

Mühazirə 10: Klasterləşdirmə metodları

tex.f.d.,dos. Yadigar İmamverdiyev 2021, payız semestri 01.12.2021



Mühazirənin planı

- Giriş
- Klasterləşdirmə metodlarının təsnifatı
- K-means metodu
- İyerarxik klasterləşdirmə
- Klasterləşdirmənin keyfiyyətinin qiymətləndirilməsi
- Bu mühazirə təqdimatı hazırlanarkan İnternetda alyetar çox sayda taqdimatdan istifada edilmişdir. Taassüf ki, onların har birinin müallifini qeyd etmak mümkün olmadı.
- Bu müəlliflərin hər birinə dərindən təşəkkür edirəm.



Klasterləşdirmənin ümumi tərifi

Klasterləşdirmə alqoritmləri öz xarakteristikalarına görə **yaxın obyektləri** bir klasterdə, **uzaq obyektləri** isə müxtəlif klasterlərdə yerləşdirməklə obyektlərin verilmiş çoxluğunu **qruplara** (klasterlərə) bölür.

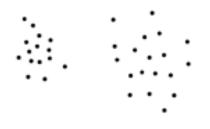


Klasterləşdirmənin məqsədləri

- Verilənlərin sonrakı emalını asanlaşdırmaq, obyektlər çoxluğunu oxşar obyektlər qruplarına ayırmaq və hər bir qrup ilə ayrılıqda işləmək (klassifikasiya, reqresiya, proqnozlaşdırma məsələləri);
- Hər bir klasterdən bir nümayəndə saxlamaqla verilənlərin həcmini azaltmaq (verilənlərin sıxılmasl məsələsi);
- Klasterlərin heç birinə yaxın düşməyən qeyri-tipik obyektləri ayırmaq (birsinifli klassifikasiya məsələsi)
- Obyektlər çoxluğunun iyerarxiyasını qurmaq (taksonomiya məsələsi).



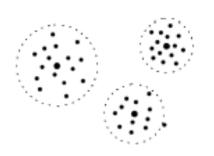
Klaster strukturlarının tipləri



Klasterlərarası məsafələr, bir qayda olaraq, klasterdaxili məsafələrdən böyükdür.



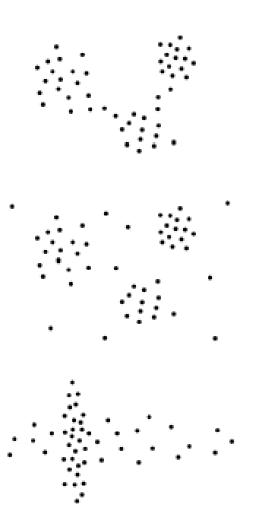
Lentşəkilli klasterlər



Mərkəzi olan klasterlər



Klaster strukturlarının tipləri



Klasterlər körpülərlə birləşə bilərlər.

Klasterlər çox seyrək yerləşmiş obyektlərin seyrək fonu üzərində yerləşə bilərlər.

Klasterlər kəsişə bilərlər.



Klaster strukturlarının tipləri



Klasterlər oxşarlığa görə deyil, müntəzəmliyin digər növünə görə yarana bilər.

Klasterlər, ümumiyyətləyə, olmaya bilərlər.

- Hər bir klasterləşdirmə metodunun öz məhdudiyyətləri var və klasterlərin yalnız müəyyən tiplərini ayırırlar.
- "Klaster strukturunun tipi" anlayışı da metoddan asılıdır və formal tərifi yoxdur.



Məsələnin formal qoyuluşu

- Test nümunələrinin $X = \{x_1, x_2, ..., x_n\}$ toplusu və nümunələr arasında *məsafə funksiyası* var.
- X-i kəsişməyən altçoxluqlara (klasterlərə) elə bölmək tələb edilir ki, hər bir altçoxluq oxşar obyektlərdən ibarət olsun, müxtəlif altçoxluqların obyektləri isə əhəmiyyətli dərəcədə fərqlənsinlər.



Klasterləşdirmə metodlarının təsnifatı

- Bölünmə (Partitioning) əsasında klasterləşdirmə
- İyerarxik klasterləşdirmə
- Sıxlığa əsaslanan klasterləşdirmə (EM, DBSCAN)



k-ortalar (k means) metodu



k-ortalar metodu (k-means)

Fərz edək ki, obyektlər çoxluğu artıq müəyyən qaydada K qrupa (klasterə) bölünüb, C_k ilə k-ci klasteri işarə edək. Tutaq ki, C(k) - k-cı obyektin aid olduğu qrupun nömrəsinə bərabərdir.

