



داده کاوی
تحلیل خوشه ای

سمیه علیزاده

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مطالب مورد بحث در کلاس

- تعاریف و مفاهیم
- انباره داده ها
- آماده سازی داده ها (پیش پردازش داده ها)

مطالب مورد بحث در کلاس

- خوشه بندی
- دسته بندی
- قوانین انجمنی
- سریهای زمانی
- وب کاوی
- متن کاوی
- پیوندکاوی و تحلیل شبکه های اجتماعی

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مطالب مورد بحث در کلاس

- متدولوژی اجرای پروژه های داده کاوی
- کاربردهای داده کاوی در بازاریابی
- کاربردهای داده کاوی در مدیریت ارتباط با مشتری
- امنیت در داده کاوی

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خوشه بندی

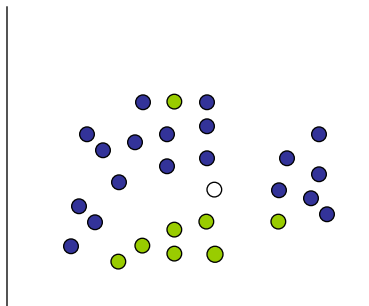


- تجزیه و تحلیل خوشه‌ای روشی برای گروه بندی داده‌ها یا مشاهدات با توجه به **شباهت** یا **درجه نزدیکی** آنها است.
- از طریق تجزیه و تحلیل خوشه‌ای، داده‌ها یا مشاهدات به دسته‌های **همگن** و **متمايز** از هم تقسیم می‌شوند.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خوشه بندی

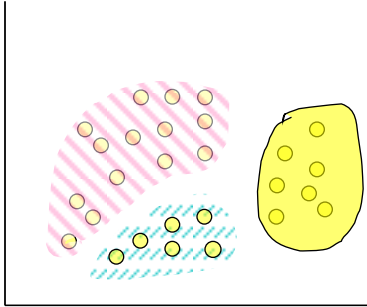
- در روش **خوشه بندی** هیچ دسته‌ای از قبل وجود ندارد و در واقع متغیرها بصورت مستقل و وابسته تقسیم نمی‌شوند. بلکه ما در اینجا بدنبال گروه‌هایی از داده‌ها هستیم که به هم **شباهت** دارند و با کشف این شباهت‌ها می‌توان رفتارها را بهتر شناسایی کرد و بر مبنای آنها طوری عمل کرد که نتیجه بهتری حاصل شود.



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خوشه بندی

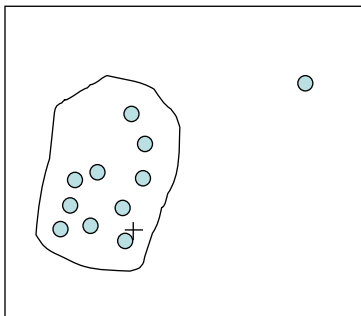
- در روش خوشه بندی هیچ دسته ای از قبل وجود ندارد و در واقع متغیرها بصورت مستقل و وابسته تقسیم نمی شوند. بلکه ما در اینجا بدنبال گروه هایی از داده ها هستیم که به هم شباهت دارند و با کشف این شباهت ها می توان رفتارها را بهتر شناسایی کرد و بر مبنای آنها طوری عمل کرد که نتیجه بهتری حاصل شود.



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

کاربردهای خوشه بندی در آماده سازی داده ها

- در بعضی موارد از خوشه بندی برای داده هایی که با سایر داده ها تفاوت چشمگیر دارد استفاده می نمایند.
- مثلاً یکسری از مشتریان همگی خریدی بالای 100 دلار در ماه دارند به غیر از یکی



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

نقاط قوت روش خوشه بندی

- روش خوشه بندی يك روش غير مستقیم است.
 - بدین معنی که این روش را می توان حتی هنگامی که هیچ نوع اطلاعات قبلی از ساختار داخلی پایگاه داده ها نداریم استفاده نمود. از این روش می توان برای کشف الگوهای پنهان و بهبود عملکرد روشهای مستقیم نیز استفاده نمود.
- خوشه بندی را می توان برای داده های گوناگون استفاده نمود.
 - با انتخاب درست اندازه فاصله های گوناگون خوشه بندی را می توان برای بیشتر انواع داده ها استفاده نمود.
- استفاده از این روش آسان است.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

نقاط ضعف روش خوشه بندی

- انتخاب اندازه های دقیق فواصل و وزنهای کار آسانی نمی باشد.
- این روش به پارامترهای اولیه نظیر تعداد خوشه ها ، حداقل نزدیکی و خوشه های اولیه حساس است.
- تفسیر نتایج این روش می تواند مشکل باشد و معمولاً نیاز به تحلیل افراد خبره در زمینه تجارت مورد نظر دارد.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

یک خوشه بندی خوب چیست؟

- یک روش خوشه بندی خوب خوشه‌های با کیفیت بالا براساس دو معیار زیر را تولید می‌کند:

شباهت بالایی نقاط داخلی هر کلاس و شباهت کم بین نقاط کلاسهای مختلف .

