ادامه داد تا هر نمونه در یک خوشه قرار بگیرد و یا اینکه شرطی را جهت پایان اجرای الگوریتم معین نمود . در هر دو دسته از الگوریتم های سلسله مراتبی کاربر می تواند تعداد خوشه های تولید شده نهایی را بعنوان یک شرط پایانی مشخص کند .

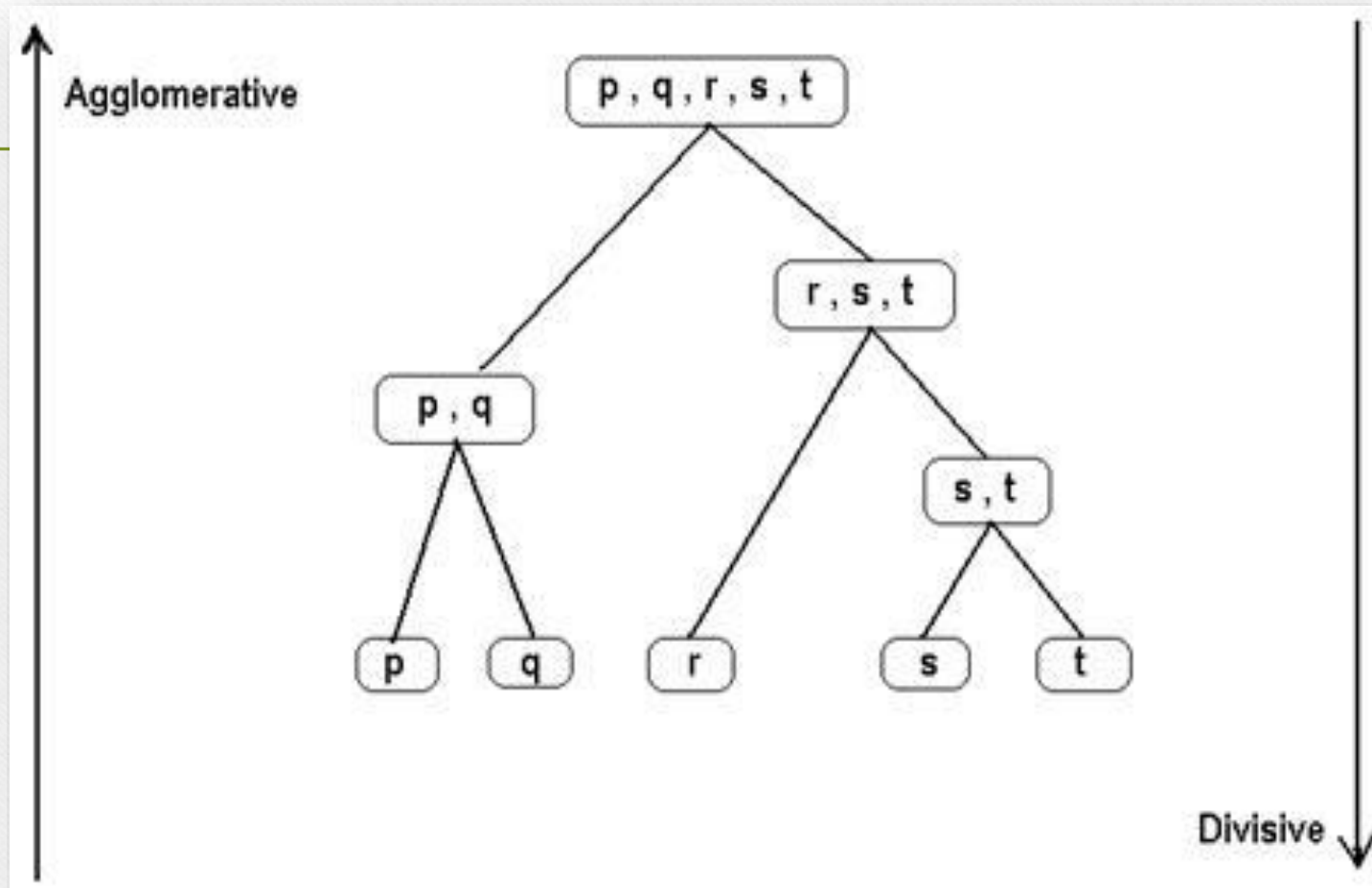
1. بالا به پایین (Top-Down) یا تقسیم کننده (Divisive)

