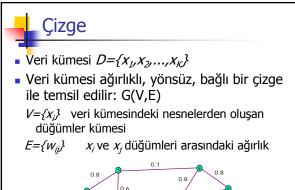
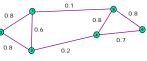


# VERİ MADENCİLİĞİ Farklı Demetleme Yöntemleri

Yrd. Doç. Dr. Şule Gündüz Öğüdücü www.cs.itu.edu.tr/~gunduz/courses/verimaden/





www.cs.itu.edu.tr/~gunduz/courses/verimaden/



### Çizge Tabanlı Demetleme

- *S={V,N,W,C}* 
  - V: veri kümesindeki nesnelerden oluşan düğümler kümesi
  - $N \subseteq V \times V$
  - W: N kümesinin elemanları için simetrik benzerlik matrisi
  - P. Demetleme kriteri
- Çizge Bölme: P demetleme kriterini enbüyütecek şekilde V kümesini k demede bölmek  $(C=\{C_1,\ldots,C_k\})$ .
  - Her demet bir altçizge G<sub>i</sub>(V<sub>i</sub>,E<sub>i</sub>)

$$\bigcup_{i=1}^{k} V_{i} = V$$

$$E_i = \{\{u, v\} \in E \land u, v \in V_i\}$$

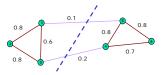
Problem: Çizge tabanlı demetleme yöntemleri için *P* demetleme

www.cs.itu.edu.tr/~gunduz/courses/verimaden/



### Çizge Tabanlı Demetleme Problemi

- Demetlemenin sağlaması gereken koşullar:
  - Aynı demetlerdeki nesnelerin birbirine daha çok banzemesi Farklı demetlerdeki nesneler birbirine daha az benzemesi
- Aynı koşullar çizge tabanlı demetlemeye uygulanırsa



- 1. Aynı grup içindeki ağırlıkları enbüyütme
- 2. Farklı gruplar arasındaki ağırlıkları enküçültme

www.cs.itu.edu.tr/~gunduz/courses/verimaden/



#### Çizge Tabanlı Demetleme için Tanımlar

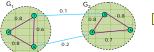
- Tanımlar:
  - uzaklık d, benzerlik s
    - *d=1-s*
  - C<sub>i</sub> ve C<sub>i</sub> demetleri arasındaki uzaklık: d(C<sub>i</sub>, C<sub>i</sub>)
    - tek bağ, tam bağ ya da ortalama
  - C<sub>i</sub> demedinin çapı: diam(C<sub>i</sub>)
    - C<sub>i</sub> demedinde bulunan en uzak iki nesne arasındaki
    - C<sub>i</sub> demedinden bulunan tüm nesneler arasındaki uzaklıkların ortalaması

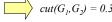


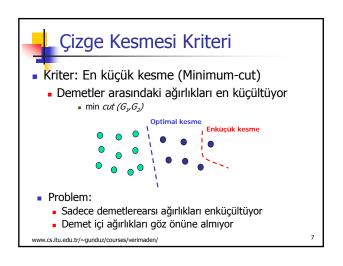
## Çizge Kesmesi

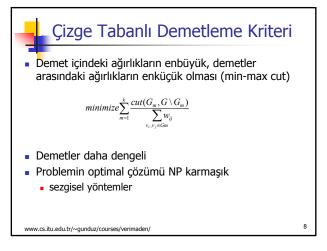
 Çizge Kesmesi: Demetleri biribirine bağlayan ayrıtların ağırlıklarının toplamı

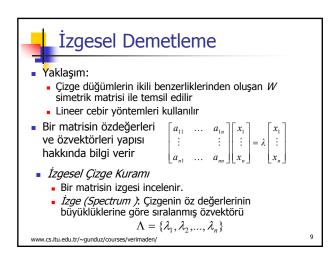
$$cut(G_{1},G_{2}) = \sum_{x_{i} \in G_{1}, x_{j} \in G_{2}} w_{ij}$$

