



Voronoi diyagramları:
 Her öğrenme örneğini
 çevreleyen dışbükey
 çokgenlerden oluşan
 karar yüzeyi

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

komşuya göre negatif

olarak sınıflandırılır



K-En Yakın Komşu Yöntemi

- Üzaklık-ağırlıklı k-en yakın komşu algoritması
 - Öğrenme kümesindeki örneklere (x_j) , sınıflandırılmak istenen örneğe (x_q) olan uzaklıklarına göre ağırlıklar verilmesi
 - yakın örneklerin ağırlığı daha fazla

 $w = \frac{1}{d(x_q, x_i)^2}$

- k-en yakın komşunun ortalaması alındığı için gürültülü veriden az etkileniyor
- İlgisiz nitelikler uzaklığı etkileyebilir
 - bu nitelikler uzaklık hesaplarken kullanılmayabilir

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

8



Konular

- Sınıflandırma yöntemleri
- Örnek tabanlı yöntemler
 - k-En Yakın Komşu Yöntemi
 - Genetik Algoritmalar
 - Karar Destek MakinalarıBulanık Küme Sınıflandırıcılar
 - Öngörü
 - Eğri Uydurma
- Model Değerlendirme
- Öğrenme, sınama, geçerleme kümelerini oluşturma
- Sınıflandırıcıları birleştirme

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Genetik Algoritmalar

- Optimizasyon amaçlı
- Bir başlangıç çözümü öneriyor, tekrarlanan her ara adımda daha iyi çözüm üretmeye çalışıyor.
- Doğal evrime ve en iyi olanın yaşamını sürdürmesine dayanıyor
- Çözümü birey olarak sunuyor.
- Birey: I=I₁,I₂,...,I_n I_i kullanılan alfabenin bir karakteri
- gen: I_i
- Toplum: Bireylerden oluşan küme

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

10



Genetik Algoritmalar

- Genetik Algoritmalar (GA) 5 parçadan oluşuyor:
 - Bireylerden oluşan bir başlangıç kümesi, P
 - Çaprazlama (Crossover): Bir anne babadan yeni bireyler üretmek için yapılan işlem
 - Mutasyon: Bir bireyi rastgele değiştirme
 - Uygunluk (fitness): En iyi bireyleri belirleme
 - Çaprazlama ve mutasyon tekniklerini uygulayan ve uygunluk fonksiyonuna göre toplum içindeki en iyi bireyleri seçen algoritma

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Çaprazlama Örnekleri



a) Single Crossover

a) Multiple Crossover

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

```
Genetik Algoritma

Input:

P //Initial Population
Output:

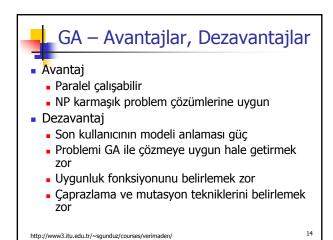
P' //Improved Population
Genetic Algorithm:

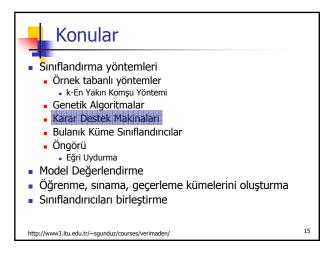
//Illustrates Genetic Algorithm

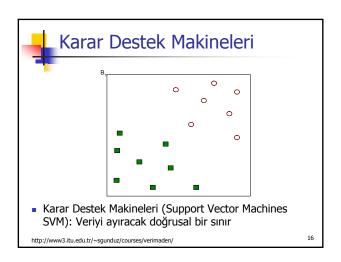
repeat

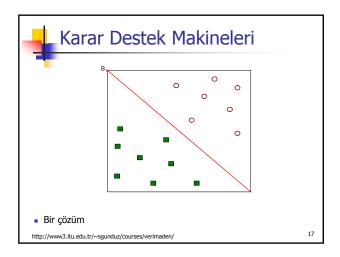
N = |P|;
P' = \emptyset;
repeat

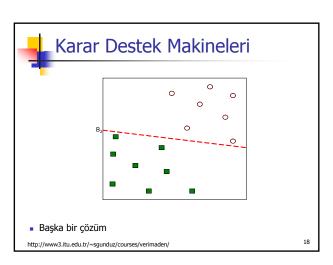
i_1, i_2 = \text{select}(P);
o_1, o_2 = \text{cross}(i_1, i_2);
o_1 = \text{mutate}(o_1);
o_2 = \text{mutate}(o_2);
P' = P' \cup \{o_1, o_2\};
\text{until } |P'| = N;
P = P';
until termination criteria satisfied;
```