در این روش ابتدا تمام داده‌ها به عنوان یک خوشه در نظر گرفته می‌شوند و سپس در طی یک فرایند تکراری در هر مرحله داده‌هایی شباهت کمتری به هم دارند به خوشه‌های مجزایی شکسته می‌شوند و این روال تا رسیدن به خوشه‌هایی که دارای یک عضو هستند ادامه پیدا می‌کند.

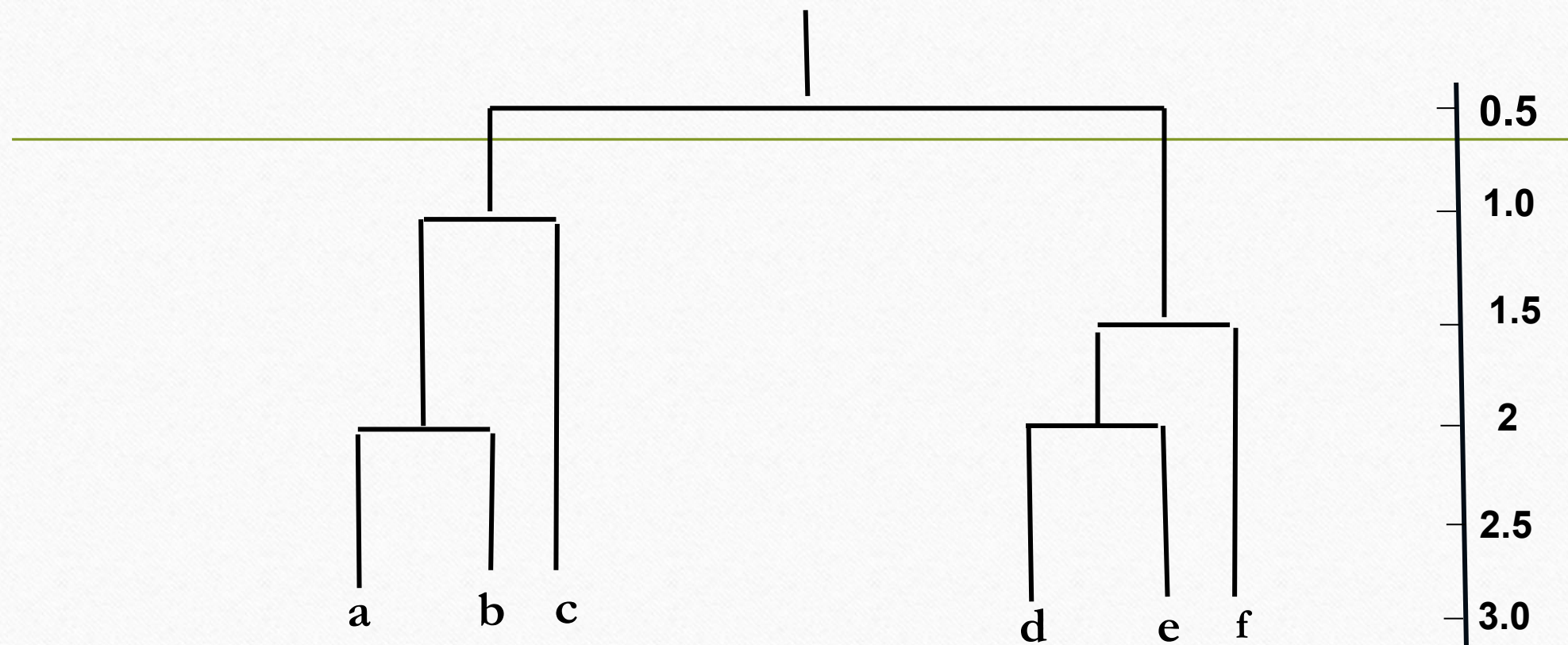
2. پایین به بالا (Bottom-Up) یا متراکم شونده (Agglomerative)

