

VERİ MADENCİLİĞİ

(Kümeleme)

Yrd.Doç.Dr. Kadriye ERGÜN
kergun@balikesir.edu.tr

İçerik

- Kümeleme İşlemleri
- Kümeleme Tanımı
- Kümeleme Uygulamaları
- Kümeleme Yöntemleri

Kümeleme (Clustering)

- Kümeleme birbirine benzeyen veri parçalarını ayırma işlemidir ve kümeleme yöntemlerinin çoğu veri arasındaki uzaklıkları kullanır.
- Nesneleri kümelere (gruplara) ayırma
- Küme: birbirine benzeyen nesnelere oluşan grup
 - Aynı kümedeki nesnelere birbirine daha çok benzer
 - Farklı kümedeki nesnelere birbirine daha az benzer

Kümeleme

- Danışmansız öğrenme: Hangi nesnenin hangi sınıfa ait olduğu ve sınıf sayısı belli değil
- Uygulamaları:
 - verinin dağılımını anlama
 - başka veri madenciliği uygulamaları için ön hazırlık

Kümeleme Uygulamaları

- Örüntü tanıma
- Görüntü işleme
- Ekonomi
- Aykırılıkları belirleme
- WWW
 - Doküman kümeleme
 - Kullanıcı davranışlarını kümeleme
 - Kullanıcıları kümeleme
- Diğer veri madenciliği algoritmaları için bir ön işleme adımı
- Veri azaltma – küme içindeki nesnelere temsil edilmesi için küme merkezlerinin kullanılması

Veri Madenciliğinde Kümeleme

- Ölçeklenebilirlik
- Farklı tipteki niteliklerden oluşan nesnelere kümeleme
- Farklı şekillerdeki kümeleri oluşturabilme
- En az sayıda giriş parametresi gereksinimi
- Hatalı veriler ve aykırılıklardan en az etkilenme
- Model oluşturma sırasında örneklerin sırasından etkilenmeme
- Çok boyutlu veriler üzerinde çalışma
- Kullanıcıların kısıtlarını göz önünde bulundurma
- Sonucun yorumlanabilir ve anlaşılabilir olması

İyi Kümeleme

- İyi kümeleme yöntemiyle elde edilen kümelerin özellikleri
 - aynı küme içindeki nesnelere arası benzerlik fazla
 - farklı kümelere bulunan nesnelere arası benzerlik az
- Oluşan kümelerin kalitesi seçilen benzerlik ölçütüne ve bu ölçütün gerçekleşmesine bağlı
 - Uzaklık / Benzerlik nesnelere nitelik tipine göre değişir
 - Nesnelere arası benzerlik: $s(i,j)$
 - Nesnelere arası uzaklık: $d(i,j) = 1 - s(i,j)$
- İyi bir kümeleme yöntemi veri içinde gizlenmiş örüntüleri bulabilmeli
- Veriyi gruplama için uygun kümeleme kriteri bulunmalı
 - kümeleme = aynı kümedeki nesnelere arası benzerliği en büyüten, farklı kümedeki nesnelere arası benzerliği en küçülten fonksiyon
- Kümeleme sonucunun kalitesi seçilen kümelerin şekline ve temsil edilme yöntemine bağlı

Kümeleme Yöntemlerinde Kullanılan Uzaklıklar

- Öklid

$$d(i, j) = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2}$$

- Minkowski

$$d(i, j) = \left[\sum_{k=1}^p (|x_{ik} - x_{jk}|^m) \right]^{\frac{1}{m}}$$

- Manhattan

$$d(i, j) = \sum_{k=1}^p (|x_{ik} - x_{jk}|)$$

Kümeleme Yöntemleri

- Hiyerarşik Kümeleme
 - Birleştirici Hiyerarşik Yöntemler
 - En yakın komşu algoritması
 - En uzak komşu algoritması
- Hiyerarşik Olmayan Kümeleme
 - K-Ortalamalar Yöntemi (K-Means)