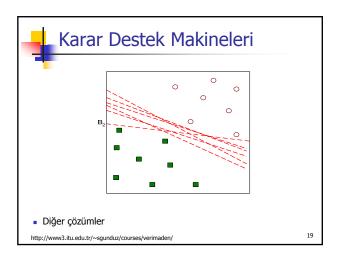


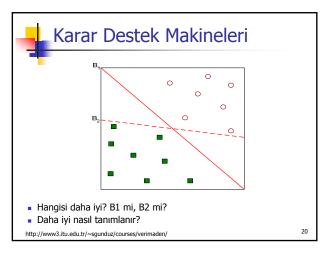


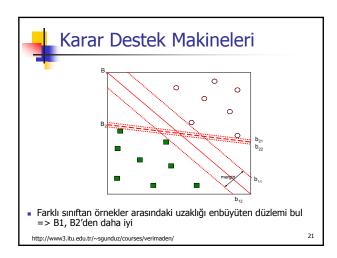


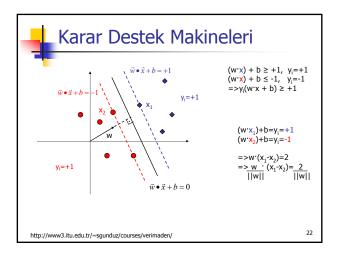


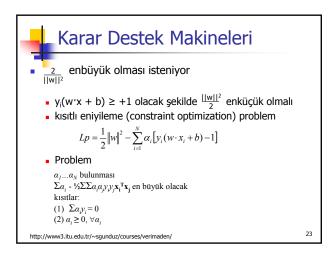


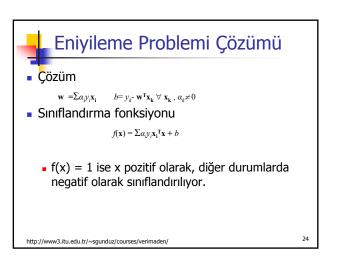


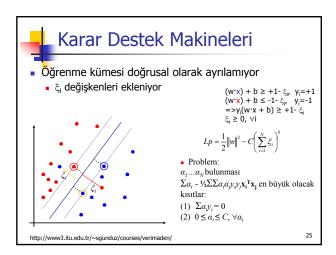


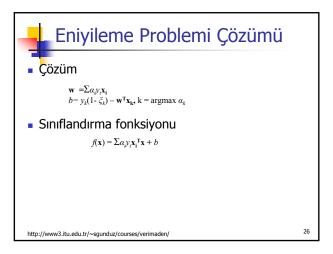




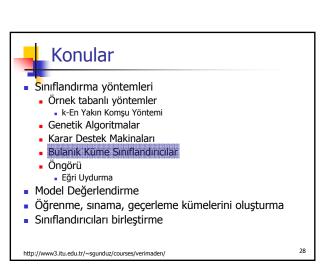












Bulanık Küme
Sınıflandırıcılar

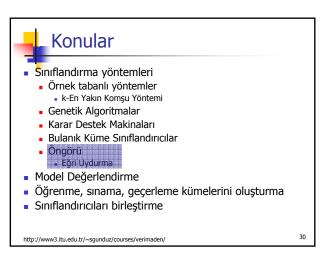
Bulanık mantık 0.0 ve 1.0 arasında gerçel değerler kullanarak üyelik dereceleri hesaplar

Nitelik değerleri bulanık değerlere dönüştürülür

Kurallar kümesi oluşturulur

Yeni bir örneği sınıflandırmak için birden fazla kural kullanılır

Her kuraldan gelen sonuç toplanır





Öngörü

- Sınıflandırma problemleriyle aynı yaklaşım
 - model oluştur
 - bilinmeyen değeri hesaplamak için modeli kullan
 - eğri uydurma
 - doğrusal
 - doğrusal olmayan
- Sınıflandırma ayrık değerli
- Öngörü sürekli değerli

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Eğri Uydurma

- Doğrusal eğri uydurma:
 - en basit eğri uydurma yöntemi
 - veri doğrusal bir eğri ile modellenir.
 - veri kümesindeki niteliklerin doğrusal fonksiyonu

$$y = w_0 + w_1 a_1 + w_2 a_2 + ... + w_k a_k$$

öğrenme kümesindeki y_i sınıfından bir x_i örneği için çıkış

$$y = w_0 x_{i0} + w_1 x_{i1} + w_2 x_{i2} + \dots + w_k x_{ik} = \sum_{i=1}^{k} w_i x_{ij}$$

karesel hatayı enküçültecek ağırlıkları bulma

$$\sum_{i=1}^{n} \left(y_i - \sum_{j=0}^{k} w_j x_{ij} \right)^2$$

32



Konular

- Sınıflandırma yöntemleri
- Model Değerlendirme
- Hata oranı
 - Anma
- Duyarlılık
- F-ölgütü ROC eğrileri
- Öğrenme, sınama, geçerleme kümelerini oluşturma
- Sınıflandırıcıları birleştirme

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



31

33

Sınıflandırma Modelini Değerlendirme

- Model başarımını değerlendirme ölçütleri nelerdir?
 - Hata oranı
 - Anma
 - Duyarlılık
 - F-ölçütü
- Farklı modellerin başarımı nasıl karşılaştırılır?
 - ROC

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Sınıflandırma Hatası



- Sınıflandırma yöntemlerinin hatalarını ölçme
 - başarı: örnek doğru sınıfa atandı
 - hata: örnek yanlış sınıfa atandı
 - hata oranı: hata sayısının toplam örnek sayısına bölünmesi
- Hata oranı sınama kümesi kullanılarak hesaplanır



Model Başarımını Değerlendirme

- Model başarımını değerlendirme ölçütleri
 - modelin ne kadar doğru sınıflandırma yaptığını ölçer hız, ölçeklenebilirlik gibi özellikleri değerlendirmez
- Karışıklık matrisi:

	ÖNGÖRÜLEN SINIF		
DOĞRU SINIF		Sınıf=1	Sınıf=-1
	Sınıf =1	a	b
	Sınıf =-1	С	d

a: TP (true positive) b: FN (false negative)

c: FP (false positive) d: TN (true negative)

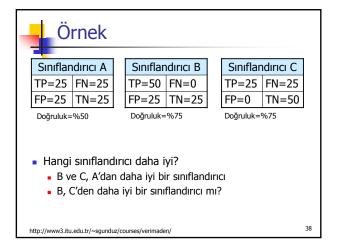


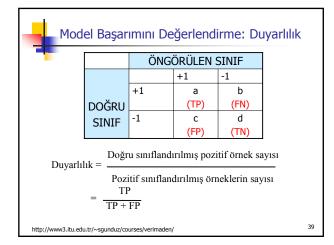
	ÖNGÖRÜLEN SINIF		
		+1	-1
	+1	a	b
DOĞRU		(TP)	(FN)
SINIF	-1	С	d
		(FP)	(TN)

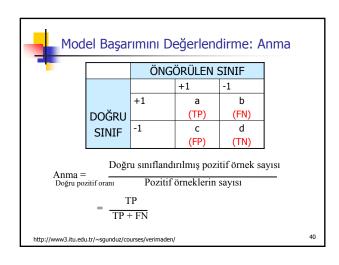
Modelin başarımı:

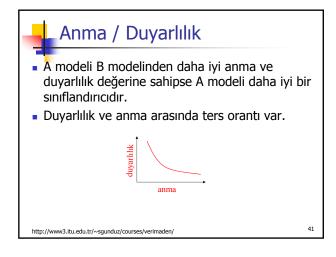
$$\begin{aligned} & Dogruluk = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{TP+TN+FP+FN} \\ & Hata & Orani = \frac{b+c}{a+b+c+d} = \frac{FN+FP}{TP+TN+FP+FN} \end{aligned}$$

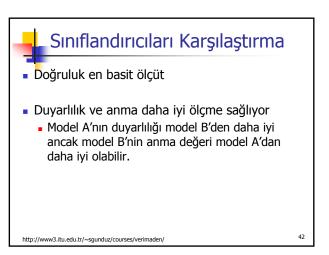
37

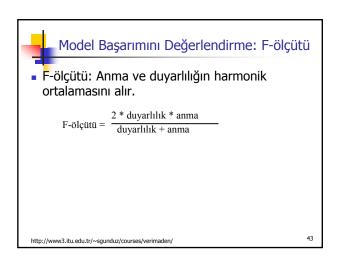








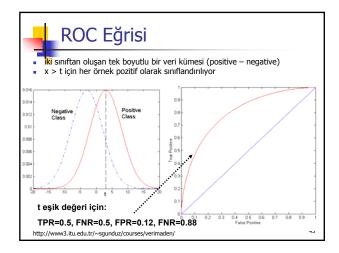


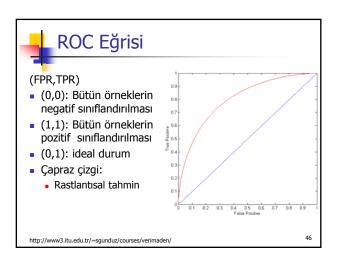


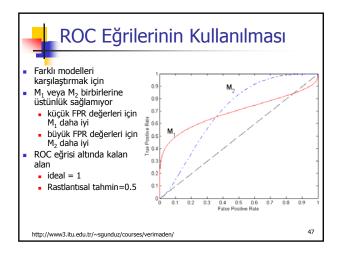
ROC (Receiver Operating Characteristic)

- İşaret işlemede bir sezicinin, gürültülü bir kanalda doğru algılama oranının yanlış alarm oranına karşı çizdirilen grafiği (algılayıcı işletim eğrisi)
- Farklı sınıflandırıcıları karşılaştırmak için ROC eğrileri
- Doğru pozitif (TPR y ekseni) oranının yanlış pozitif (FPR - x ekseni) oranına karşı çizdirilen grafiği
 - TPR = TP / (TP + FN)
 - FPR = FP / (TN + FP)
- ROC üzerindeki her nokta bir sınıflandırıcının oluşturduğu bir modele karşı düşer

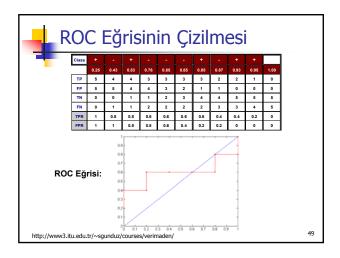
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/











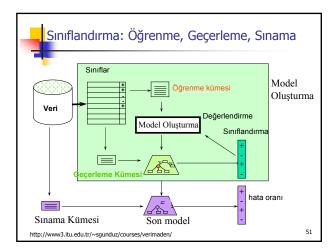


Model Parametrelerini Belirleme

- Sınama kümesi sınıflandırıcı oluşturmak için kullanılmaz
- Bazı sınıflandırıcılar modeli iki aşamada oluşturur
 - modeli olustur
 - parametreleri ayarla
- Sınama kümesi parametreleri ayarlamak için kullanılmaz
- Uygun yöntem üç veri kümesi kullanma: öğrenme, geçerleme, sınama
 - geçerleme kümesi parametre ayarlamaları için kullanılır
 - model oluşturulduktan sonra öğrenme ve geçerleme kümesi son modeli oluşturmak için kullanılabilir