- کیفیت نتایج خوشه بندی بستگی به روش اندازه‌گیر شباهت به کار رفته و همچنین پیاده سازی آن روش دارد

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مراحل تجزیه و تحلیل خوشه‌ای

- انتخاب معیار شباهت یا نزدیکی مشاهدات
- انتخاب روش تجزیه و تحلیل خوشه‌ای
- تصمیم‌گیری در مورد تعداد خوشه‌ها
- تفسیر دسته‌ها یا گروه‌های تشکیل شده



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

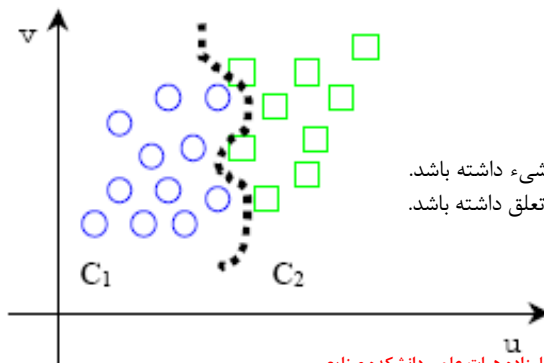
رویکردهای اصلی خوشه بندی

- روش افراز بندی
- روش سلسله مراتبی
- روش مبتنی بر چگالی
- روش Grid-based
- روش مبتنی بر مدل

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

روش افراز بندی

فرض کنید یک پایگاه داده از n شیء داشته باشیم. یک روش افراز بندی، K افراز از این داده‌ها درست می‌کند بطوریکه هر افراز یک خوشه را نشان می‌دهد. پس داده‌ها در k گروه گروه بندی می‌شوند که باید دارای دو شرط زیر باشند:



الف) هر گروه بایستی حداقل یک شیء داشته باشد.
ب) هر شیء باید تنها به یک گروه تعلق داشته باشد.

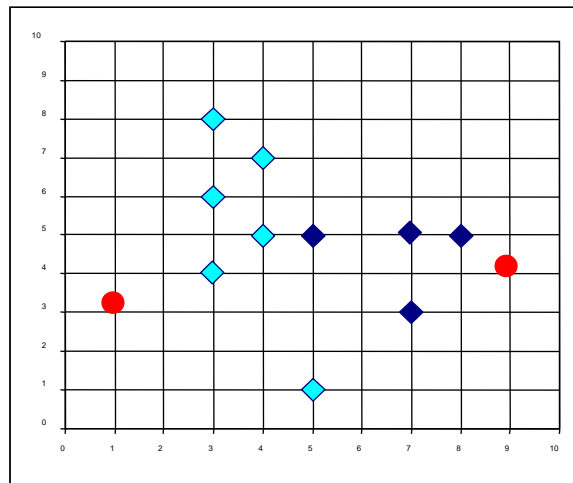
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

روش افزار بندی

- الگوریتم *K-means* که هر خوشه با میانگین اشیاء آن خوشه ، نمایش داده می شود. (با مرکز خوشه)
- الگوریتم *K-medoids* که هر خوشه با یکی از اشیاء که در نزدیکی مرکز خوشه جای گرفته است ، نمایش داده می شود.

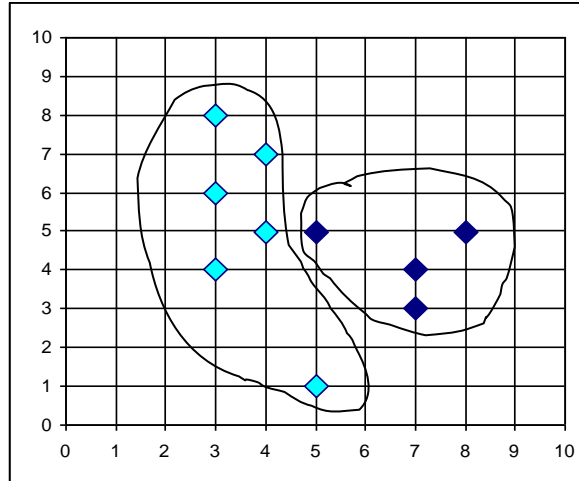
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