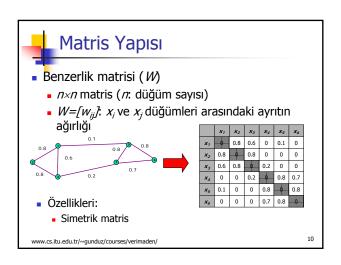


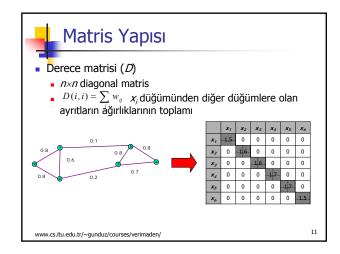


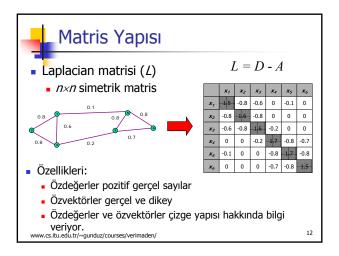














#### Optimal Enküçük Kesme Bulma (Hall'70, Fiedler'73)

- İki altçizgeye (G<sub>1</sub>, G<sub>2</sub>) bölünen çizge bir vektörle temsil  $p_i = \begin{cases} +1 & \text{if } x_i \in G_I \end{cases}$ -1 if  $x_i \in G_2$
- Bölmenin kesmesini enküçültmek için f(p) fonksiyonunu enküçültecek p vektörü bulunur:

$$f(p) = \sum_{i,j \in V} w_{ij} (p_i - p_j)^2 = p^T L p$$
• Rayleigh Kuramına göre:

- - f(p)'nin enküçük değeri L matrisinin ikinci enküçük özdeğeri ile elde edilir.
  - p için optimal çözüm Fiedler vektörü olarak bilinen  $\lambda_2$ vektörüdür.

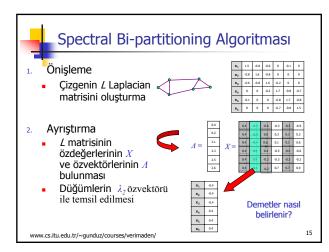
www.cs.itu.edu.tr/~gunduz/courses/verimaden/



## İzgesel Çizge Demetleme

- A. Pothen, H.D. Simon and K. Paul Liou. *Partitioning* Sparse Matrices with Eigenvectors of Graphs, SIAM J. Mat. Theory and Appl., Vol. 11, No. 3, pp. 430 - 452, 1990.
  - Önişleme
    - veri kümesinin matris olarak temsil edilmesi
  - Avristirma
    - Matrisin özvektörlerinin ve özdeğerlerinin bulunması
    - Veri kümesindeki her nesnenin bir veya daha çok özvektörü kullanılarak daha küçük bir boyuta taşınması
  - Gruplama
    - Yeni boyutlardan yararlanarak nesnelerin iki veya daha fazla demedé ayrılması

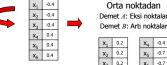
www.cs.itu.edu.tr/~gunduz/courses/verimaden/





### Spectral Bi-partitioning Algoritmasi

- Gruplama
  - Tek boyutlu vektörde bulunan elemanlar sıralanır
  - Vektör ikiye bölünür
- Bölme noktası nasıl belirlenir?
  - Ortalamadan ya da orta noktadan bölünür





www.cs.itu.edu.tr/~gunduz/courses/verimaden/



# K-Yönlü İzgesel Demetleme

- Biz çizge k adet altçizgeye bölünmek isteniyor.
- İki vaklasım
  - Yinelemeli ikiye demetleme (L. Hagen, A.B. Kahng, New spectral methods for ratio cut partitioning and clustering, IEEE Trans. Comput. Aided Des. 11,1992)
    - Yinelemeli olarak ikiye demetleme algoritmasını hiyerarşıik olarak uygulanması
  - Daha fazla sayıda özvektörü kullanarak demetleme (J. Shi and J. Malik. Normalized cuts and image segmentation. IEEE Trans. on Pattern Analysis and Machine Intelligence, 22(8):888-905, 2000.=
    - Özvektörleri kullanarak veriyi daha az boyutlu bir uzaya taşır

www.cs.itu.edu.tr/~gunduz/courses/verimaden/



## K-Yönlü İzgesel Demetleme

- Çizgeler arasındaki optimal kesmeyi yaklaşık olarak bulabilir (Shi & Malik, 2002).
- Veri içindeki grupları belirgin hale getirir (M. Brand and K. Huang. A unifying theorem for spectral embedding and clustering, Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics, January 2003.)
  - Benzer nesneler arasındaki ilişki kuvvetleniyor, daha az benzer nesneler arasındaki ilişki zayıflıyor.