k- c_l klasterə aid olan nöqtələr sisteminin ağırlıq mərkəzini m_k ilə işarə edək:

$$m_k = \frac{1}{N_k} \sum_{x_i \in C_k} x_i, k = 1, ..., K$$

Məqsəd: $\min_{C,m_k} \sum_{i=1}^n \rho(x_i,m_k)$



k-ortalar (k means) metodu

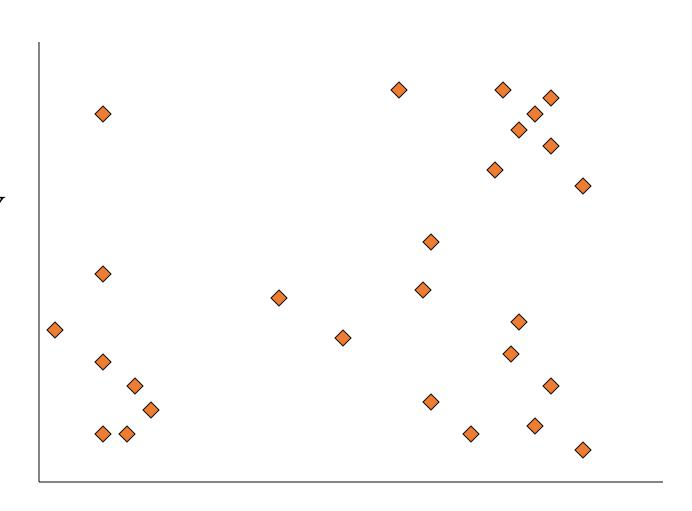
- (1) Klasterlərin sayını (*K*) qərarlaşdırıb onların başlanğıc ağırlıq mərkəzlərini təsadüfi seçirik və obyektləri yaxın olduqları mərkəzə uyğun klasterlərə yerləşdiririk.
- ② Klasterlərin m_k ağırlıq mərkəzlərini yenidən hesablayırıq.
- ③ x_i nöqtəsindən bütün m_k -lara $\rho(x_i, m_k)$ məsafələrini hesablayırıq.

- 4 Addım 3-ü bütün x_i , i=1,...,n obyektləri üçün təkrarlayırıq.
- \bigcirc Əgər heç olmasa, bir C_k klasteri dəyişirsə, onda $Addım\ 2$ -yə keçirilir, əks halda klasterləşdirmə prosesi sona çatır.



k-ortalar alqoritminə misal (1)

Neçə klasterin olmasını qərarlaşdırırıq. Tutaq ki, *k*=3.

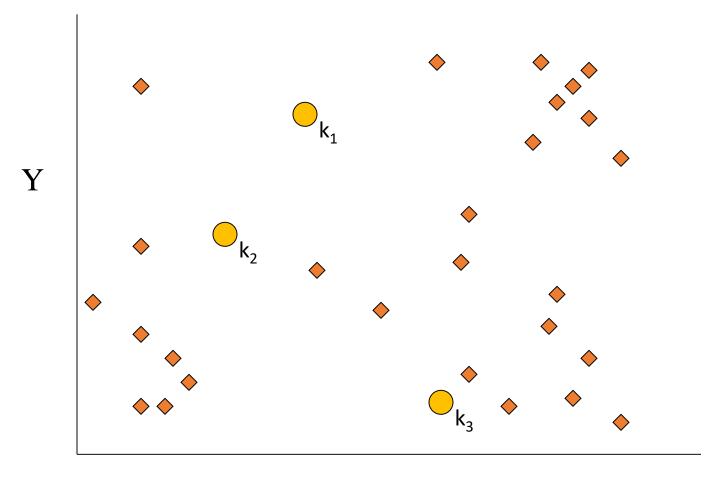


X



k-ortalar alqoritminə misal (2)

3 klasterin ağırlıq mərkəzini təsadüfi seçirik.

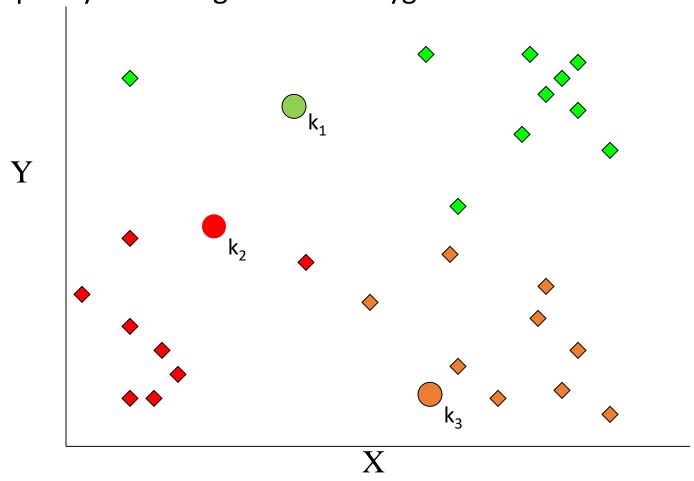


X



k-ortalar alqoritminə misal(3)

Seçilmiş məsafə metrikası ilə məsafələri hesablayırıq. Hər bir nöqtəni yaxın olduğu mərkəzə uyğun klasterə aid edirik.