انتخاب K عضو دلخواه اولیه



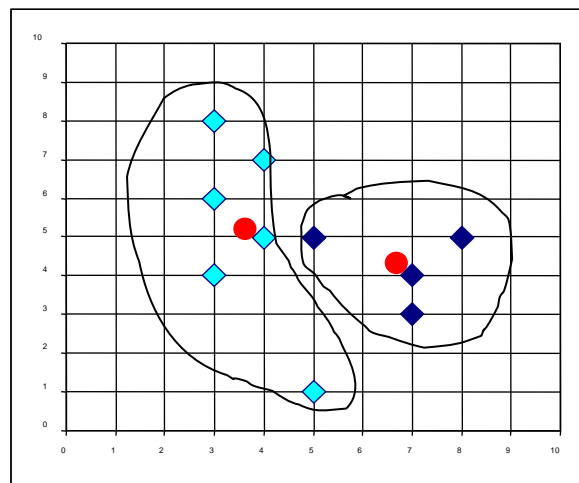
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

هر عضو به شبیه ترین می پیوندد



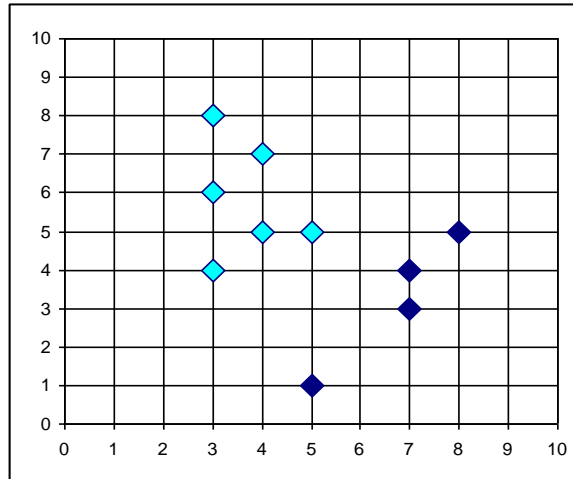
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

به روز کردن میانگین داده ها



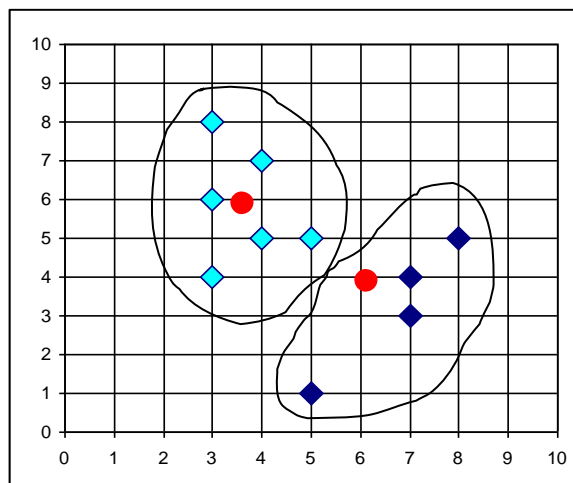
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

بازنگری مجدد



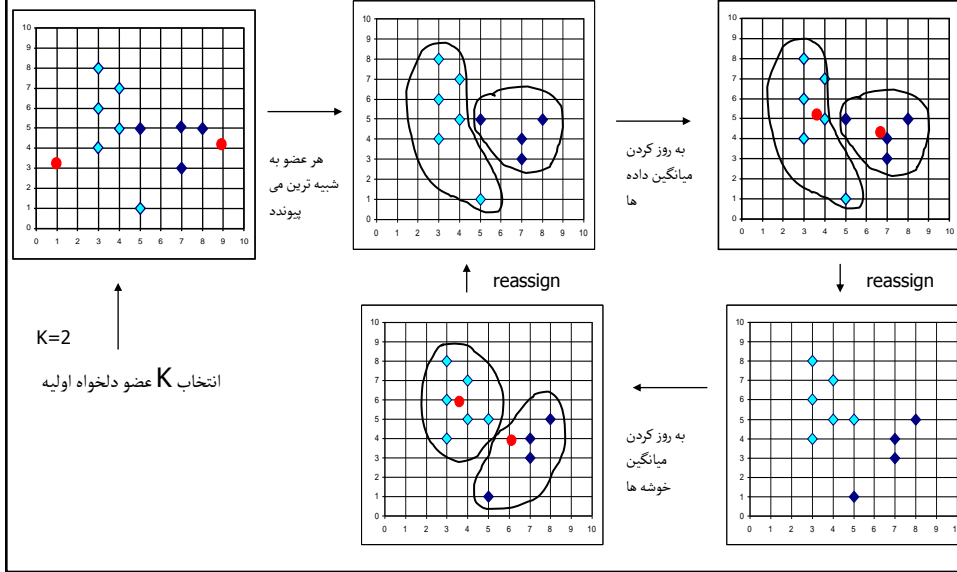
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