### K-Özvektör Demetleme

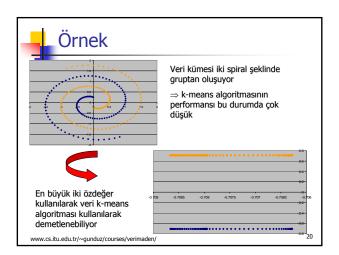
- k özvektör kullanarak demetleme yapılıyor (A. Ng, M. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm, In Advances in Neural Information Processing Systems 14: Proceedings of the 2001.)
  - Önişleme: ölçeklendirilmiş ağırlık matrisi oluşturulur

$$A' = D^{-1/2} A D^{-1/2}$$

- Ayrıştırma: A' matrisinin özvektörleri bulunur. Veri kümesi en büyük k özdeğer ile temsil edilir
- Demetleme: k-means algoritması kullanılarak nxk boyutundaki veri k demede ayrılır.

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

19





# Çizge Tabanlı Demetleme Yöntemi (Kannan'00)

- En küçük kesme bulunarak demetlenirse altçizgeler arasındaki nesneler dengeli dağılmayabilir.
- Demetlerin kalitesi önemli (Ravi Kannan and Santosh Vempala and Adrian Vetta, On Clusterings: Good, Bad, and Spectral, Proceedings of the 41st Annual Symposium on the Foundation of Computer Science, 2000.)





www.cs.itu.edu.tr/~gunduz/courses/verimaden/



# Çizge Tabanlı Demetleme Yöntemi (Kannan'00)

Bir kesme (S, S) için genişlik

$$\psi(S) = \frac{\sum_{x_i \in S, x_j \in \overline{S}} w_{ij}}{\min\{|S| | \overline{S}|\}}$$

Bir kesme (S, S̄) için iletkenlik

$$\phi(S) = \frac{\sum_{x_i \in S, x_j \in \overline{S}} w_{ij}}{\min\{c(S), c(\overline{S})\}}$$

c(S) şu şekilde tanımlanmıştır

$$c(S) = c(S, V) = \sum_{x, \in S} \sum_{x, \in V} w_{ij}$$

- Bir demedin genişliği (iletkenliği) demet içindeki kesmelerin genişliklerinin (iletkenliklerinin) en küçüğü
- Demetlemenin genişliği (İletkenliği) demetlerin genişliklerinin (İletkenliklerinin) en küçüğü
- Genişliğin (iletkenliğin) büyük olması iyi bir demetleme olduğunu gösteriyor

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

22



# Çizge Tabanlı Demetleme Yöntemi (Kannan'00)

- Çizgeyi demetlemek için iki kriter beraber kullanılıyor:
  - Her demedin iletkenliği (genişliği) en az α değerinde olmalı
  - Demetler arası ayrıtların ağırlıklarının toplamının bütün ayrıtların ağırlıklarının toplamına oranı ε değerinden büyük olmamalı
- Problemin çözümü NP-karmaşık olduğu için yaklaşık bir çözüm önerilyor.

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

23

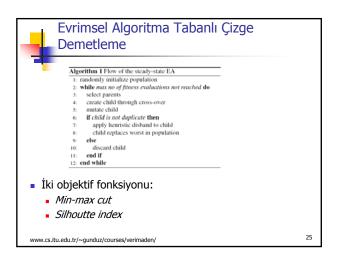


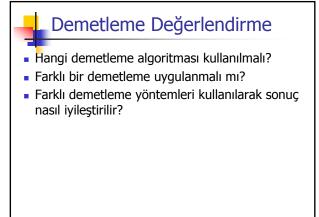
#### Evrimsel Algoritma Tabanlı Çizge Demetleme

- Çizge demetleme problemi NP-karmaşık bir problem olduğundan doğa esinli algoritmalar kullanılarak problem çözülebilir (Ş.Uyar and Ş.Oguducu, A New Graph-Based Evolutionary Approach to Sequence Clustering, The Fourth International Conference on Machine Learning and Applications, 2005)
- Amaç:
  - Aynı demetteki nesneler arasındaki ayrıtların ağırlıklarının toplamının, demetler arasındaki ayrıtların ağırlıklarının toplamına oranını enbüyütmek
  - Demet sayısını adaptif olarak belirlemek.

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

24

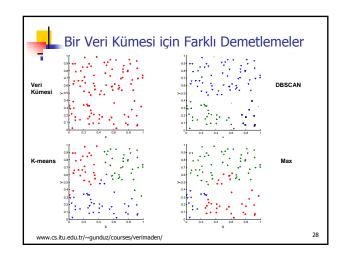




www.cs.itu.edu.tr/~gunduz/courses/verimaden/

Pemetleme Değerlendirme

Farklı demetleme teknikleri
En iyi demetleme algoritmasını seçmek zor
Verinin istatistiksel özelliğine, önişleme tekniklerine, nitelik sayısına bağlı olarak algoritmaların avantajları ve dezavantajları var
Aynı veri kümesi üzerinde farklı algoritmalar farklı demetleme sonuçları üretebilir. Hangi demetlemenin daha iyi olduğuna karar vermek gerekiyor
uygulama alanını iyi incelemek gerekiyor
demetleme sonucunu iyi anlamak gerekiyor