در این روش ابتدا هر داده‌ها به عنوان خوشه‌ای مجزا در نظر گرفته می‌شود و در طی فرایندی تکراری در هر مرحله خوشه‌هایی که شباهت بیشتری با یکدیگر با یکدیگر ترکیب می‌شوند تا در نهایت یک خوشه و یا تعداد مشخصی خوشه حاصل شود. از انواع الگوریتم‌های خوشه‌بندی سلسله مراتبی متراکم شونده رایج می‌توان از الگوریتم‌های Single-Link ، Average-Link و Complete-Link نام برد. تفاوت اصلی در بین تمام این روشها به نحوه محاسبه شباهت بین خوشه‌ها مربوط می‌شود.

تفاوت بین روشهای بالا به پایین با روشهای پایین به بالا



معمولا فرآیند خوشه ای سلسله مراتبی توسط یک نمودار با نام دندوگرام نمایش داده می شود که مثالی از آن در شکل زیر ملاحظه می شود :



این نمودار ادغام و تقسیم خوشه ها در هر مرحله نمایش می دهد . محور عمودی کنار نمودار، مقادیر مقیاس میان خوشه ها را نشان می دهد . برای مثال همانطور که مشاهده می شود هنگامی تشابه میان خوشه های $\{a, b, c\}$ و $\{d, e, f\}$ تقریباً برابر با 0.5 این دو خوشه با هم ادغام می شوند.

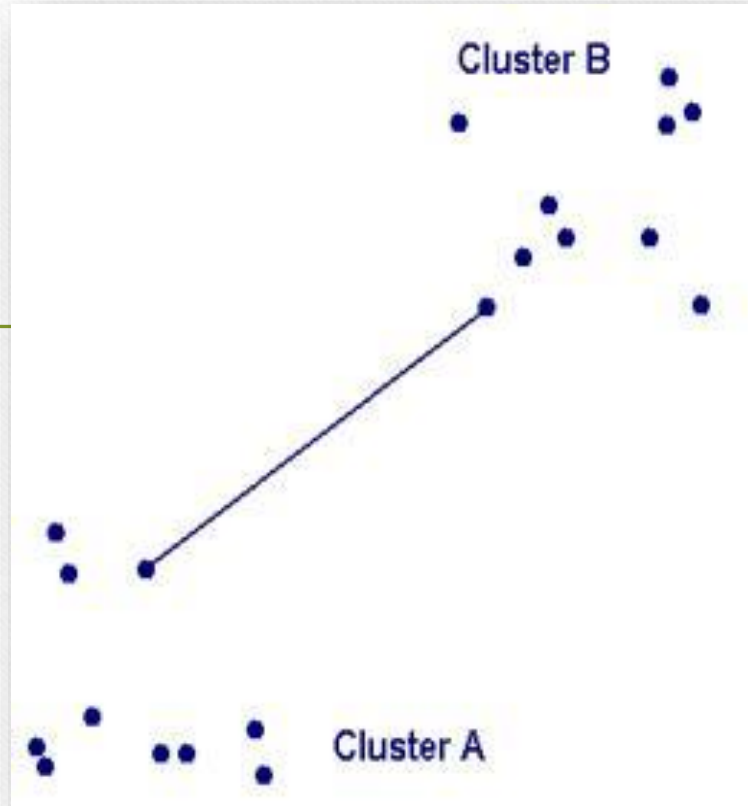
خوشه‌بندی با روش Single-Link

این روش یکی از قدیمی‌ترین و ساده‌ترین روشهای خوشه‌بندی است و جزء روشهای خوشه‌بندی سلسله مراتبی و انحصاری محسوب می‌شود. به این روش خوشه‌بندی، تکنیک نزدیکترین همسایه (Nearest Neighbour) نیز گفته می‌شود. در این روش برای محاسبه شباهت بین دو خوشه A و B از معیار زیر استفاده می‌شود:

$$d_{AB} = \min d_{ij} \quad i \in A, j \in B$$

که i یک نمونه داده متعلق به خوشه A و j یک نمونه داده متعلق به خوشه B می‌باشد.