به روز کردن میانگین خوشه ها



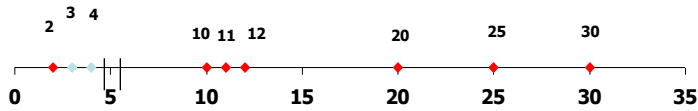
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

الگوریتم K-means

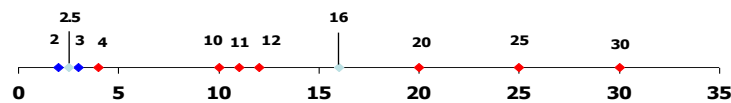


مثال K-Means

- $\{2, 4, 10, 12, 3, 20, 30, 11, 25\}$, $k=2$
- $m_1=3, m_2=4$

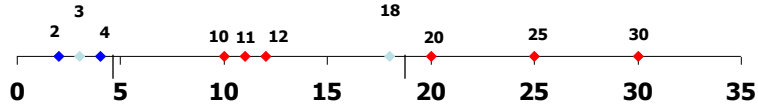


- $K_1=\{2, 3\}$, $K_2=\{4, 10, 12, 20, 30, 11, 25\}$, $m_1=2.5, m_2=16$

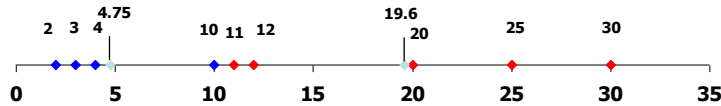


سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

- $K_1=\{2,3,4\}, K_2=\{10,12,20,30,11,25\}, m_1=3, m_2=18$

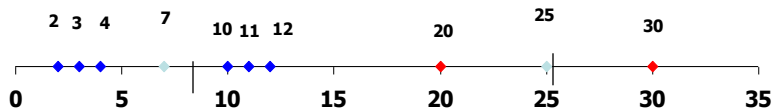


- $K_1=\{2,3,4,10\}, K_2=\{12,20,30,11,25\}, m_1=4.75, m_2=19.6$



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

- $K_1=\{2,3,4,10,11,12\}, K_2=\{20,30,25\}, m_1=7, m_2=25$



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

K-Means

نقاط قوت

- این روش نسبتاً برای پایگاه‌های داده بزرگ کارا و ارتقا پذیر می باشد زیرا پیچیدگی محاسباتی اش عبارتست از $O(tkn)$ که: n تعداد کل اشیاء، K تعداد خوشه‌ها و t تعداد تکرارهای الگوریتم است. این روش اغلب به یک بهینه محلی ختم می شود نه یک بهینه سراسری. بهینه سراسری از طریق روشهایی مانند ژنتیک بدست می آید.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

K-Means

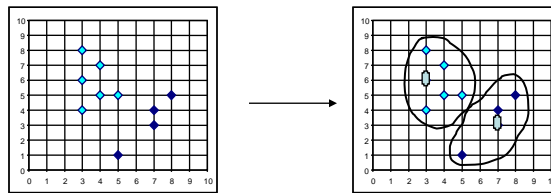
نقاط ضعف

- روش K-means تنها هنگامی کاربرد دارد که بتوان مراکز خوشه‌ها را تعریف نمود. مثلاً برای داده‌هایی با صفات رده‌ای این روش کارا نیست.
- از معایب این روش تعیین K می باشد که می‌بایست کاربر ابتدا آنرا معین کند و راه خاصی برای تعیین آن مشخص نشده است.
- همچنین این روش برای کشف خوشه‌هایی با شکل‌های پیچیده مناسب نیست.
- یکی از مهمترین نقاط ضعف این روش این است که در برابر نویزها و داده‌های دوراز مرکز حساس است زیرا این داده‌ها به راحتی مراکز را تغییر می دهند و ممکن است نتایج مطلوبی حاصل نشود.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

روش K -Medoids

- K -Medoids: در این روش بجای استفاده از مرکز یک خوشه به عنوان مرجع، می‌توان از **medoid** ها استفاده کرد. یعنی شی که در مرکزترین محل یک خوشه می‌باشد. بنابراین روش افراز هنوز می‌تواند مبتنی بر اصل مینیمم سازی مجموع عدم شباهتها میان هر شی و شیء مرجع متناظرش شکل بگیرد.



