

## مباحث تکمیلی در خوشه بندی

### تعیین K – تعداد خوشه ها

- اگر تمام متغیرها کاملا مستقل باشند هیچ خوشه ای ایجاد نمی شود (تمام فضا بصورت تصادفی با نقاط داده ها پر می شود) بر عکس اگر تمام متغیرها وابسته باشند آنگاه تمام داده ها تشکیل یک خوشه می دهند.
- در شرایط بین استقلال و وابستگی کامل ما نمی دانیم که واقعا چند خوشه وجود دارد.
- معمولا در انتخاب مقدار K نقش شخص تحلیل گر بسیار بیشتر از کامپیوتر می باشد.
- برای همین با توجه به کاربردهای متفاوت روشهای خوشه بندی و در جاهای متفاوت ممکن است به تعداد بیشتر یا کمتری از خوشه ها نیاز باشد.

## تعیین K – تعداد خوشه ها (ادامه)

- در بسیاری از موارد با يك مقدار K خوشه بندي را انجام داده ، نتایج را بررسی می کنند و دوباره به سراغ يك K دیگر می روند.
- بعد از هر بار انجام این کار قدرت و ارزش نتایج را بوسیله اندازه گیری میزان متوسط فواصل در داخل خوشه ها و میزان متوسط فواصل بین مراکز خوشه ها و بسیاری روشهای دیگر بررسی می کنند.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## تعیین K – تعداد خوشه ها (ادامه)

- تعداد خوشه ها را می توان به اندازه تعداد داده های داخل پایگاه داده در نظر گرفت که در این حالت داده های هر خوشه دقیقا به هم متشابه هستند و داده های خوشه های متفاوت هم با هم تفاوت دارند (در هر خوشه تنها يك داده وجود دارد).
- ولي در این حالت دیگر استفاده ای از خوشه بندي نمی توان کرد زیرا علت استفاده از خوشه بندي یافتن الگوهای مناسب در پایگاه داده و تلخیص آنها بمنظور فهم بهتر پایگاه داده ها می باشد.
- پس تعداد خوشه ها از تعداد داده ها کمتر می باشد .
- اگر تعداد خوشه ها یکی در نظر گرفته شود تمام داده ها در واقع در آن خوشه قرار دارند و اطلاعات مفیدی کسب نمی شود .

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## تعیین K – تعداد خوشه ها (ادامه)

مزیت خوشه بندی سلسله مراتبی این است که به استفاده کننده و تحلیل گر اجازه می دهد که از بین حالات مختلف یک عدد برای تعداد خوشه ها انتخاب نماید.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## ارزیابی کیفیت خوشه بندی ۱

• دو سؤال اصلی:

- چه تعداد خوشه مناسب است؟
- کدام تخصیص بهترین تخصیص نمونه ها به خوشه ها است؟

➤ دو مشکل عمده:

- تنظیمات اولیه پارامترهای الگوریتم ها
- مقدار دهی اولیه الگوریتم های خوشه بندی

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## ارزیابی کیفیت خوشه‌بندی ۲

- کمترین شباهت بین خوشه‌های و بیشترین شباهت درون خوشه‌های

شاخص دیویس-بولدین

- داشتن خوشه‌های متراکم با مرزهای مشخص

شاخص دان

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## شاخص دان (داشتن خوشه‌های متراکم با مرزهای مشخص)

- در شاخص دان دو معیار در نظر گرفته میشود:

• حداکثر فاصله درون خوشه ای **Intercluster**

• حداقل فاصله بیرون خوشه ای **Intracluster**

$$D = \frac{d_{min}}{d_{max}},$$

- صورت کسر حداقل فاصله بین دو عنصر از دو خوشه متفاوت را بررسی می کند.
- مخرج کسر بیشترین فاصله بین دو عنصر درون یک خوشه را نشان میدهد.
- **D** در بازه  $[0, \infty]$  است و هر چه بیشتر باشد بهتر است.

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## شاخص دان (داشتن خوشه‌های متراکم با مرزهای مشخص)

$$D_{nc} = \min_{t=1, \dots, nc} \left\{ \min_{j=t+1, \dots, nc} \left( \frac{d(c_t, c_j)}{\max_{k=1, \dots, nc} \text{diam}(c_k)} \right) \right\}$$

$$d(c_t, c_j) = \min_{x \in c_t, y \in c_j} d(x, y),$$

$$\text{diam}(C) = \max_{x, y \in C} d(x, y)$$

$$\max_{nc=1, 2, \dots, n} \{D_{nc}\}$$

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## شاخص دیویس-بولدین 1

- شاخص شباهت بین خوشه‌ها و بر اساس شاخص پراکندگی و شاخص تفاوت بین دو خوشه تعریف می‌شود.

$$R_{ij} \geq 0.$$

$$R_{ij} = R_{ji}.$$

$$\text{if } s_i = 0 \text{ and } s_j = 0, \text{ then } R_{ij} = 0.$$

$$\text{if } s_j > s_k \text{ and } d_{ij} = d_{ik} \text{ then } R_{ij} > R_{ik}.$$

$$\text{if } s_j = s_k \text{ and } d_{ij} < d_{ik} \text{ then } R_{ij} > R_{ik}.$$

$$s_i = \frac{1}{C_i} \sum_{x \in C_i} \|x - z_i\|,$$

$$d_{ij} = \|z_i - z_j\|,$$

$$R_{ij} = \frac{(s_i + s_j)}{d_{ij}}$$

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی

## شاخص دیویس-بولدین ۲

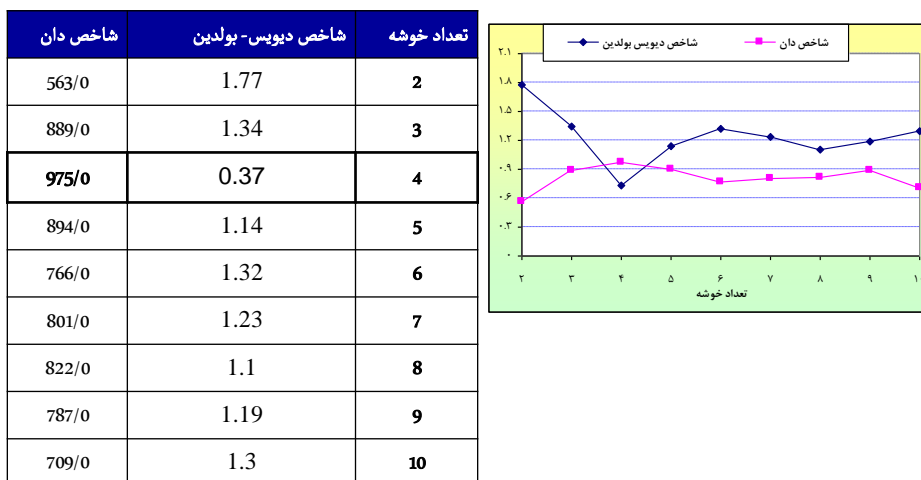
• میانگین شباهت بین هر خوشه و شبیه‌ترین خوشه به آن می‌باشد.

$$R_i = \max_{i=1, \dots, n_c, i \neq j} R_{ij},$$

$$DB_{nc} = \frac{1}{n_c} \sum_{i=1}^{n_c} R_i,$$

$$\frac{1}{C} \sum_{k=1}^c \max_{l \neq k} \left\{ \frac{S_c(Q_k) + S_c(Q_l)}{d_{ce}(Q_k, Q_l)} \right\}$$

سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی



سمیه علیزاده هیات علمی دانشکده صنایع  
دانشگاه خواجه نصیر طوسی