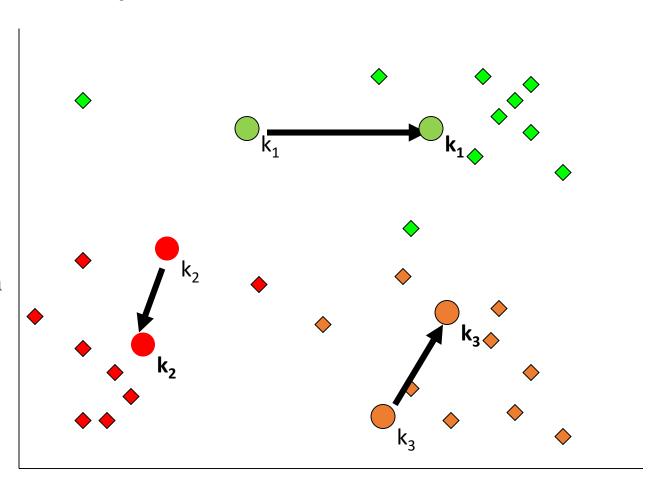




k-ortalar alqoritminə misal (4)

Y

Hər bir klasterin mərkəzini yenidən hesablayırıq və mərkəzin yerini dəyişirik.



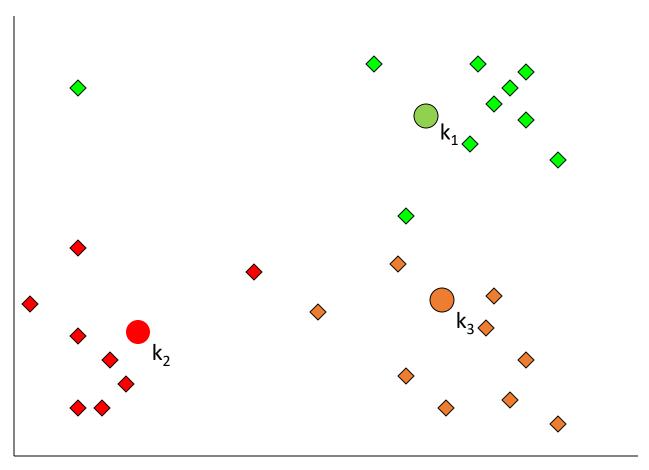
X



k-ortalar alqoritminə misal (5)

Nöqtələrdən mərkəzlərə məsafələr yenidən hesablanır və mərkəzlərə ən yaxın nöqtələr yenidən paylanır

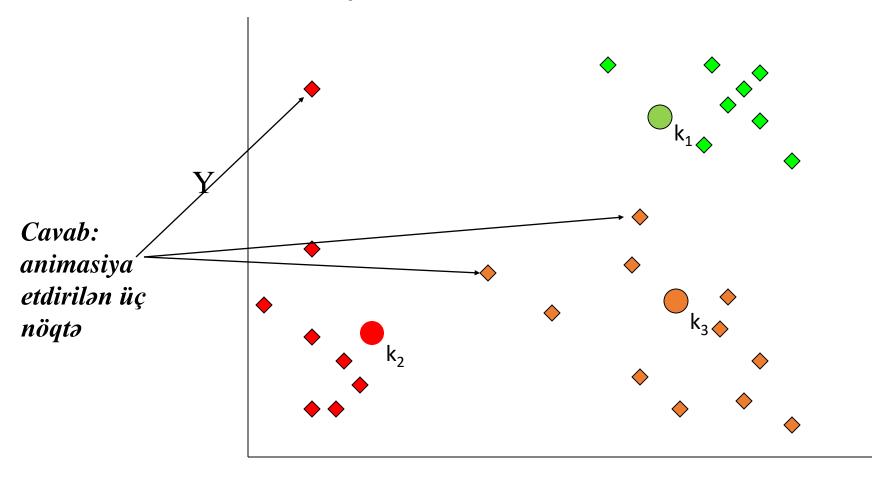
Sual: Hansı nöqtələr yenidən paylandı?







k-ortalar alqoritminə misal (6)



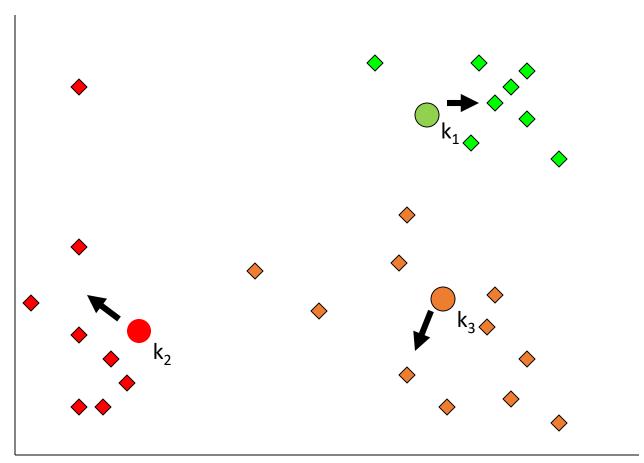
X



k-ortalar alqoritminə misal (7)