Gözetimli öğrenme için kullanılan yöntemler:
 Doğruluk, kesinlik, duyarlılık
 Demetleme yöntemlerinde değerlendirilmesi gerekenler:
 Hatalı veriler için örüntü bulunmaması
 Farklı demetleme algoritmalarını karşılaştırma
 Farklı demetlemeleri karşılaştırma

Demetleme Değerlendirme

Farklı demetleri karşılaştırma

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

Veri Kümesi Demetlemeye Uygun mu?

• Veri kümesi içinde gruplar olmayabilir.

• Nesneler rasgele dağılmış

• Her demetleme algoritması veri kümesi üzerinde demetleme yapar

• Hopkins istatistiği: Veri kümesi içinde demetler bulunup bulunmadığını test etmek için kullanılır

• Veri uzayında rasgele dağılmış p nokta üretilir

• Veri kümesinden örnekleme ile p nokta seçilir

• Her iki küme için veri kümesinden en yakın nesnler seçilir.

• u, yapay olarak üretilmiş noktalara olan uzaklık, wı, veri kümesinden seçilmiş olan noktalara olan uzaklık

• Hopkins istatistiği  $H = \sum_{i=1}^{p} w_i$   $H = \sum_{i=1}^{p} w_i + \sum_{i=1}^{p} w_i$ \*\*www.cs.ltu.edu.tt/~\*gunduz/courses/verimaden/\*



### Demetleme Değerlendirme Ölçütleri

- Üç yaklaşım:
  - Harici Gösterge: Veri kümesi için öngörülen bir yapıya dayanarak değerlendirme
  - Dahili Gösterge: Ek bir bilgi kullanmadan veri kümesinden elde edilen bilgiye dayanarak değerlendirme
  - Göreceli Değerlendirme: Aynı algoritmanın farklı parametrelerini kullanarak elde edilen demetleme sonuçlarını değerlendirme
- İki kriter:
  - Sıkılık: Her demette bulunan nesneler birbirine mümkün olduğunca yakın olmalı
  - Uzaklık: Demetler birbirinden mümkün olduğunca uzak olmalı
    - Tek bağ
    - Tam bač
    - Demet merkezleri arasındaki uzaklık

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

31



## Harici Değerlendirme

- Demetleme algoritması kullanılarak elde edilen demetleme C={C<sub>1</sub>,...,C<sub>k</sub>}
- Veri içinden bulunan gruplar P={P<sub>1</sub>,...,P<sub>m</sub>}
- Demetleme sonucundan elde edilen dağılım
  - SS: Eğer iki nesne Ciçin aynı demette ve Piçin aynı grupta ise (a)
  - SD: Eğer iki nesne Ciçin aynı demette ancak Piçin farklı gruplarda ise (b)
  - DS: Eğer iki nesne Ciçin farklı demette ancak Piçin aynı grupta ise (c)
  - DS: Eğer iki nesne  $\mathcal C$ için farklı demette ancak  $\mathcal P$ için aynı grupta ise (d)

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

2



## Harici Değerlendirme Ölçütleri

- Rand Statistics:  $R = \frac{a+d}{d}$
- Jaccard katsayısı:  $J = \frac{a}{a+b+c}$
- Folkes ve Mallows göstergesi:  $FM = \sqrt{\frac{a}{a+b} \cdot \frac{a}{a+c}}$

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

33



## Harici Değerlendirme Ölçütleri

Entropi: Her demette sınıfların nasıl dağıldığı

$$e_{i} = -\sum_{j=1}^{m} p_{ij} \log_{2} p_{ij}$$

$$e = \sum_{i=1}^{k} \frac{n_{i}}{n} e_{i}$$

m: sınıf sayısı k: demet sayısı p<sub>ii</sub>= n<sub>ii</sub>/n<sub>i</sub>

 $p_{ij} = n_{ij}/n_i$   $n_i$ : i demedindeki nesne sayısı  $n_{ij}$ : i demedinde j sınfından nesne sayısı n: toplam nesne sayısı

 Saflık: Bir demette ne kadar tek sınıftan örnek bulunduğu

$$p_i = \max p_i$$

$$purity = \sum_{i=1}^{k} \frac{n_i}{n_i} p$$

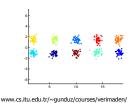
www.cs.itu.edu.tr/~gunduz/courses/verimaden/

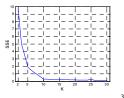
3



## Dahili Değerlendirme Ölçütleri

- Sadece veri kümesi özellikleri kullanılarak yapılan değerlendirme
  - Hataların karelerinin toplamı (SSE)
- İki farklı demetlemeyi ya da iki demedi karşılaştırmak için iyi bir yöntem
- Demet sayısını tahmin etmek için de kullanılabilir.