در واقع در این روش شباهت بین دو خوشه، کمترین فاصله بین یک عضو از یکی با یک عضو از دیگری است. در شکل روبرو این مفهوم بهتر نشان داده شده است.



شباهت بین دو خوشه در روش Single-Link برابر است با کمترین فاصله بین داده‌های دو خوشه

مثال: فرض می شود که 5 نمونه داده موجود است که هر یک از آنها دارای دو

ویژگی X و Y می باشد

	O_1	O_2	O_3	O_4	O_5
X	1	5	6	1	4
Y	2	2	1	1	1

	O_1	O_2	O_3	O_4	O_5
O_1	0				
O_2	4	0			
O_3	5.1	1.4	0		
O_4	1	4.1	5.0	0	
O_5	3.1	1.4	2.0	3.0	0

معیار ارزیابی فاصله اقلیدسی است. در میان مقادیر کوچکترین عدد یک است که نشان دهنده ی تشابه بالای O_1 و O_4 می باشد بنابراین در این مرحله این دو نمونه می توانند در یک خوشه قرار گرفته و بعبارتی ادغام شوند .

$$\text{Distance}(\{O_1, O_4\}, O_2) = \min(d(O_1, O_2), d(O_4, O_2)) = \min(4.0, 4.1) = 4.0$$

$$\begin{aligned} \text{Distance}(\{O_1, O_4\}, O_3) &= \min(d(O_1, O_3), \\ &\quad d(O_4, O_3)) = \min(5.1, 5.0) = 5.0 \end{aligned}$$

$$\begin{aligned} \text{Distance}(\{O_1, O_4\}, O_5) &= \min(d(O_1, O_5), \\ &\quad d(O_4, O_5)) = \min(3.1, 3.0) = 3.0 \end{aligned}$$

	$\{O_1, O_4\}$	O_2	O_3	O_5
$\{O_1, O_4\}$	0			
O_2	4	0		
O_3	5.0	1.4	0	
O_5	3.0	1.4	2.0	0

در این مرحله دو مقدار 1.4 در ماتریس نشان دهنده ی نزدیکی نمونه های O_3 و O_5 به نمونه ی O_2 است می توان این سه نمونه را در یک خوشه ادغام کرد و کار را ادامه داد. اما اگر شرط الگوریتم این است که در هر مرحله فقط یکی از نمونه ها با نمونه یا خوشه ی دیگر ادغام شود ، مجبور خواهیم بود تا نمونه ی O_2 را با یکی نمونه های O_3 یا O_5 ادغام کنیم .

$$\text{Distance}(\{O_1, O_4\}, \{O_2, O_3\}) = \min(d(O_1, O_2), d(O_4, O_2), d(O_4, O_3), d(O_1, O_3)) = 4.0$$

$$\text{Distance}(\{O_1, O_4\}, O_5) = \min(d(O_1, O_5), d(O_4, O_5)) = 3.0$$

$$\text{Distance}(\{O_2, O_3\}, O_5) = \min(d(O_2, O_5), d(O_3, O_5)) = 1.4$$

نمونه ی O_2 را با O_3 ادغام و پس از آن ماتریس تشابه را بروز رسانی می کنیم

	$\{O_1, O_4\}$	$\{O_2, O_3\}$	O_5
$\{O_1, O_4\}$	0		
$\{O_2, O_3\}$	4.0	0	
O_5	3.0	1.4	0

ماتریس تشابه پس از اجرای مرحله دوم

با توجه به محتوای ماتریس تشابه در این مرحله نمونه ی O_5 با خوشه ی $\{O_2, O_3\}$ ادغام خواهد شد و ماتریس تشابه پس از این ادغام اصلاح می شود. اما فاصله همان 1.4 را نشان می دهد .

با توجه محتوای ماتریس تشابه در این مرحله نمونه ی O_5 با خوشه ی $\{O_2, O_3\}$ ادغام خواهد شد و ماتریس تشابه پس از این ادغام اصلاح می شود. اما فاصله همان 1.4 را نشان می دهد.