En yakın komşu algoritması

- En yakın komşu yöntemine «tek bağlantı kümeleme yöntemi» adı da verilmektedir. Başlangıçta tüm gözlem değerleri birer küme olarak değerlendirilir. Adım adım bu kümeler birleştirilerek yeni kümeler elde edilir.
- Bu yöntemde öncelikle gözlemler arasındaki uzaklıklar belirlenir. Öklid uzaklık bağıntısı kullanılabilir.
- Uzaklıklar göz önüne $\text{Min } d(i,j)$ seçilir. Söz konusu uzaklıkla ilgili satırlar birleştirilerek yeni bir küme elde edilir. Bu duruma göre uzaklıkların yeniden hesaplanması gerekir.
- Tek bir gözlemden oluşan kümeler arasındaki uzaklıkları doğrudan hesaplayabiliriz. Ancak birden fazla gözlem değerine sahip olan iki küme arasındaki uzaklığın belirlenmesi gerektiğinde farklı bir yol izlenir. İki kümenin içerdiği gözlemler arasında «birbirine en yakın olanların uzaklığı» iki kümenin birbirine olan uzaklığı olarak kabul edilir.

Örnek 1.

- Aşağıdaki tabloda verilen beş gözlem değeri, en yakın komşu algoritması ile kümelenmek isteniyor.

Gözlemler	X₁	X₂
1	4	2
2	6	4
3	5	1
4	10	6
5	11	8

- Adım1. Öncelikle uzaklık tablosu oluşturulur. Her bir gözlemin birbiriyle arasındaki öklid uzaklığı hesaplanır.

Örnek 1.

$$d(1,2) = \sqrt{(4 - 6)^2 + (2 - 4)^2} = 2,83$$

$$d(1,3) = \sqrt{(4 - 5)^2 + (2 - 1)^2} = 1,41$$

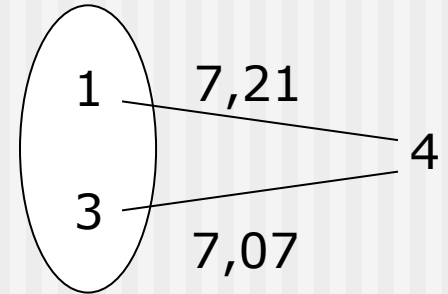
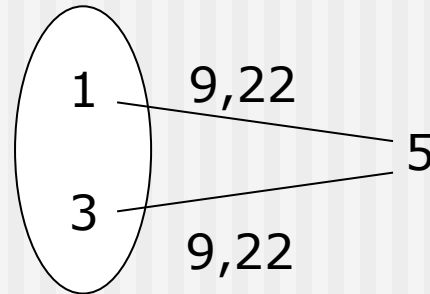
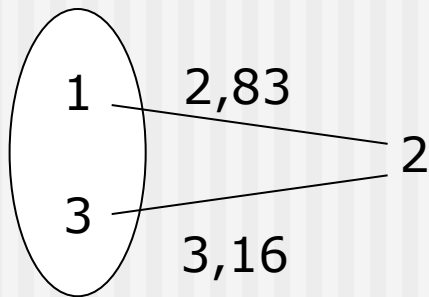
$$d(1,4) = \sqrt{(4 - 10)^2 + (2 - 6)^2} = 7,21$$

...

Gözlemler	1	2	3	4	5
1					
2	2,83				
3	1,41	3,16			
4	7,21	4,47	7,07		
5	9,22	6,4	9,22	2,24	

Örnek 1.

- Adım 2. Uzaklıklar tablosunda $\text{Min } d(i,j)$ değerinin 1,41 olduğu görülmektedir. İlgili gözlemler 1 ve 3 gözlemleridir. Bu iki değer birleştirilerek (1,3) kümesi elde edilir. Sonrasında bu kümeye göre uzaklıklar matrisi yeniden incelenir.



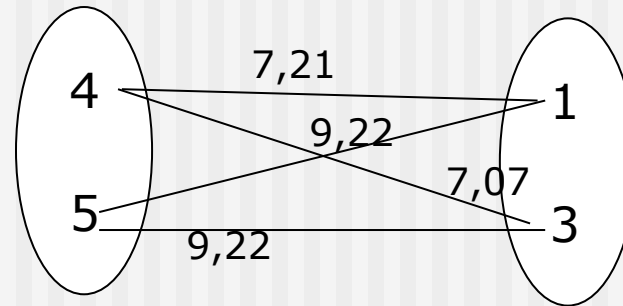
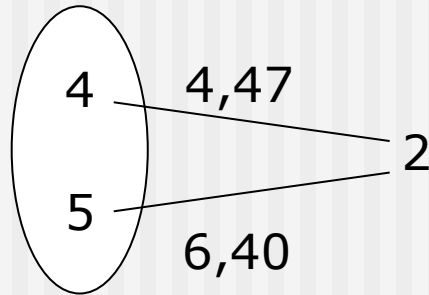
Örnek 1.