# Dahili Değerlendirme Ölçütleri

- Silhouette Göstergesi:
  - $x_i$  nesnesi  $C_i$  demedinde
  - Ortalama uzaklığa göre  $x_i$  nesnesine en yakın demet  $C_h$
  - *x<sub>i</sub>* nesnesi için silhouette göstergesi

$$s(x_i) = \frac{d(x_i, C_h) - d(x_i, C_j)}{\max(d(x_i, C_h), d(x_i, C_j))}$$

- $-1 \leq s(v_i) \leq 1$
- 1'e yakın olursa x<sub>i</sub> doğru demette

• Demetleme için silhouette göstergesi:

- $S_{j} = \frac{\sum_{i=1}^{j} S(i t_{i})}{|V_{j}|}$
- $GS = \frac{\sum_{j=1}^{k} S_j}{I}$

www.cs.itu.edu.tr/~gunduz/courses/verimaden/



# Göreceli Değerlendirme

- $P_{alg}$  seçilen demetleme algoritmasının parametreleri
- $P_{ak}$ deki parametrelerin farklı değerleri ile elde edilen demetlemeler  $C_{ii}$  i=1,...,nc arasında veriye en çok uyanı seçme
- iki durum:
  - Demet sayısı nc  $P_{alg}$ 'deki parametereler arasında değil:
    - Palg'deki parametrelerin değerleri geniş bir aralıkta değiştirilerek demetleme algoritması çalıştırılır. nc. << W (nesne sayısı) sabit kaldığı en geniş aralık seçilir. Paramettre değerleri olarak bu aralığın orta noktası seçilir. Bu yöntemle demet sayısı da belirlenmiş olur.
  - Demet sayısı  $nc P_{alg}$ 'deki parametereler arasında:

    - En iyi demetleme, demetleme göstergesi q kullanarak seçilir.
       nc<sub>min</sub> ve nc<sub>max</sub> arasında değişen farklı demet sayıları için algoritma çalıştırılır.
       nc'nin her farklı değeri için algoritma diğer parametreleri değiştirerek r defa çalıştırılır.
       Her nc'için q'nun en büyük değeri seçilir ve nc'nin fonksiyonu olarak çizilir. Bu çizim kullanılarak nc değeri belirlenir.

www.cs.itu.edu.tr/~gunduz/courses/verimaden/



## Göreceli Değerlendirme Ölçütleri

- Birbiri ile örtüşmeyen demetler için tanımlanmış göstergeler:
  - Hubert istatistiği:

 $\Gamma = (1/M) \sum_{i=1}^{N-1} \sum_{j=1}^{N} P(i, j) \cdot Q(i, j)$ 

N= veri kümesindeki nesne sayısı M=N(N-1)/2 P: yakınlık matrisi

Q: (i,j) elemanı x<sub>i</sub> ve x<sub>j</sub> nesnelerinin bulundukları demetler arasındaki uzaklık olan matris

Dunn göstergesi:

$$Dnc = \min_{i=1,\dots,nc} \left\{ \min_{j=i+1,\dots,nc} \frac{d(c_i, c_j)}{\max_{k=1,\dots,nc} diam(c_k)} \right\}$$

 $d(c_i,c_j) = \min_{x \in ci, y \in cj} d(x,y)$  $diam(C) = \max_{x,y \in C} d(x,y)$ 

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

38



### Hiyerarşik Demetleme için Gösterge

- Hiyerarşik demetleme için 4 gösterge
  - Demetlerin standard sapmalarının karakökü (RMSSTD)
  - Semi-partial R-squared (SPR)
  - R-Squared (RS)
  - İki demet arası uzaklık (CD)

www.cs.itu.edu.tr/~gunduz/courses/verimaden/

39



## Örtüşen Demetleri Değerlendirme

- Örtüşen demetleme için  $U=[u_{ij}]$  matrisi:  $x_i$  nesnesinin  $c_i$ demedine dahil olma olasılığı
- Bölme katsayısı:

$$PC = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{nc} u_{ij}^{2}$$

[1/nc,1] arasında değişir.

www.cs.itu.edu.tr/~qunduz/courses/verimaden/