	$\{O_1, O_4\}$	$\{O_2, O_3, O_5\}$
$\{O_1, O_4\}$	0	
$\{O_2, O_3, O_5\}$	3.0	0

ماتریس تشابه پس از اجرای مرحله سوم

بالاخره در مرحله نهایی دو خوشه در یک خوشه ادغام می شوند. در واقع کلیه نمونه ها در یک نمونه قرار می گیرند.

مثال: در این مثال 6 نمونه داده و ماتریس فاصله بین آنها در جدول نشان داده شده است.

جدول 1: ماتریس فاصله بین 6 نمونه داده

6	5	4	3	2	1	
8	12	24	13	4	0	1
10	11	22	10	0		2
9	3	7	0			3
18	6	0				4
8.5	0					5
0						6

در ابتدا هر داده به عنوان یک خوشه در نظر گرفته می شود و یافتن نزدیکترین خوشه در واقع یافتن کمترین فاصله بین داده های بالا خواهد بود. با توجه به جدول 1 مشخص است که داده های 3 و 5 کمترین فاصله را دارا هستند. و در نتیجه آنها را با هم ترکیب کرده و خوشه جدیدی حاصل می شود که فاصله آن از سایر خوشه ها برابر است با کمترین فاصله بین 3 و 5 از سایر خوشه ها. نتیجه در جدول 2 نشان داده شده است.

جدول 2: ماتریس فاصله بین 5 خوشه حاصل از تکرار اول

6	4	(3 و 5)	2	1	
8	24	12	4	0	1
10	22	10	0		2
8.5	6	0			(3 و 5)
18	0				4
0					6

با توجه به جدول 2 مشخص است که داده‌های 1 و 2 کمترین فاصله را دارا هستند. و در نتیجه آنها را با هم ترکیب کرده و خوشه جدیدی حاصل می‌شود که فاصله آن از سایر خوشه‌ها برابر است با کمترین فاصله بین 1 و 2 از سایر خوشه‌ها.

جدول 3: ماتریس فاصله بین 4 خوشه حاصل از تکرار دوم

6	4	(3 و 5)	(1 و 2)	
8	22	10	0	(1 و 2)
8.5	6	0		(3 و 5)
18	0			4
0				6

با توجه به جدول 3 مشخص است که خوشه‌های (3 و 5) و 4 کمترین فاصله را دارا هستند. و در نتیجه آنها را با هم ترکیب کرده و خوشه جدیدی حاصل می‌شود که فاصله آن از سایر خوشه‌ها برابر است با کمترین فاصله بین (3 و 5) و 4 از سایر خوشه‌ها.

با توجه به جدول 4 مشخص است که خوشه‌های (1 و 2) و 6 کمترین فاصله را دارا هستند. و در نتیجه آنها را با هم ترکیب کرده و خوشه جدیدی حاصل می‌شود که فاصله آن از سایر خوشه‌ها برابر است با کمترین فاصله بین (1 و 2) و یا 6 از سایر خوشه‌ها.