PAM (Partitioning Around Medoids, 1987)

استراتژی اساسی الگوریتم خوشه بندی K -medoids پیدا کردن K شیء نماینده آغازین (**medoid**) بطور دلخواه از n شیء پایگاه داده‌ها می‌باشد.

— هر شیء باقیمانده با **medoid** ای هم خوشه می‌شود که بیشترین شباهت را به آن داشته باشد. سپس این استراتژی مکرراً یکی از اشیاء **medoid** را با یکی از اشیاء غیر **medoid** جایگزین می‌کند به طوری که کیفیت نتیجه خوشه‌بندی بهبود بخشیده شود.

— این کیفیت با بکارگیری تابع هزینه تخمین زده می‌شود که میانگین عدم تشابه بین یک شیء و **medoid** آن خوشه را اندازه‌گیری می‌کند.

الگوریتم k-medoids

- شیء به عنوان medoid های اولیه به صورت دلخواه اختیار کن.
- تکرار کن تا اینکه هیچ تغییری رخ ندهد.
- هر کدام از اشیاء باقیمانده را به خوشه‌ای با نزدیکترین medoid تخصیص بده
- بطور تصادفی یک شیء غیر medoid را انتخاب کن ، .
- هزینه نهایی S را از عوض کردن (medoid آن خوشه) و محاسبه کن
- اگر $s < 0$ آنگاه جای عناصر را عوض کن تا مجموعه K تا medoid جدید شکل بگیرد.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

هزینه تغییرات

- برای اندازه‌گیری اینکه شیء o بهتر از o' به عنوان یک medoid هست یا خیر ، کفایت حاصل معادله زیر را بدست آوریم و اگر جابه جایی $E(o) < E(o')$ باشد :

$$E = \sum_{i=1}^n \sum_{p \in C_i} d(p, o_i)^2$$

- در اینجا E در اصل میزان کل فاصله ها از هر نقطه را نشان می دهد و \bar{C} میزان هزینه تعویض می کرد که منفی بودن آن برابر سود است.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

K-Medoids مثال

- 1, 2, 6, 7, 8, 10, 15, 17, 20 – break into 3 clusters
 - Cluster = 6 – 1, 2
 - Cluster = 7
 - Cluster = 8 – 10, 15, 17, 20
- Random non-medoid – 15 replace 7 (total cost=-13)
 - Cluster = 6 – 1 (cost 0), 2 (cost 0), 7(1-0=1)
 - Cluster = 8 – 10 (cost 0)
 - New Cluster = 15 – 17 (cost 2-9=-7), 20 (cost 5-12=-7)
- Replace medoid 7 with new medoid (15) and reassign
 - Cluster = 6 – 1, 2, 7
 - Cluster = 8 – 10
 - Cluster = 15 – 17, 20

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

K-Medoids مثال

- Random non-medoid – 1 replaces 6 (total cost= -1)
 - Cluster = 8 – 6 (cost 2-0=2), 7 (cost 1-1=0), 10 (cost 0)
 - Cluster = 15 – 17(cost 0), 20(cost 0)
 - New Cluster = 1 – 2 (cost 1-4= -3)
- Replace medoid 6 with new medoid (1) and reassign
 - Cluster = 1 – 2
 - Cluster = 8 – 6, 7, 10
 - Cluster = 15 – 17, 20
- Random non-medoid – 10 replaces 8 (total cost=2)
don't replace
 - Cluster = 1– 2(cost 0)
 - Cluster = 15 – 17 (cost 0), 20(cost 0)
 - New Cluster = 10 – 6 (cost 0), 7 (cost 0), 8 (cost 2-0=2)

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

K-Medoids مثال

- Random non-medoid – 17 replaces 15 (total cost=0) don't replace
 - Cluster = 1 – 2(cost 0)
 - Cluster = 8 – 6 (cost 0), 7 (cost 0), 10 (cost 0)
 - New Cluster = 17 – 15 (cost 2-0=2), 20(cost 3-5=-2)
- Random non-medoid – 20 replaces 15 (total cost=6) don't replace
 - Cluster = 1 – 2(cost 0)
 - Cluster = 8 – 6 (cost 0), 7 (cost 0), 10 (cost 0)
 - New Cluster = 20 – 15 (cost 5-0=5), 17(cost 3-2=1)
- Other possible changes all have high costs
 - 1 replaces 15, 2 replaces 15, 1 replaces 8, ...
- No changes, final clusters
 - Cluster = 1 – 2
 - Cluster = 8 – 6, 7, 10
 - Cluster = 15 – 17, 20