Y

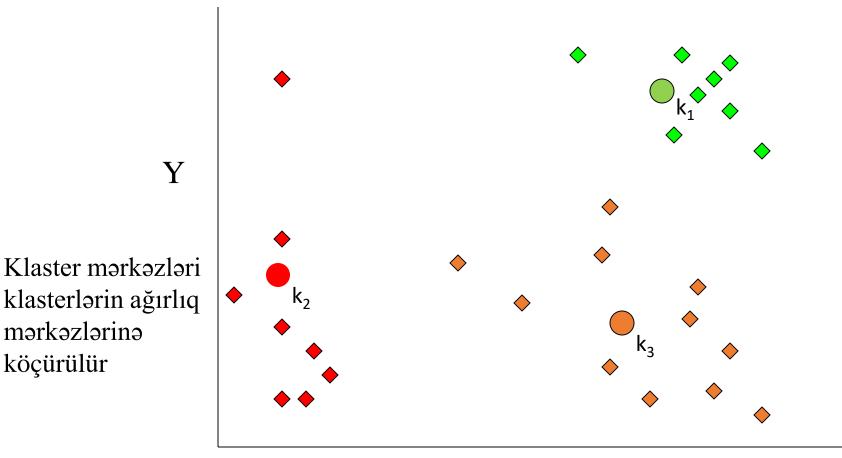
Klasterlərin ağırlıq mərkəzləri yenidən hesablanır



X



k-ortalar alqoritminə misal (8)



klasterlərin ağırlıq mərkəzlərinə köçürülür



k-ortalar alqoritminin üstünlükləri

- Sadədir, başa düşüləndir
- Mürəkkəbliyi aşağıdır
 - $M\ddot{u}r \ni kk \ni blik O(nkt)$, burada t iterasiyaların sayıdır



k-ortalar alqoritminin nöqsanları

- Kvadratik meyllərin cəminin qlobal minimumunun tapılmasına zəmanət vermir, lokal minimumlardan biri tapılır
- Nəticə klasterlərin başlanğıc mərkəzlərinin seçilməsindən asılıdır, onların optimal seçimi hələ də məlum deyil.
- Klasterlərin sayını əvvəlcədən bilmək lazımdır.



Universiteti K-ortalar alqoritminin nöqsanlarının aradan qaldırılması

- (1) Bir neçə təsadüfi klasterləşdirmə edilir: keyfiyyət funksionalına görə ən yaxşısı seçilir.
- \bigcirc Klasterlərin k sayı tədricən artırılır.



Medoidlər metodu (k-medoids)

- K-medoids ədədi ortanın əvəzinə hər bir klasterin medianı istifadə edilir.
 - 1, 3, 5, 7, 9-un ədədi ortası: 5
 - 1, 3, 5, 7, 1009 -un ədədi ortası: **205**
 - 1, 3, 5, 7, 1009 -un medianı: **5**
 - Medianın üstünlüyü: ekstremal qiymətlərin təsiri yoxdur



Medoidlər metodu (k-medoids)

- Hər bir klasterdə ağırlıq mərkəzlərinin əvəzinə medoid hesablanır.
- Medoid verilənlər çoxluğunun və ya klasterin elə obyektidir ki, digər obyektlərdən ona qədər olan orta məsafə minimaldır.
- Obyektlərin klasterdən klasterə yerdəyişməsi alqoritmi *k*-ortalar alqoritminə analojidir.
- Medoid obyektdir, buna görə bütün obyektlər arasındakı məsafələr matrisini bilmək kifayətdir (üstünlük!)



İyerarxik klasterləşdirmə

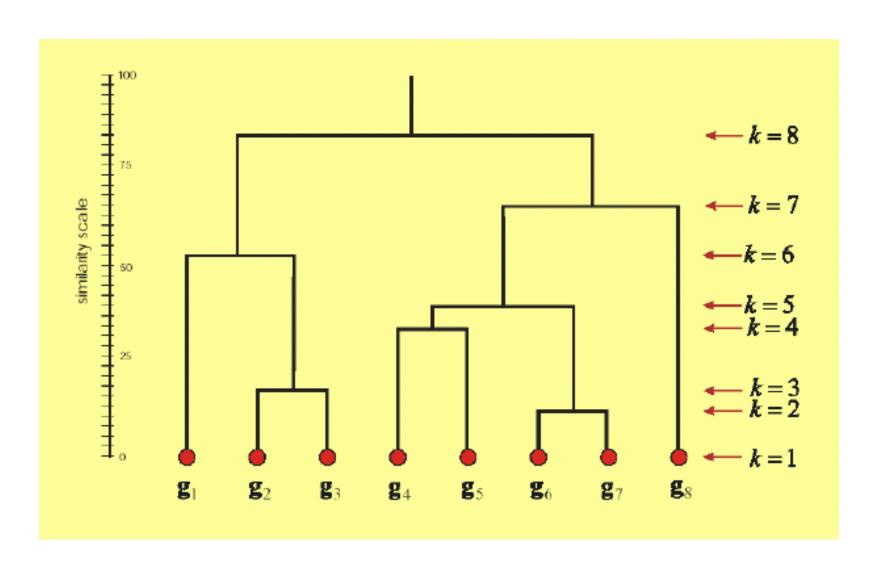


İyerarxik klasterləşdirmə (taksonomiya)