جدول 4: ماتریس فاصله بین 3 خوشه حاصل از تکرار سوم

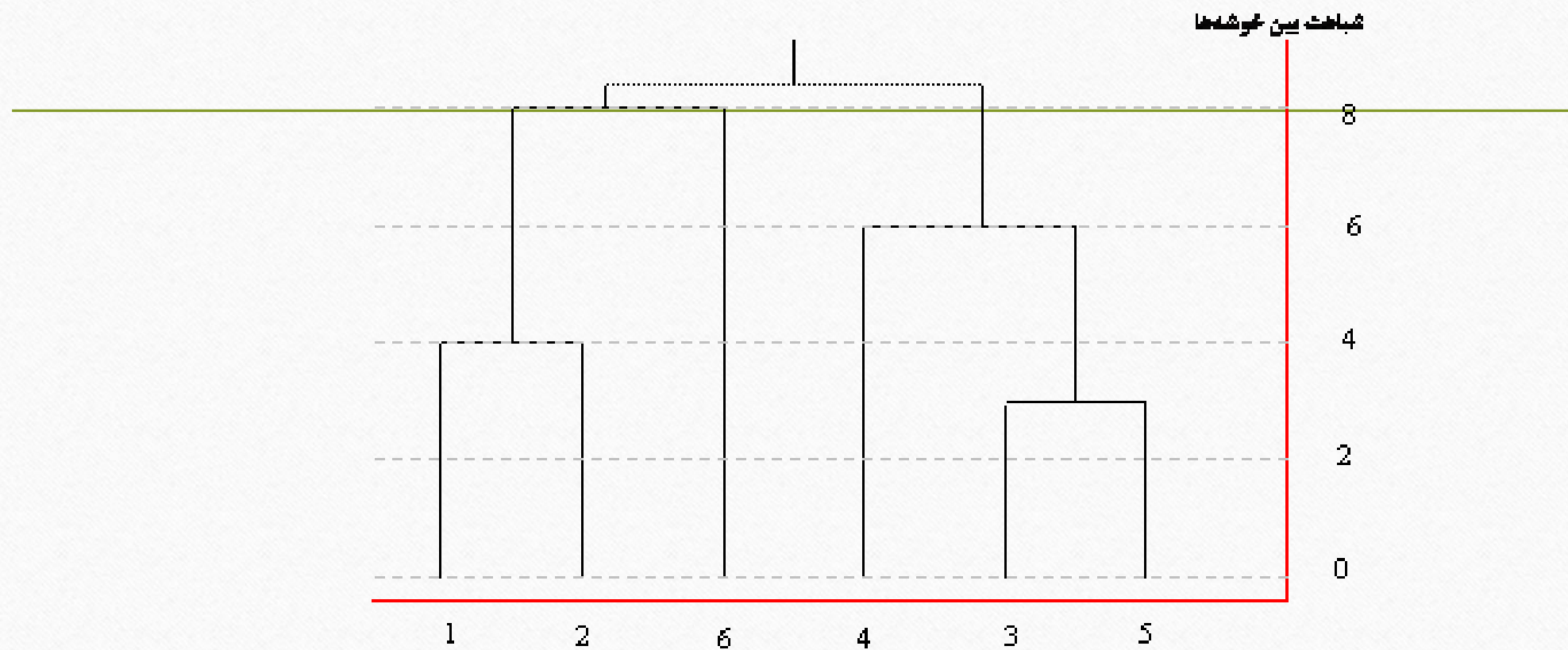
6	(3, 4 و 5)	(1 و 2)	
8	10	0	(1 و 2)
8.5	0		(3, 4 و 5)
0			6

جدول 5: ماتریس فاصله بین 2 خوشه حاصل از تکرار چهارم

(3, 4 و 5)	(1, 2 و 6)	
8.5	0	(1, 2 و 6)
0		(3, 4 و 5)

در نهایت این دو خوشه حاصل با هم ترکیب می‌شوند.

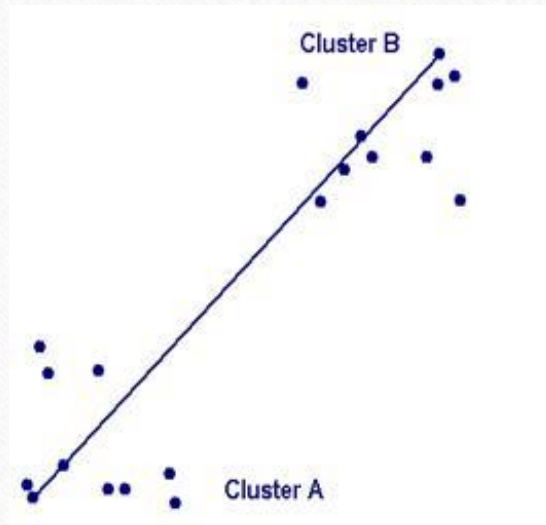
نتیجه ترکیب نهایی در دندوگرام زیر نشان داده شده است .



خوشه‌بندی با روش Complete-Link

این روش همانند Single-Link جزء روش‌های خوشه‌بندی سلسله مراتبی و انحصاری محسوب می‌شود. به این روش خوشه‌بندی، تکنیک دورترین همسایه (furthest Neighbour) نیز گفته می‌شود. در این روش برای محاسبه شباهت بین دو خوشه A و B از معیار زیر استفاده می‌شود:

$$d_{AB} = \max_{i \in A, j \in B} d_{ij}$$



که i یک نمونه داده متعلق به خوشه A و j یک نمونه داده متعلق به خوشه B می‌باشد. در واقع در این روش شباهت بین دو خوشه بیشترین فاصله بین یک عضو از یکی با یک عضو از دیگری است.