- Yeni uzaklık tablosu,

Gözlemler	(1,3)	2	4	5
(1,3)				
2	2,83			
4	7,07	4,47		
5	9,22	6,4	2,24	

- Bu tabloya bakıldığında $\text{Min } d(i,j)=2,24$ olduğu görülür. Bu değer 4 ve 5 gözlemleri arasındaki uzaklığı görülür. (4,5) yeni bir küme oluşturur. Bu durumda (1,3), 2 ve (4,5) kümeleri arasındaki uzaklık tablosu yeniden oluşturulur.

Örnek 1.



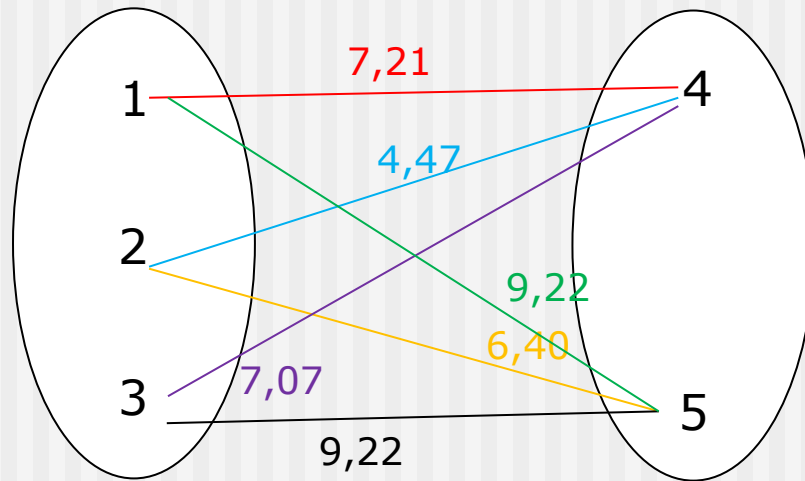
Örnek 1.

- Bu durumdaki uzaklık tablosu,

Gözlemler	(1,3)	2	(4,5)
(1,3)			
2	2,83		
(4,5)	7,07	4,47	

- Adım 4. En son uzaklıklar tablosu incelendiğinde $\text{Min } d(i,j)=2,83$ olduğu görülür. O halde bu uzaklık ile ilgili olan 2 gözlemi ile (1,3) kümesi birleştirilecektir. Elde edilen (1,2,3) kümesi ile (4,5) kümesi arasındaki uzaklığı belirlemek için kümeler içindeki her bir değer eşlenir ve en küçük olan belirlenir. En küçük uzaklık 4,47 olduğuna göre iki küme arasındaki uzaklığın bu değer olduğu kabul edilir.

Örnek 1.



- Adım 5. Elde edilen iki küme birleştirilerek sonuç küme bulunur. Bu küme (1,2,3,4,5) gözlemlerinden oluşan kümedir. Uzaklık düzeyi göz

önüne alınarak kümeler şu şekilde belirlenir.

Uzaklık	Kümeler
1,41	(1,3)
2,24	(4,5)
2,83	(1,2,3)
4,47	(1,2,3,4,5)

En uzak komşu algoritması

- En yakın komşu algoritması ile benzer adımları içerir. Gözlemler arasındaki uzaklıklar hesaplanır ve minimum değerli olan birleştirilir. Sonraki küme uzaklıkları tablosu oluşturulurken en uzak mesafe kullanılır.

K-Ortalamlar Yöntemi (K-Means) (1/2)

- Bu yöntemde daha başlangıçta belli sayıdaki küme için toplam ortalama hatayı minimize etmek amaçlanır.
- N boyutlu uzayda N örnekli kümelerin verildiğini varsayalım. Bu uzay $\{C_1, C_2, \dots, C_k\}$ biçimde K kümeye ayrılsın. O zaman $\sum n_k = N$ ($k=1, 2, \dots, k$) olmak üzere C_k kümesinin ortalama vektörü M_k şu şekilde hesaplanır.

$$M_k = \frac{1}{n_k} \sum_{i=1}^{n_k} X_{ik}$$

- Burada X_k değeri C_k kümesine ait olan i . örnektir. C_k kümesi için kare-hata, her bir C_k örneği ile onun merkezi (centroid) arasındaki Öklid uzaklıkları toplamıdır. Bu hataya «küme içi değişme» adı da verilir.