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

50





Model Başarımını Tahmin Etme

- Örnek: Doğruluğu %25 olan bir modelin gerçek başarımı ne kadardır?
 - Sınama kümesinin büyüklüğüne bağlı
- Sınıflandırma (hileli) yazı tura atmaya benziyor
 - tura doğru sınıflandırma (başarı), yazı yanlış sınıflandırma (başarısızlık)
- İstatistikte birbirinden bağımsız olayların başarı ya da başarısızlıkla sonuçlanmaları Bernoulli dağılımı ile modellenir.
- Gerçek başarı oranını belirlemek için istatistikte güven aralıkları tanımlanmıştır.

http://www3.itu.edu.tr/~squnduz/courses/verimaden/

52



Güven Aralığı

- p belli bir güvenle belli bir aralıkta bulunmaktadır.
- Örnek: N=1000 olayda S=750 başarı sağlanmış.
 - Tahmin edilen başarı oranı: 75%
 - Gerçek başarıya ne kadar yakın
 - %80 güven ile p∈[73,2 76,7]
- Örnek: N=100 olayda S=75 başarı sağlanmış.
 - Tahmin edilen başarı oranı: 75%
 - Gerçek başarıya ne kadar yakın
 - %80 güven ile *p*∈[69,1 80,1]

http://www3.itu.edu.tr/~sgunduz/courses/verimaden



Ortalama Değer ve Varyans

- Başarı oranı p olan tek bir Bernoulli olayının ortalama değeri ve varyansı: p, p (1–p)
- N kere tekrarlanan Bernoulli olayının beklenen başarı oranı f=S/N
- Büyük N değerleri için, f normal dağılım
- fiçin ortalama değer ve varyans: p, p (1−p)/N
- Ortalama değeri 0 ve varyansı 1 olan X rastlantı değişkeninin %c güven aralığı:

 $Pr[-z \le X \le z] = c$

Simetrik bir dağılım için:

 $Pr[-z \le X \le z] = 1 - 2*Pr[X \ge z]$



Güven Sınırları

Örtalama değeri 0 ve varyansı 1 olan bir normal dağılımın güven sınırları



Pr[<i>X</i> ≥ <i>z</i>]	Z
0.1%	3.09
0.5%	2.58
1%	2.33
5%	1.65
10%	1.28
20%	0.84
40%	0.25

- $Pr[-1,65 \le X \le 1,65] = 90\%$
- fin ortalama değerinin 0, varyansının 1 olacak şekilde dönüştürülmesi gerekir.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Dönüşüm

- fin ortalama değerinin 0, varyansının 1 olacak şekilde dönüştürülmesi için f p
- Güven aralığı

$$\Pr\left[-z \le \frac{f-p}{\sqrt{p(1-p)/N}} \le z\right] = c$$

p'nin değeri

$$p = \left(f + \frac{z^2}{2N} \pm z \sqrt{\frac{f}{N} - \frac{f^2}{N} + \frac{z^2}{4N^2}} \right) / \left(1 + \frac{z^2}{N} \right)$$

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Örnek

- f = 75%, N = 1000, c = 80% (z = 1.28): p $\in [0,732 - 0,767]$
- f = 75%, N = 100, C = 80% (Z = 1.28): p ∈ [0,691 – 0,801]

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Konular

- Sınıflandırma yöntemleri
- Model Değerlendirme
- Öğrenme, sınama, geçerleme kümelerini oluşturma
 - holdout
 - k-kat çapraz geçerleme
 - Bootstrap
- Sınıflandırıcıları birleştirme

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

58



Verinin Dengesiz Dağılımı

- Küçük veya dengesiz veri kümeleri için örnekler tanımlayıcı olmayabilir
- Veri içinde bazı sınıflardan çok az örnek olabilir
 - tibbi veriler: %90 sağlıklı, %10 hastalık
 - elektronik ticaret: %99 alışveriş yapmamış, %1 alışveriş yapmış
 - güvenlik: %99 sahtekarlık yapmamış, %1 sahtekarlık yapmış
- Örnek: Sınıf1: 9990 örnek, Sınıf2: 10 örnek
 - bütün örnekleri sınıf1'e atayan bir sınıflandırıcının hata oranı: 9990 / 10000= %99,9
 - hata oranı yanıltıcı bir ölçüt olabilir

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

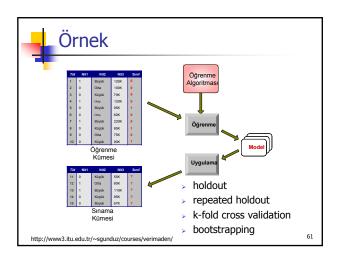


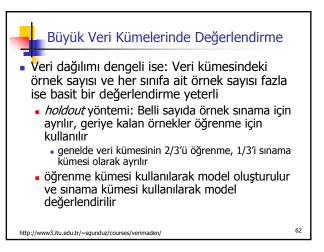
57

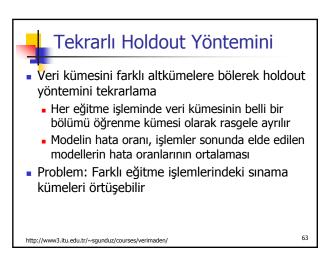
Dengeli Dağılım Nasıl Sağlanır?