شباهت بین دو خوشه در روش Complete-Link برابر است با بیشترین فاصله بین داده‌های دو خوشه.

مثال: در این قسمت سعی شده است تا در مثالی با فرض داشتن 5 نمونه داده که هر یک از آنها

دارای دو ویژگی X و Y می باشد

	O_1	O_2	O_3	O_4	O_5
X	1	5	6	1	4
Y	2	2	1	1	1

	O_1	O_2	O_3	O_4	O_5
O_1	0				
O_2	4	0			
O_3	5.1	1.4	0		
O_4	1	4.1	5.0	0	
O_5	3.1	1.4	2.0	3.0	0

معیار ارزیابی فاصله اقلیدسی است. در میان مقادیر کوچکترین عدد یک است که نشان دهنده ی تشابه بالای O_1 و O_4 می باشد بنابراین در این مرحله این دو نمونه می توانند در یک خوشه قرار گرفته و بعبارتی ادغام شوند.

$$\text{Distance}(\{O_1, O_4\}, O_2) = \max(d(O_1, O_2), d(O_4, O_2)) = \max(4.0, 4.1) = 4.1$$

$$\text{Distance}(\{O_1, O_4\}, O_3) = \max(d(O_1, O_3), d(O_4, O_3)) = \max(5.1, 5.0) = 5.1$$

$$\text{Distance}(\{O_1, O_4\}, O_5) = \max(d(O_1, O_5), d(O_4, O_5)) = \max(3.1, 3.0) = 3.1$$

	$\{O_1, O_4\}$	O_2	O_3	O_5
$\{O_1, O_4\}$	0			
O_2	4.1	0		
O_3	5.1	1.4	0	
O_5	3.1	1.4	2.0	0

	$\{O_1, O_4\}$	$\{O_2, O_3, O_5\}$
$\{O_1, O_4\}$	0	
$\{O_2, O_3, O_5\}$	5.1	0

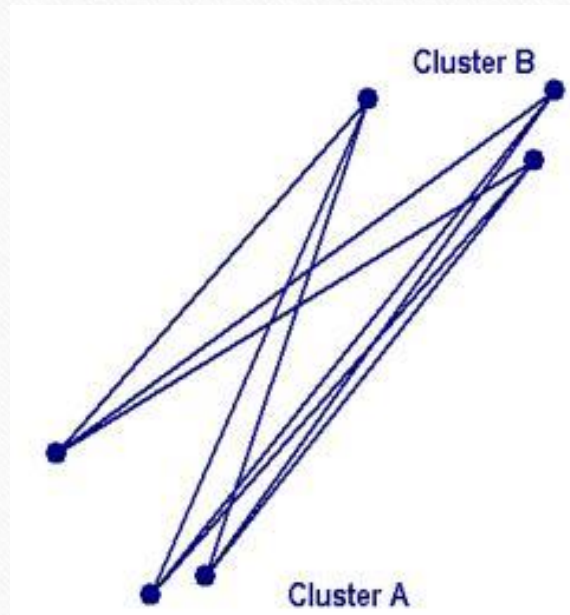
ماتریس تشابه نهایی

در ماتریس فوق کوچکترین عدد 1.4 است که مربوط به تشابه میان نمونه O_2 با نمونه های O_3 و O_5 است. همانند مثال قبل در این لحظه بدلیل تشابه یکسان با دو انتخاب روبرو هستیم بدلیل اینکه نشان دهیم الگوریتم می تواند در هر مرحله چندین نمونه را در یک خوشه واحد قرار دهد، سه نمونه را در این مرحله در یک خوشه ادغام می کنیم. در مرحله پایانی همانند الگوریتم قبل اعضا دو خوشه ی باقیمانده در یک خوشه ادغام می شوند.