- Əsas ideya: aşağı səviyyələrdəki klasterlər yuxarıdakı klasterləri bölməklə alınır.
- Klassifikasiyanın təpəsində bütün obyektləri əhatə edən bir klaster durur.
- Ən aşağı səviyyədə hər birində bir obyekt olan n klaster olur.
- Belə iyerarxik sturukturları kök ağacları (dendroqramlar) şəklində göstərmək əlverişlidir.



İyerarxik Klasterləşdirmə





İyerarxik klasterləşdirmə alqoritmlərinin üstünlüyü

- İyerarxik klasterləşdirmə alqoritmlərinin girişinə klasterlərin sayını vermək lazım deyil.
- Dendroqramın əsasında istifadəçi özü müəyyən səviyyədə onu kəsərək müəyyən sayda klasterlər ala bilər.



İyerarxik klasterləşdirmə alqoritmlərinin iki sinfi

- İyerarxik klasterləşdirmə alqoritmlərinin iki növü vardır:
 - Aqqlomerativ ("aşağıdan yuxarıya"):
 - Bir elementdən ibarət çoxluqlardan başlayırlar
 - Onlar bütün obyektlər bir klasterdə olana kimi birləşdirilir.
 - Çox geniş yayılmış metoddur.
 - Diviziv ("yuxarıdan aşağıya"):
 - Obyektlər çoxluğu tək elementlərdən ibarət klasterlər alınana kimi rekursiv bölünür.

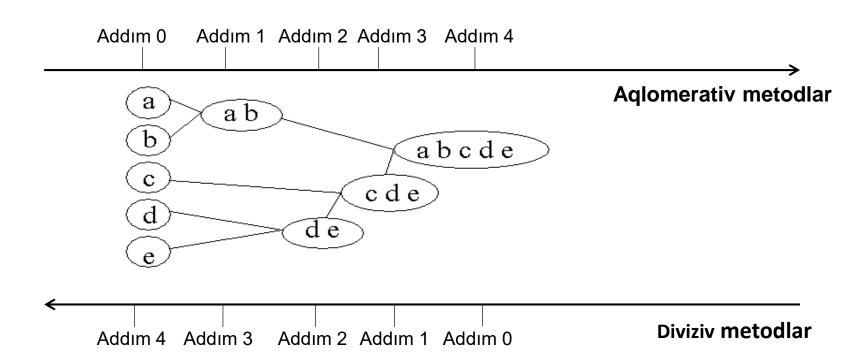


Aqlomerativ metodlar

- Başlanğıc elementlərin ardıcıl birləşdirilməsi və müvafiq olaraq klasterlərin sayının azalması
- Başlanğıcda bütün obyektlər ayrıca klasterlər olur.
- Birinci addımda ən oxşar obyektlər klasterdə birləşdirilir.
- Sonrakı addımlarda birləşdirmə o vaxta kimi davam etdirilir ki, bütün obyektlər bir klaster təşkil etsin.



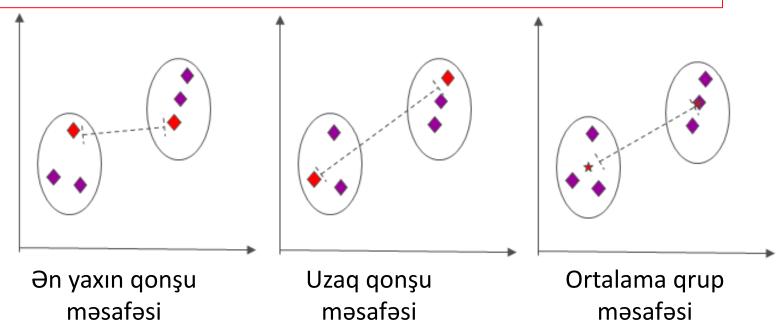
Aqlomerativ metodlar





Klasterlər arasındakı məsafə (1)

- Aqlomerativ metodlarda klasterləri birləşdirmək üçün klasterlər arasında məsafə funksiyası müəyyən edilməlidir.
- Üç yanaşma geniş istifadə edilir.