K-Ortalamlar Yöntemi (K-Means)

(2/2)

- Küme içi değişmeler şu şekilde hesaplanır.

$$e_i^2 = \sum_{i=1}^{n_k} (X_{ik} - M_k)^2$$

- K kümesini içeren bütün kümeler uzayı için kare-hata içindeki değişmelerin toplamıdır. O halde söz konusu kare-hata şu şekilde hesaplanır.

$$E_k^2 = \sum_{k=1}^K e_k^2$$

- Kare-hata kümeleme yönteminin amacı verilen K değeri için E_k^2 değerini minimize eden K kümelerini bulmaktır. O halde k-ortalamlar algoritmasında E_k^2 değerinin bir önceki iterasyona göre azalması beklenir.

K-Means Algoritmasının Adımları

- K-Means algoritmasına başlamadan önce k küme sayısının belirlenmesi gerekir. Sonra aşağıdaki işlemler gerçekleştirilir.
 1. Her bir kümenin merkezi belirlenir. Bu merkezler M_1, M_2, \dots, M_k biçimindedir.
 2. e_1, e_2, \dots, e_k küme içi değişmeler hesaplanır. Bu değişmelerin toplamı olan E_k^2 değeri bulunur.
 3. M_k merkez değerleri ile gözlem değerleri arasındaki uzaklıklar hesaplanır. Bir gözlem değeri hangi yakın ise o merkez ile ilgili küme içine dahil edilir.
 4. Yukarıdaki 2. ve 3. adımlar kümelerde değişiklik olmayıncaya kadar devam ettirilir.

K-Means Algoritmasının Özellikleri

- Gerçeklemesi kolay
- Karmaşıklığı diğer kümeleme yöntemlerine göre az
- K-Means algoritması bazı durumlarda iyi sonuç vermeyebilir
 - Veri grupları farklı boyutlarda ise
 - Veri gruplarının yoğunlukları farklı ise
 - Veri gruplarının şekli küresel değilse
 - Veri içinde aykırılıklar varsa

Örnek 2.

- Aşağıdaki gözlem değerleri k-ortalamlar yöntemi ile kümelenmek isteniyor.

Gözlemler	Değişken1	Değişken2
X_1	4	2
X_2	6	4
X_3	5	1
X_4	10	6
X_5	11	8

- Kümelerin sayısı başlangıçta $k=2$ kabul edilir. Rasgele iki küme belirlenir.

$$C_1 = \{X_1, X_2, X_4\}$$

$$C_2 = \{X_3, X_5\}$$

Örnek 2.

Gözlemler	Değişken1	Değişken2	Küme Üyeliği
X ₁	4	2	C ₁
X ₂	6	4	C ₁
X ₃	5	1	C ₂
X ₄	10	6	C ₁
X ₅	11	8	C ₂

- Adım 1. a) Belirtilen iki kümenin merkezleri şu şekilde hesaplanır.

$$M_1 = \left\{ \frac{4 + 6 + 10}{3}, \frac{2 + 4 + 6}{3} \right\} = \{6.67, 4.0\}$$

$$M_2 = \left\{ \frac{5 + 11}{2}, \frac{1 + 8}{2} \right\} = \{8.0, 4.5\}$$

Örnek 2.

- b) Küme içi değişmeler şu şekilde hesaplanır.

$$e_1^2 = [(4 - 6,67)^2 + (2 - 4,0)^2] + [(6 - 6,67)^2 + (4 - 4,0)^2] + [(10 - 6,67)^2 + (6 - 4,0)^2] = 26,67$$

$$e_2^2 = [(5 - 8)^2 + (1 - 4,5)^2] + [(11 - 8)^2 + (8 - 4,5)^2] = 42,50$$

- Bu durumda toplam kare-hata şu şekilde hesaplanır.

$$E^2 = e_1^2 + e_2^2 = 26,67 + 42,50 = 69,17$$

Örnek 2.