خوشه‌بندی با روش Average-Link

$$d_{AB} = \frac{\sum_{i \in A, j \in B} d_{ij}}{N_A N_B}$$

این روش همانند Single-Link جزء روشهای خوشه‌بندی سلسله مراتبی و انحصاری محسوب می‌شود. از آنجا که هر دو روش خوشه‌بندی Single-link و Complete-link بشدت به نویز حساس می‌باشد، این روش که محاسبات بیشتری دارد، پیشنهاد شد. در این روش برای محاسبه شباهت بین دو خوشه A و B از معیار زیر استفاده می‌شود:



که i یک نمونه داده متعلق به خوشه A و j یک نمونه داده متعلق به خوشه B می‌باشد. و N_A تعداد اعضاء خوشه A و N_B تعداد اعضاء خوشه B است. در واقع در این روش، شباهت بین دو خوشه میانگین فاصله بین تمام اعضاء یکی با تمام اعضاء دیگری است.

در روش میانگین، باید میانگین کلیه فواصل زوج نمونه ها را طوری که هر نمونه متعلق به یک خوشه است را محاسبه کرد.

$$\text{Distance}(C_i, C_j) = \frac{1}{|C_i| \times |C_j|} \sum_{o_i \in C_i, o_j \in C_j} \text{Distance}(o_i, o_j)$$

منظور از $|C_i|$ تعداد نمونه هایی است که در خوشه C_i قرار دارند .

	O_1	O_2	O_3	O_5
O_1	0			
O_2	1.0	0		
O_3	11.0	2.0	0	
O_5	5.0	3.0	4.0	0

با توجه به محتوای ماتریس، نمونه های O_1 و O_2 در مرحله ی اول ادغام می شوند . پس از آن چنانچه ماتریس تشابه تشکیل شود مشاهده می شود که نمونه های O_3 و O_4 می توانند در مرحله ی بعدی ادغام شوند . در نهایت دو خوشه ی $\{O_1 \text{ و } O_2\}$ و $\{O_3 \text{ و } O_4\}$ در یک خوشه قرار می گیرند .

نحوه ی محاسبه ی فاصله میان این دو خوشه پایانی به صورت زیر است:

$$\text{Distance}(\{O_1 \text{ و } O_2\}, \{O_3 \text{ و } O_4\}) = \frac{1}{4} \times (d(O_1 \text{ و } O_3) + d(O_1 \text{ و } O_4) + d(O_2 \text{ و } O_3) + d(O_2 \text{ و } O_4)) = \frac{1}{4} (11+5+2+3) = 5.25$$

روش های پیشرفته خوشه بندی

تحلیل خوشه‌ای

- در روش‌های خوشه‌بندی که تاکنون مطرح شده‌اند، هر داده می‌توانست تنها به یک خوشه اختصاص یابد. چنین روشی برای تخصیص داده‌ها به خوشه‌ها برای برخی از کاربردهای خوشه‌بندی مانند بازاریابی که هدف تخصیص هر مشتری به یک سیستم بازاریابی می‌باشد، مفید است.

✓ اما در برخی از کاربردهای خوشه‌بندی این امر چندان خوشایند نیست. لذا از روش‌های منطبق با ایده **فازی** برای خوشه‌بندی و تخصیص داده‌ها استفاده می‌شود.

✓ از جمله این روش‌ها می‌توان خوشه‌بندی مبتنی بر مدل احتمالی و خوشه‌بندی فازی و ... را نام برد.

- **مثال** فرض کنید شرکت الکترونیکی فروشگاه آنلاینی دارد که علاوه بر خرید و فروش آنلاین، مشتریان می‌توانند نظرات خود را درباره محصولات در آنجا قرار دهند. این نظرات می‌توانند شامل چند محصول باشند و ممکن است محصولش حتی یک نظر هم نداشته باشد.

✓ هدف اصلی بررسی‌کننده نظرها، خوشه‌بندی نظرهای مختلف است. در این صورت تخصیص هر نظر به یکی از خوشه‌ها هدف این شرکت نمی‌باشد.

✓ در این صورت نیاز به روش خوشه‌بندی احساس می‌شود که نظرها شامل بیش از یک موضوع باشند. لذا آن‌ها را به بیش از یک خوشه اختصاص می‌دهد.