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

PAM (Partitioning Around Medoids)

- روش ذکر شده (partitioning around medoids) PAM نام دارد که یکی از اولین الگوریتمهای K-medoids است و تلاش می‌کند k افراز برای n شیء معین کند.
- در این روش همه زوجهای ممکن از اشیاء آنالیز می‌شوند که یکی medoid و دیگری غیر medoid است.
- یک شیء با شیئیء جابه جا می‌شود که بیشترین کاهش را در خطای مربع داشته باشد.
- لذا این روش برای پایگاه داده‌های بزرگ مشکل است. برای رفع این اشکال از الگوریتم‌های CLARA , CLARANS استفاده می‌شود.

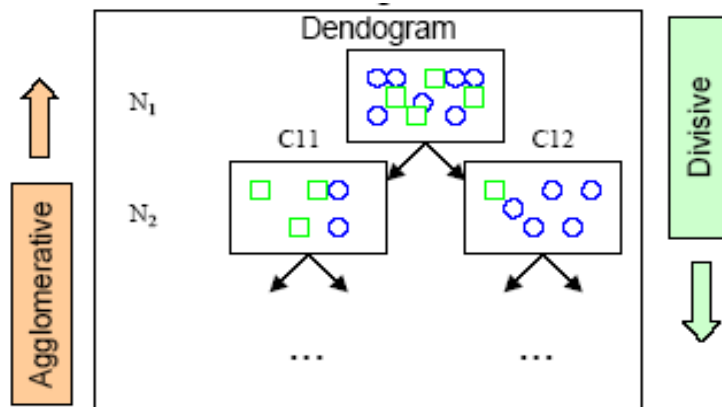
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

CLARA (Clustering Large Applications) (1990)

- توسط Kaufmann , Rousseeuw در 1990 ارائه شد. این الگوریتم برای پایگاه داده‌های بزرگ بکار می‌رود. به این صورت که چندین نمونه تصادفی از این پایگاه داده بر می‌دارد و الگوریتم PAM را روی هر نمونه اجرا می‌کند و آن نمونه را خوشه بندی می‌کند. سپس عناصر باقیمانده پایگاه داده را به نزدیکترین خوشه تخصیص می‌دهد.

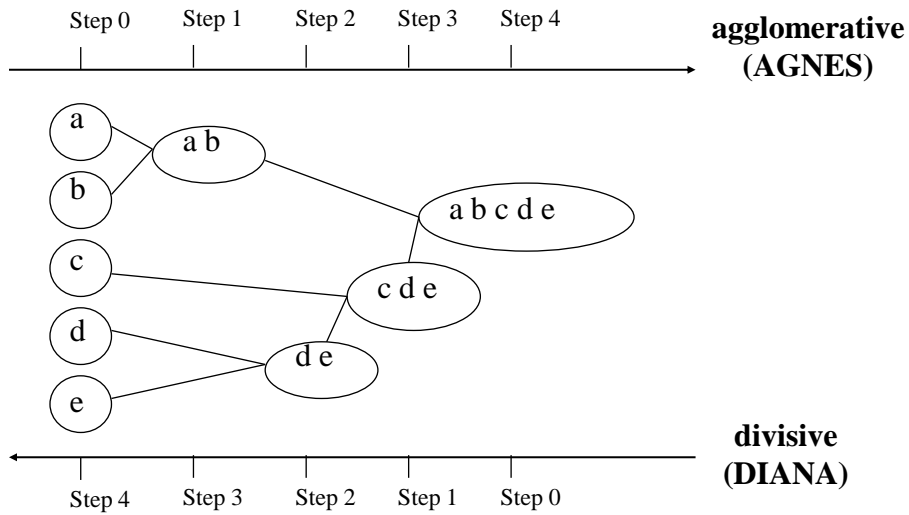
سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خوشه بندی سلسله مراتبی