- Veri kümesinde iki sınıf varsa
 - iki sınıfın eşit dağıldığı bir veri kümesi oluştur
 - Az örneği olan sınıftan istenen sayıda rasgele örnekler seç
 - Çok örneği olan sınıftan aynı sayıda örnekleri ekle
- Veri kümesinde iki sınıftan fazla sınıf varsa
 - Öğrenme ve sınama kümesini farklı sınıflardan aynı sayıda örnek olacak şekilde oluştur

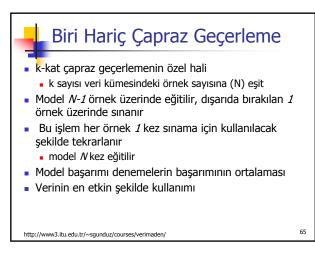
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

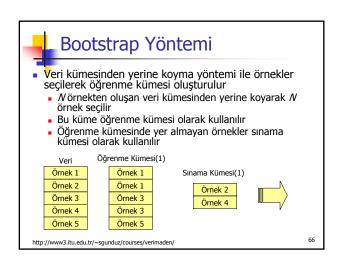














0.632 bootstrap

- N örnekten oluşan bir veri kümesinde bir örneğin seçilmeme olasılığı: $1-\frac{1}{N}$
- Sınama kümesinde yer alma olasılığı:

$$\left(1 - \frac{1}{n}\right)^n \approx e^{-1} = 0.368$$

 Öğrenme kümesi veri kümesindeki örneklerin %63,2'sinden oluşuyor

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



67

Bootstrap Yönteminde Model Hatasını Belirleme

- Model başarımını sadece sınama kümesi kullanarak belirleme kötümser bir yaklaşım
 - model örneklerin sadece ~%63'lük bölümüyle eğitiliyor
- Model başarımı hem öğrenme kümesindeki hem de sınama kümesindeki başarım ile değerlendirilir hata = 0,632 hata_(sınama) + 0,368 hata_(öğrenme)
- İşlem birkaç kez tekrarlanarak hatanın ortalaması alınır.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

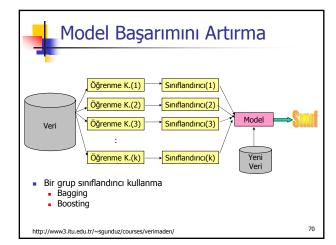
68



Konular

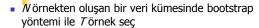
- Sınıflandırma yöntemleri
- Model Değerlendirme
- Öğrenme, sınama, geçerleme kümelerini oluşturma
- Sınıflandırıcıları birleştirme
 - Bagging
- Boosting

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/





Bagging



- ullet Bu işlemi k öğrenme kümesi oluşturmak üzere tekrarla
- Aynı sınıflandırma algoritmasını k öğrenme kümesi üzerinde kullanarak k adet sınıflandırıcı oluştur
- Yeni bir örneği sınıflandırmak için her sınıflandırıcının sonucunu öğren
- Yeni örnek en çok hangi sınıfa atanmışsa o sınıfın etiketiyle etiketlendir.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden



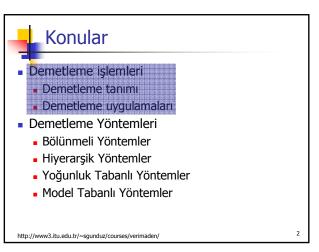
Boosting

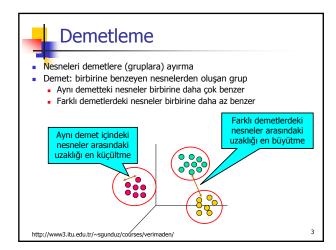
- Öğrenme kümesindeki her örneğin bir ağırlığı var
- Her öğrenme işleminden sonra, her sınıflandırıcı için yapılan sınıflandırma hatasına bağlı olarak örneklerin ağırlığı güncelleniyor
- Yeni bir örneği sınıflandırmak için her sınıflandırıcının doğruluğuna bağlı olarak ağırlıklı ortalaması alınıyor.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

/2

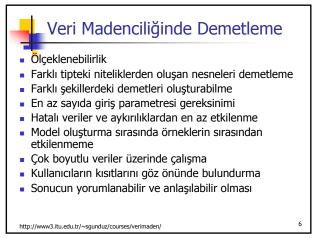








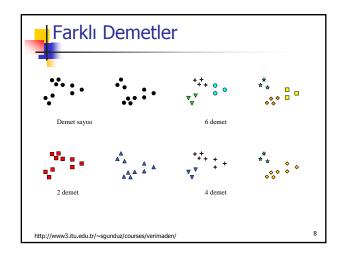




İyi Demetleme İyi demetleme yöntemiyle elde edilen demetlerin özellikleri aynı demet içindeki nesneler arası benzerlik fazla farklı demetlerde bulunan nesneler arası benzerlik az Oluşan demetlerin kalitesi seçilen benzerlik ölçütüne ve bu ölçütün gerçeklenmesine bağlı Uzaklık / Benzerlik nesnelerin nitelik tipine göre değişir Nesneler arası benzerlik: s(i,j)Nesneler arası uzaklık: d(i,j) = 1 - s(i,j)