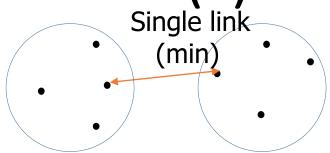
Klasterlər arasındakı məsafə (2)

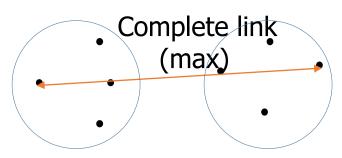
- ① Ən yaxın qonşu məsafəsi: $D(C_i, C_j) = \min \{ \rho(x, z) | x \in C_i, z \in C_j \}$
- 2 Uzaq qonşu məsafəsi:

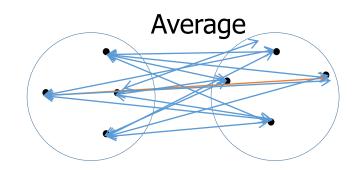
$$D(C_i, C_j) = \max \{ \rho(x, z) | x \in C_i, z \in C_j \}$$

3 Ortalama qrup məsafəsi:

$$D(C_{i}, C_{j}) = \frac{1}{|C_{i}||C_{j}|} \sum_{x \in C_{i}} \sum_{z \in C_{j}} \rho(x, z)$$









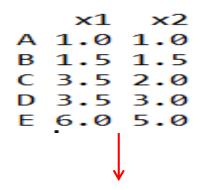
Aqlomerativ metodlar

- *n* verilənlər obyekti verilib.
- İyerarxik aqlomerativ klasterləşdirmə alqoritmi aşağıdakı addımlarla həyata keçirilir:
- **Addım 1.** n verilənlər obyekti üçün məsafələr matrisi hesablanır
- Addım 2. Hər bir obyekt klaster kimi götürülür
- Addım 3. Klasterlərin sayı 1 olana kimi təkrarlanır
 - Addım 3.1. İki ən yaxın klaster birləşdirilir
 - Addım 3.2. Məsafə matrisi yenilənir



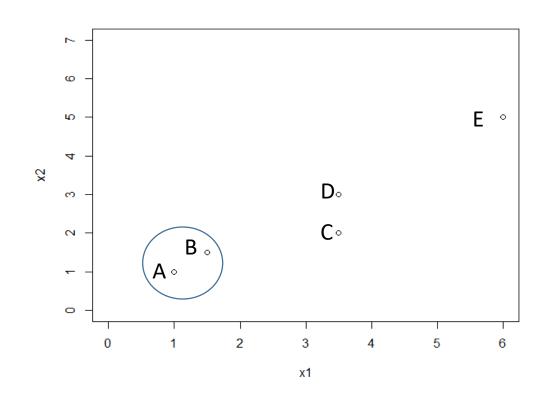
İyerarxik aqlomerativ klasterləşdirmə - Misal

Misal. 5 obyekt verilib:



Məsafə Matrisi

| | \boldsymbol{A} | B | C | D | $\boldsymbol{\mathit{E}}$ |
|------------------|------------------|--------|------|------|---------------------------|
| A | 0 | (0.71) | 2.69 | 3.20 | 6.4 |
| В | 0.71 | 0 | 2.06 | 2.6 | 5.7 |
| C | 2.69 | 3.2 | 0 | 1 | 3.9 |
| D | 3.20 | 2.6 | 1 | 0 | 3.2 |
| \boldsymbol{E} | 6.4 | 5.7 | 3.9 | 3.2 | 0 |

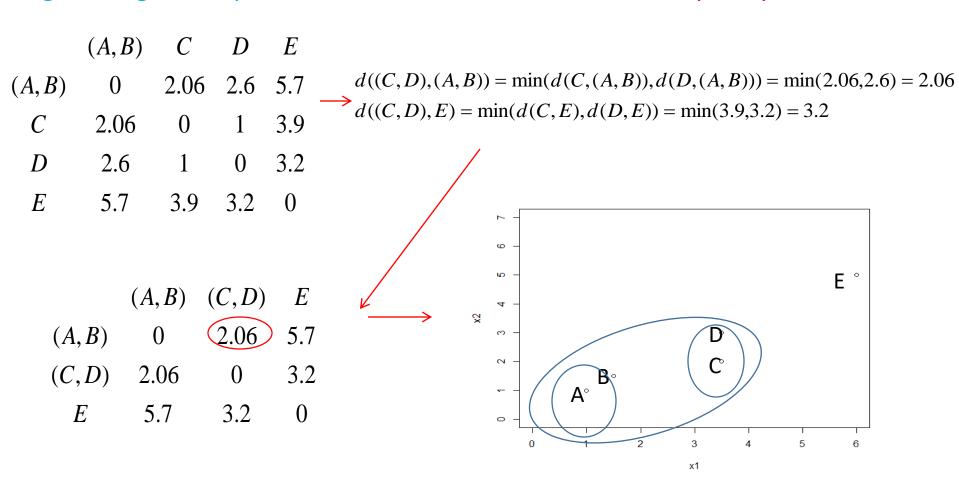




Iyerarxik aqlomerativ klasterləşdirmə — Misal Single Link(age) funksiyasından istifadə edərək, məsafə matrisini yeniləyirik.