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

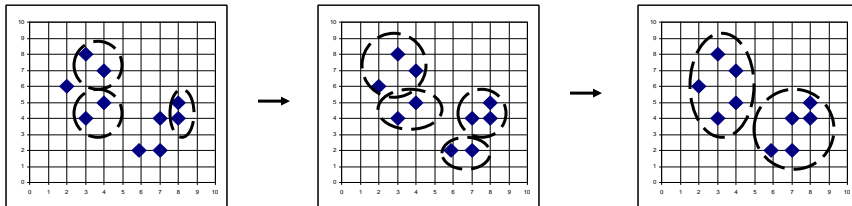
خوشه بندی سلسله مراتبی



سمیه علیزاده هیات علمی دانشکده صنایع

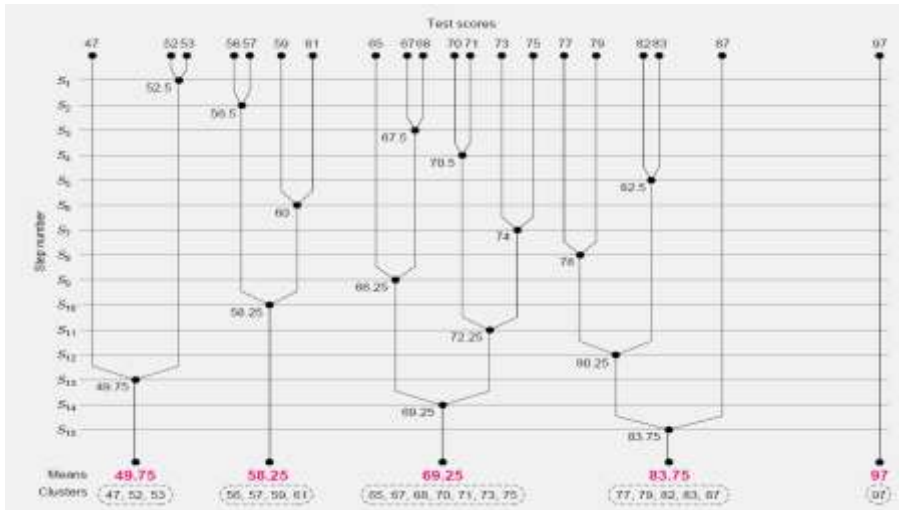
AGNES (Agglomerative Nesting)

- خوشه‌ها را مکرراً با هم ترکیب می‌کند. به این صورت که ابتدا هر یک از اشیاء را در داخل یک خوشه قرار می‌دهد و سپس این خوشه‌ها را با ترکیب کردن به خوشه‌های بزرگ و بزرگتر تبدیل می‌کند تا اینکه همه اشیاء در یک خوشه قرار گیرند و یا به شرط پایان برسد.



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مثال Agglomerative Clustering



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

(top-down) Divisive

- خوشه‌ها را مکرراً تقسیم می‌کند. دقیقاً بر عکس روش ترکیبی عمل می‌کند به این صورت که ابتدا همه اشیاء در یک خوشه قرار دارند و الگوریتم این خوشه را به خوشه‌های کوچک و کوچکتر تجزیه می‌کند تا اینکه هرشیء در یک خوشه قرار گیرد. این روش معمولاً مناسب نیست و خیلی کم مورد استفاده قرار می‌گیرد زیرا پیچیدگی محاسباتش بالاست

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

فاصله در خوشه بندی سلسله مراتبی

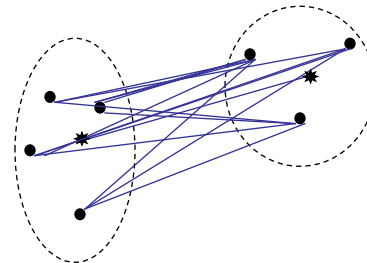
- محاسبه فاصله ها مهم می باشد
- روش محاسبه فاصله نیز مهم است

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

فاصله در خوشه بندی سلسله مراتبی

- معیارهای گوناگونی که در روشهای سلسله مراتبی برای فاصله بین خوشه ها بکار می رود ، عبارتند از :

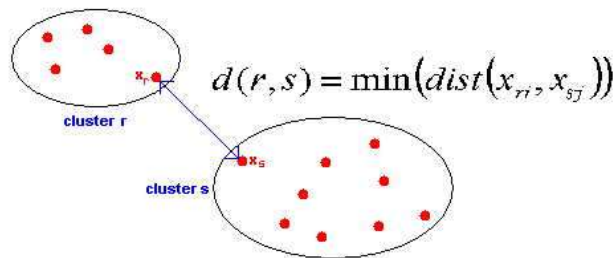
- **Single Link**
- **Complete Link**
- **Average Link**
- **Centroid**