- İyi bir demetleme yöntemi veri içinde gizlenmiş örüntüleri bulabilmeli
- Veriyi gruplama için uygun demetleme kriteri bulunmalı
- demetleme = aynı demetteki nesneler arası benzerliği enbüyüten, farklı demetlerdeki nesneler arası benzerliği enküçülten fonksiyon
- Demetleme sonucunun kalitesi seçilen demetlerin şekline ve temsil edilme yöntemine bağlı

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/





Temel Demetleme Yaklaşımları

- Bölünmeli yöntemler: Veriyi bölerek, her grubu belirlenmiş bir kritere göre değerlendirir
- Hiyerarşik yöntemler: Veri kümelerini (ya da nesneleri) önceden belirlenmiş bir kritère göre hiyerarşik olarak ayırır
- Yoğunluk tabanlı yöntemler: Nesnelerin yoğunluğuna göré demetleri oluşturur
- Model tabanlı yöntemler: Her demetin bir modele uyduğu varsayılır. Amaç bu modellere uyan verileri gruplamak

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Konular

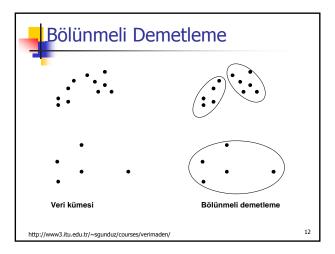
- Demetleme işlemleri
 - Demetleme tanımı
 - Demetleme uygulamaları
- Demetleme Yöntemleri
 - Bölünmeli Yöntemler
 - K-means demetleme yöntemi
 - K-medoids demetleme yöntemi
 - Hiyerarşik Yöntemler
 - Yoğunluk Tabanlı Yöntemler
 - Model Tabanlı Yöntemler

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Bölünmeli Yöntemler

- Amaç: n nesneden oluşan bir veri kümesini (D) k (k≤n) demete ayırmak
 - her demette en az bir nesne bulunmalı
 - her nesne sadece bir demette bulunmalı
- Yöntem: Demetleme kriterini enbüyütücek şekilde D veri kümesi k gruba ayırma
 - Global çözüm: Mümkün olan tüm gruplamaları yaparak en iyisini seçme (NP karmaşık)
 - Sezgisel çözüm: k-means ve k-medoids
 - ------ yozum. Kamednis ve Kamedolds kameans (MacQueen'67): Her demet kendi merkezi ile temsil edilir
 - k-medoids veya PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Her demet, demette bulunan bir nesne ile temsil edilir



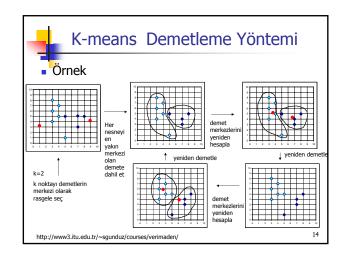
K-means Demetleme

- Bilinen bir k değeri için k-means demetleme algoritmasının 4 aşaması vardır:
 - veri kümesi *k* altkümeye ayrılır (her demet bir altküme)
 - Her demetin ortalaması hesaplanır: merkez nokta (demetteki nesnelerin niteliklerinin ortalaması)
 - Her nesne en yakın merkez noktanın olduğu demete dahil edilir
 - Nesnelerin demetlenmesinde değişiklik olmayana kadar adım 2'ye geri dönülür.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

13

15

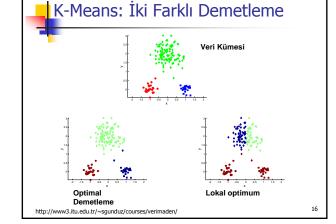


4

k-means Demetleme Yöntemi

- Demet sayısının belirlenmesi gerekir
- Başlangıçta demet merkezleri rasgele belirlenir
 - Her uygulamada farklı demetler oluşabilir
- Benzerlik Öklid uzaklığı, kosinüs benzerliği gibi yöntemlerle ölçülebilir
- Az sayıda tekrarda demetler oluşur
 - Yakınsama koşulu çoğunlukla az sayıda nesnenin demet değiştirmesi şekline dönüştürülür
- Karmaşıklığı:
 - Yer karmaşıklığı O((n+k) d)
 - Zaman karmaşıklığı O(ktnd)
 - k: demet sayısı, t: tekrar sayısı, n: nesne sayısı, d: nitelik sayısı

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



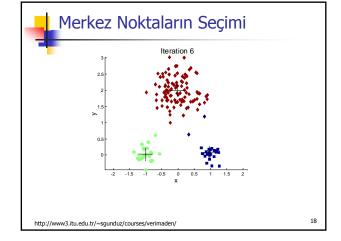


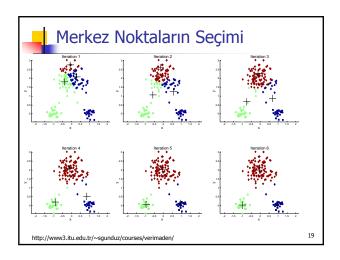
K-Means Demetleme Yöntemini Değerlendirme