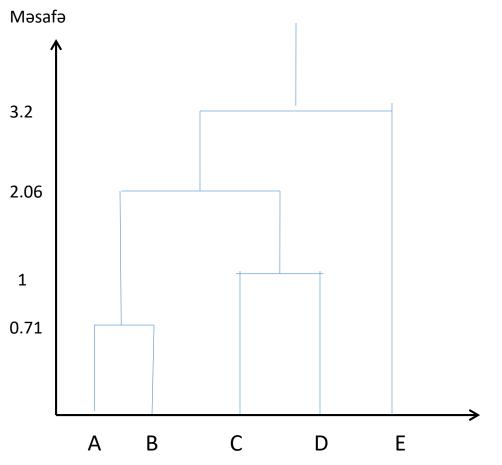
İyerarxik Aqlomerativ Klasterləşdirmə — Misal Single Linkage funksiyasından istifadə edərək, məsafə matrisini yeniləyirik.





İyerarxik aqlomerativ klasterləşdirmə – Misal

Dendrogram

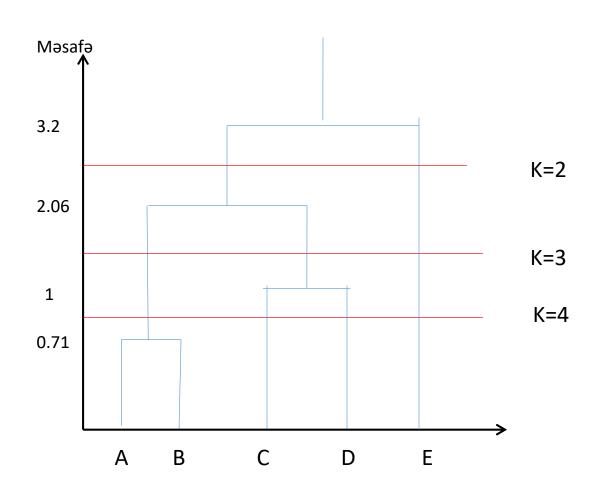


- 1. Başlanğıcda 5 klaster var.
- 2. A və B klasterlərini 0.71 məsafəsində (A, B) klasterində birləşdiririk.
- 3. C və D klasterlərini 1 məsafəsində (C, D) klasterində birləşdiririk.
- 4. (A,B) və (C, D) klasterlərini 2.06 məsafəsində ((A, B), (C, D)) klasterində birləşdiririk.
- 5. ((A, B), (C, D)) və E klasterlərini 3.2 məsafəsində birləşdiririk.
- 6. Axırıncı klasterə bütün obyektlər daxildir, bununla da hesablamalar dayanır.



İyerarxik aqlomerativ klasterləşdirmə – Misal

• Klasterlərin sayını necə qərarlaşdırırıq? Ağacı haradasa kəsirik.





Diviziv ("yuxarıdan aşağıya") metodlar

Ağac kökdən yarpaqlara doğru qurulur.

Yanaşmalardan biri:

- ①Birinci addımda obyektlər çoxluğuna hər hansı alqoritm(ağırlıq mərkəzləri, medoidləri) tətbiq edərək bu çoxluğu iki klasterə bölürlər
- 2 Sonra alınmış klasterlərdən hər birini bölürlər və s.



Klasterləşdirmənin keyfiyyətinin qiymətləndirilməsi



Klasterləşdirmənin keyfiyyətinin qiymətləndirilməsi

- Vahid (hamılıqla qəbul edilən, bütün hallarda tətbiq edilə bilən) qiymətləndirmə metodu yoxdur.
- Qiymətləndirmə nəzərdə tutur ki, toplu (və ya toplunun bir hissəsi) insan tərəfindən işarələnib.
 - Klasterlər klasterləşdirmənin nəticəsidir.
 - Siniflər insan tərəfindən işarələnmənin nəticəsidir.



Səhvlər matrisi (Confusion matrix)

- n = obyektlarin say
- m_i = i klasterində obyektlərin sayı
- $c_i = j sinfində obyektlərin sayı$
- n_{ij} = i klasterində j sinfindən olan obyektlərin sayı
- $p_{ij} = n_{ij}/m_i$ = i klasterində olan obyektin j sinfinə mənsub edilməsi ehtimalı

| | Class 1 | Class 2 | Class 3 | |
|-----------|----------|----------|----------|-------|
| Cluster 1 | n_{11} | n_{12} | n_{13} | m_1 |
| Cluster 2 | n_{21} | n_{22} | n_{23} | m_2 |
| Cluster 3 | n_{31} | n_{32} | n_{33} | m_3 |
| | c_1 | c_2 | c_3 | n |

| | Class 1 | Class 2 | Class 3 | |
|-----------|----------|----------|----------|-------|
| Cluster 1 | p_{11} | p_{12} | p_{13} | m_1 |
| Cluster 2 | p_{21} | p_{22} | p_{23} | m_2 |
| Cluster 3 | p_{31} | p_{32} | p_{33} | m_3 |
| | c_1 | c_2 | c_3 | n |