- C) M_1 ve M_2 merkezlerinden olan uzaklıkların minimum olması istendiğinden aşağıdaki hesaplamalar yapılır. Öklid uzaklık formülü kullanılarak söz konusu mesafeler hesaplanır. Örneğin (M_1, X_1) noktaları arasındaki uzaklık $M_1 = \{6.67, 4.00\}$ ve $X_1 = \{4, 2\}$ olduğuna göre şu şekilde hesaplanır.

$$d(M_1, X_1) = \sqrt{(6,67 - 4)^2 + (4 - 2)^2} = 3,33$$

$$d(M_2, X_1) = \sqrt{(8 - 4)^2 + (4,5 - 2)^2} = 4,72$$

- Bu işlemler sonucunda X_1 gözlem değerinin M_1 ve M_2 merkezlerine olan uzaklıkları göz önüne alındığında $d(M_1, X_1) < d(M_2, X_1)$ olduğu görülür. Bu durumda M_1 merkezinin X_1 gözlem değerine daha yakın olduğu anlaşılır. O halde $X_1 \in C_1$ olarak kabul edilir. Benzer biçimde tüm gözlem değerleri için tablo oluşturulur.

Örnek 2.

Gözlemler	M_1 'den uzaklık	M_2 'den uzaklık	Küme Üyeliği
X_1	$d(M_1, X_1) = 3,33$	$d(M_2, X_1) = 4,72$	C_1
X_2	$d(M_1, X_2) = 0,67$	$d(M_2, X_2) = 2,06$	C_1
X_3	$d(M_1, X_3) = 3,43$	$d(M_2, X_3) = 4,61$	C_1
X_4	$d(M_1, X_4) = 3,89$	$d(M_2, X_4) = 2,50$	C_2
X_5	$d(M_1, X_4) = 5,90$	$d(M_2, X_4) = 4,61$	C_2

Örnek 2.

- Bu durumda yeni kümeler şu şekilde olacaktır.

$$C_1 = \{X_1, X_2, X_3\}$$

$$C_2 = \{X_4, X_5\}$$

- Adım 2. Yukarıda belirtilen iki kümenin merkezleri şu şekilde hesaplanır.

$$M_1 = \left\{ \frac{4 + 6 + 5}{3}, \frac{2 + 4 + 1}{3} \right\} = \{5, 2.33\}$$

$$M_2 = \left\{ \frac{10 + 11}{2}, \frac{6 + 8}{2} \right\} = \{10.5, 7\}$$

Örnek 2.

- b) Küme içi değişmeler şu şekilde hesaplanır.

$$e_1^2 = [(4 - 5)^2 + (2 - 2,33)^2] + [(6 - 5)^2 + (4 - 2,33)^2] \\ + [(5 - 5)^2 + (1 - 2,33)^2] = 9,33$$

$$e_2^2 = [(10 - 10,5)^2 + (6 - 7)^2] + [(11 - 10,5)^2 + (8 - 7)^2] = 2,50$$

- Bu durumda toplam kare-hata şu şekilde hesaplanır.

$$E^2 = e_1^2 + e_2^2 = 9,33 + 2,50 = 11,83$$

- Bu değer bir önceki iterasyonda elde edilen $E^2 = 69,17$ değerinden daha küçük olduğu anlaşılmaktadır.

Örnek 2.

- M_1 ve M_2 merkezlerinden gözlem değerlerine olan uzaklıklar hesaplanır. Bunun sonucunda $d(M_1, X_1) < d(M_2, X_1)$ olduğu görülür. Bu durumda M_1 merkezinin X_1 gözlem değerine daha yakın olduğu anlaşılır. O halde $X_1 \in C_1$ olarak kabul edilir. Benzer biçimde tüm gözlem değerleri için tablo oluşturulur.

Gözlemler	M_1 'den uzaklık	M_2 'den uzaklık	Küme Üyeliği
X_1	$d(M_1, X_1) = 1,05$	$d(M_2, X_1) = 8,20$	C_1
X_2	$d(M_1, X_2) = 1,94$	$d(M_2, X_2) = 5,41$	C_1
X_3	$d(M_1, X_3) = 1,33$	$d(M_2, X_3) = 8,14$	C_1
X_4	$d(M_1, X_4) = 6,20$	$d(M_2, X_4) = 1,12$	C_2
X_5	$d(M_1, X_4) = 8,25$	$d(M_2, X_4) = 1,12$	C_2

Örnek 2.

- Bu durumda yeni kümeler şu şekilde oluşacaktır.

$$C_1 = \{X_1, X_2, X_3\}$$

$$C_2 = \{X_4, X_5\}$$

- Kümelerde önceki adıma göre herhangi bir değişme olmadığı için iterasyona son verilir.