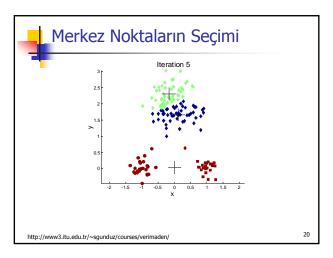
- Yaygın olarak kullanılan yöntem hataların karelerinin toplamı (Sum of Squared Error SSE)
 - Nesnelerin bulundukları demetin merkez noktalarına olan uzaklıklarının karelerinin toplamı

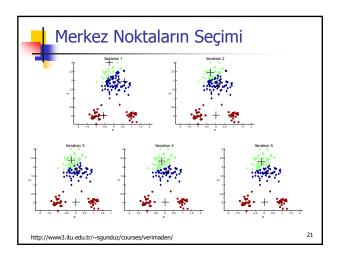
$$SSE = \sum_{i=1}^{K} \sum_{i=1}^{K} dist^{2}(m_{i}, x)$$

- x: C_i demetinde bulunan bir nesne, m_i : C_i demetinin merkez noktası
- Hataların karelerinin toplamını azaltmak için k demet sayısı artırılabilir
 - Küçük k ile iyi bir demetleme, büyük k ile kötü bir demetlemeden daha az SSE değerine sahip olabilir.
- Başlangıç için farklı merkez noktaları seçerek farklı demetlemeler oluşturulur
- En az SSE değerini sahip olan demetleme seçilir













K-Means Demetleme Algoritmasının Özellikleri

- Gerçeklemesi kolay
- Karmaşıklığı diğer demetleme 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

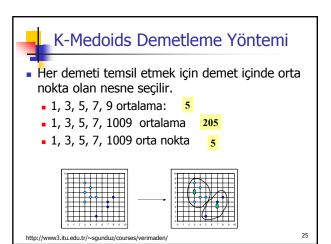
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

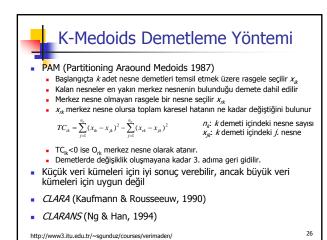


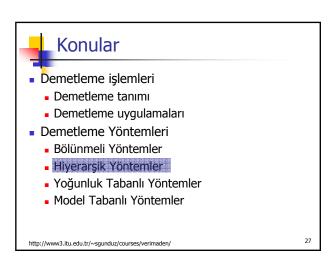
Konular

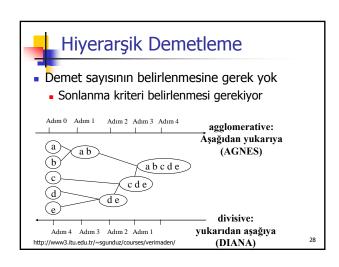
- Demetleme işlemleri
 - Demetleme tanımı
 - Demetleme uygulamaları
- Demetleme Yöntemleri
 - Bölünmeli Yöntemler
 - K-means demetleme yöntemi
 - K-medoids demetleme yöntemi
 - Hiyerarşik Yöntemler
 - Yoğunluk Tabanlı Yöntemler
 - Model Tabanlı Yöntemler

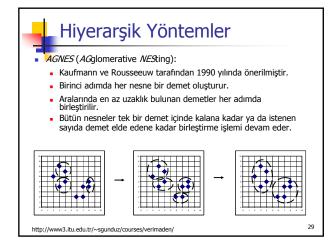
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