Metrikalar

| | Class 1 | Class 2 | Class 3 | |
|-----------|----------|----------|----------|-------|
| Cluster 1 | p_{11} | p_{12} | p_{13} | m_1 |
| Cluster 2 | p_{21} | p_{22} | p_{23} | m_2 |
| Cluster 3 | p_{31} | p_{32} | p_{33} | m_3 |
| | c_1 | c_2 | c_3 | n |

Entropiya:

- i klasterinin entropiyası: $e_i = -\sum_{j=1}^L p_{ij} \log p_{ij}$
 - Müntəzəm paylandıqda ən yüksəkdir, bir sinif olduqda sıfra bərabərdir
- Bütün klasterlərin entropiyası: $e = \sum_{i=1}^{K} \frac{m_i}{n} e_i$
- Saflıq (ing. Purity):
 - i klasterinin saflığı: $p_i = \max_j p_{ij}$
 - Bütün klasterlərin saflığı: $p(C) = \sum_{i=1}^K \frac{m_i}{n} p_i$



Metrikalar

| | Class 1 | Class 2 | Class 3 | |
|-----------|----------|----------|----------|-------|
| Cluster 1 | p_{11} | p_{12} | p_{13} | m_1 |
| Cluster 2 | p_{21} | p_{22} | p_{23} | m_2 |
| Cluster 3 | p_{31} | p_{32} | p_{33} | m_3 |
| | c_1 | c_2 | c_3 | n |

- Precision:
 - i klasterinin j sinfinə nəzərən: $Prec(i,j) = p_{ij}$
- Recall:
 - i klasterinin j sinfinə nəzərən: $Rec(i,j) = \frac{n_{ij}}{c_j}$
- F-measure:
 - Precision və Recall-un Harmonik ortasıdır:

$$F(i,j) = \frac{2 * Prec(i,j) * Rec(i,j)}{Prec(i,j) + Rec(i,j)}$$



Metrikalar

| | Class 1 | Class 2 | Class 3 | |
|-----------|----------|-----------------|----------|-------|
| Cluster 1 | n_{11} | n_{12} | n_{13} | m_1 |
| Cluster 2 | n_{21} | n_{22} | n_{23} | m_2 |
| Cluster 3 | n_{31} | n ₃₂ | n_{33} | m_3 |
| | c_1 | c_2 | c_3 | n |

Klasterlər və klasterləşdirmə üçün Precision/Recall

- i klasterinə elə k_i sinfi təyin olunur ki, $k_i = \arg\max_i n_{ij}$
- **Precision:**
 - i klasteri: $Prec(i) = \frac{n_{ik_i}}{m_i}$
 - Bütün klasterlər üçün: $Prec(C) = \sum_{i} \frac{m_i}{n} Prec(i)$
- Recall:
 - i klasteri: $Rec(i) = \frac{n_{ik_i}}{c_{k_i}}$
 - Bütün klasterlər üçün: $Rec(C) = \sum_{i} \frac{m_i}{n} Rec(i)$
- F-measure:
 - Precision va Recall-un Harmonik ortası



Metrikalar barəsində bəzi qeydlər

Entropiya – sinfin klasterlər üzrə «yayılmasıdır». Nə qədər kiçikdirsə, bir o qədər yaxşıdır, ideal halda *Entropy*=0.

Saflığın qiyməti nə qədər böyükdürsə, bir o qədər yaxşıdır. İdeal halda, *Purity*=1.

F-ölçü klasterləşdirmənin ümumi keyfiyyətini göstərir, lakin klasterlərin özlərinin quruluşunu göstərmir.



Yaxşı və pis klasterləşdirmə: Misal

| | Class 1 | Class 2 | Class 3 | |
|-----------|---------|---------|---------|-----|
| Cluster 1 | 2 | 3 | 85 | 90 |
| Cluster 2 | 90 | 12 | 8 | 110 |
| Cluster 3 | 8 | 85 | 7 | 100 |
| | 100 | 100 | 100 | 300 |

| | Class 1 | Class 2 | Class 3 | |
|-----------|---------|---------|---------|-----|
| Cluster 1 | 20 | 35 | 35 | 90 |
| Cluster 2 | 30 | 42 | 38 | 110 |
| Cluster 3 | 38 | 35 | 27 | 100 |
| | 100 | 100 | 100 | 300 |

Purity: (0.94, 0.81, 0.85)

– Bütün klasterlər üçün: 0.86

Precision: (0.94, 0.81, 0.85)

Bütün klasterlər üçün: 0.86

Recall: (0.85, 0.9, 0.85)

- Bütün klasterlər üçün: 0.87

Purity: (0.38, 0.38, 0.38)

Bütün klasterlər üçün: 0.38

Precision: (0.38, 0.38, 0.38)

Bütün klasterlər üçün: 0.38

Recall: (0.35, 0.42, 0.38)

Bütün klasterlər üçün: 0.39