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

Single linkage

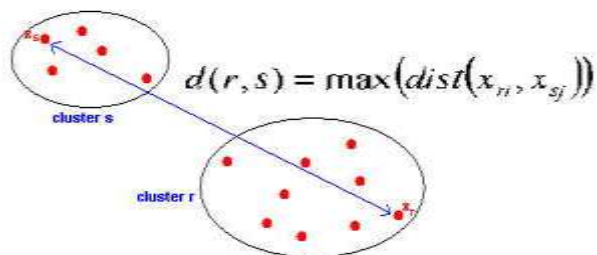
- فاصله بين دسته‌ها بر حسب حداقل فاصله ممكنه بين عناصر آنها محاسبه مي‌شود.
- كليه فاصله بين زوجهاي عناصر دو دسته محاسبه شده و حداقل آنها فاصله بين دو دسته را تعيين مي‌كند.



سميه عليزاده هيات علمي دانشكده صنايع
دانشگاه خواجه نصير طوسي

Complete linkage

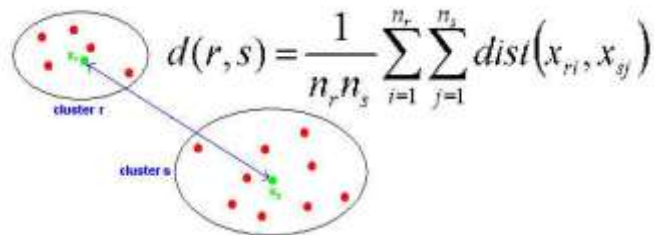
- فاصله بين دسته‌ها بر حسب دورترين فاصله ممكنه بين عناصر آنها محاسبه مي‌شود.



سميه عليزاده هيات علمي دانشكده صنايع
دانشگاه خواجه نصير طوسي

Average linkage

- فاصله بین دو دسته مساوی مقادیر متوسط کلیه فاصله‌های ممکنه بین عناصر دو دسته است.

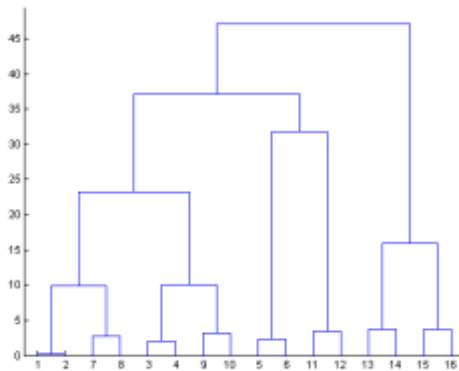


$$d(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} \text{dist}(x_{ri}, x_{sj})$$

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

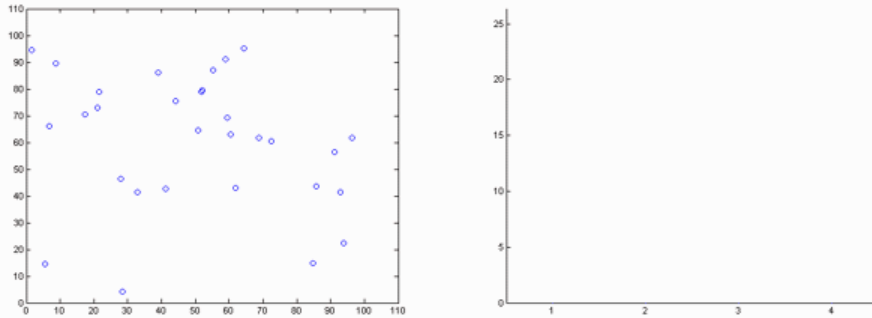
خوشه بندی سلسله مراتبی

- به نمودار درختی تشکیل شده **dendrogram** گفته می‌شود.



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

خوشه بندی سلسله مراتبی



سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی

مقایسه خوشه بندی سلسله مراتبی و غیر سلسله مراتبی

- روشهای خوشه بندی غیر سلسله مراتبی معمولاً سریعتر عمل می کنند ولی نیاز به یکسری تصمیم گیری از طرف تحلیل گر و استفاده کننده دارد.
- انتخاب تعداد خوشه ها از این گونه تصمیم گیریها می باشد.
- در این گونه روشها معمولاً یکسری خوشه های اولیه ایجاد شده و سپس در مراحل بعدی بهبود صورت می گیرد.
- از آنجاییکه در این روشها مناسب بودن خوشه ها به تعداد خوشه ها و یا حتی خوشه های اولیه بستگی دارد این روشها کمتر از روشهای خوشه بندی سلسله مراتبی انجام می شوند.

سمیه علیزاده هیات علمی دانشکده صنایع
دانشگاه خواجه نصیر طوسی