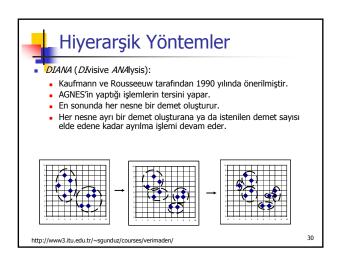


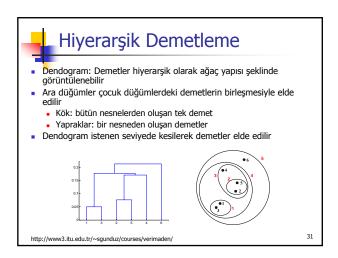


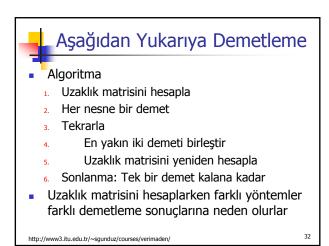


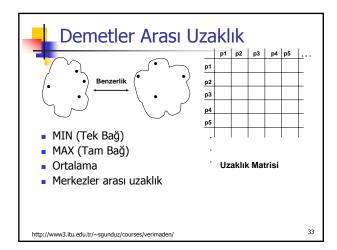


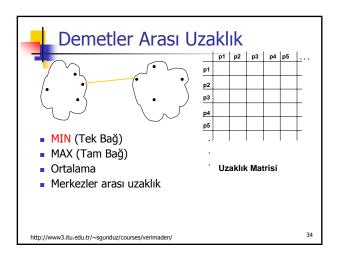


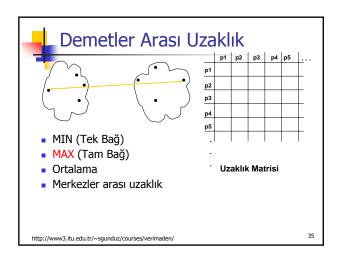


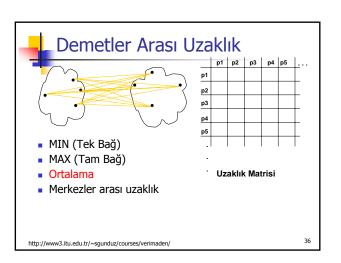


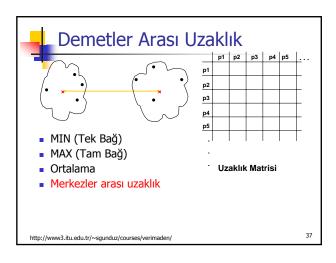


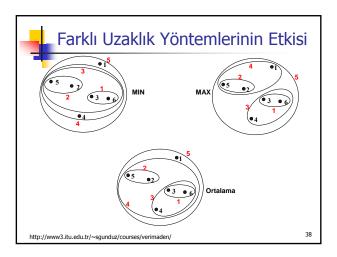














Hiyerarşik Demetleme Yöntemlerinin Özellikleri

- Demetleme kriteri yok
- Demet sayılarının belirlenmesine gerek yok
- Aykırılıklardan ve hatalı verilerden etkilenir
- Farklı boyuttaki demetleri oluşturmak problemli olabilir
- Yer karmaşıklığı O(n²)
- Zaman karmaşıklığı O(n²logn)

n: nesne sayısı

http://www3.itu.edu.tr/~squnduz/courses/verimaden/



Konular

- Demetleme işlemleri
 - Demetleme tanımı
 - Demetleme uygulamaları
- Demetleme Yöntemleri
 - Bölünmeli Yöntemler
 - Hiyerarşik Yöntemler
 - Yoğunluk Tabanlı Yöntemler
 - Model Tabanlı Yöntemler

http://www3.itu.edu.tr/~squnduz/courses/verimaden/



Yoğunluk Tabanlı Yöntemler

- Demetleme nesnelerin yoğunluğuna göre yapılır.
- Başlıca özellikleri:
 - Rasgele şekillerde demetler üretilebilir.
 - Aykırı nesnelerden etkilenmez.
 - Algoritmanın son bulması için yoğunluk parametresinin verilmesi gerekir.
- Başlıca yoğunluk tabanlı yöntemler:
 - <u>DBSCAN:</u> Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - <u>DENCLUE</u>: Hinneburg & D. Keim (KDD'98)
 - CLIQUE: Agrawal, et al. (SIGMOD'98)

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



DBSCAN

- İki parametre:

 - Eps: En büyük komşuluk yarıçapı
 MinPts: Eps yarıçaplı komşuluk bölgesinde bulunan en az nesne sayısı
- N_{eps}(p): {q∈D | d(p,q)≤Eps}
 Doğrudan erişilebilir nesne: Eps ve MinPts koşulları altında bir q nesnesinin doğrudan erişilebilir bir p nesnesi şu şartları sağlar:

 - p∈N_{eps}(q)
 q nesnesinin çekirdek nesne koşulunu sağlaması

 $N_{eps}(q) \ge MinPts$



MinPts = 5

Eps = 1 cm



DBSCAN

- Erisilebilir nesne:
 - Eps ve MinPts kosulları altında a nesnesinin erişilebilir bir p nesnesi olması icin:
 - p₁,p₂...,p_n nesne zinciri olması,
 - p1=q, pn=p,
 - p_i nesnesinin doğrudan erişilebilir nesnesi: p_{i+1}
- Yoğunluk bağlantılı Nesne:
 - Eps ve MinPts koşulları altında q nesnesinin yoğunluk bağlantılı nesnesi p şu koşulları sağlar:
 - p ve q nesneleri Eps ve MinPts koşulları altında bir o nesnesinin erişilebilir nesnesidir.



http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Yoğunluk Tabanlı Yöntemler: DBSCAN

- Veri tabanındaki her nesnenin Eps yarıçaplı komşuluk bölgesi araştırılır.
- Bu bölgede *MinPts*'den daha fazla nesne bulunan *p* nesnesi çekirdek nesne olacak şekilde demetler oluşturulur.
- Çekirdek nesnelerin doğrudan erişilebilir nesneleri
- Yoğunluk bağlantılı demetler birleştirilir.
- Hiçbir yeni nesne bir demete eklenmezse işlem sona
- Yer karmaşıklığı O(n)
- Zaman karmaşıklığı O(nlogn) n: nesne savisi

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

44



Konular

- Demetleme işlemleri
 - Demetleme tanımı
 - Demetleme uygulamaları
- Demetleme Yöntemleri
 - Bölünmeli Yöntemler
 - Hiyerarşik Yöntemler
 - Yoğunluk Tabanlı Yöntemler
 - Model Tabanlı Yöntemler

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Model Tabanlı Demetleme Yöntemleri

- Veri kümesi için öngörülen matematiksel model en uygun hale getiriliyor.
- Verinin genel olarak belli olasılık dağılımlarının karışımından geldiği kabul edilir.
- Model tabanlı demetleme yöntemi
 - Modelin yapısının belirlenmesi
 - Modelin parametrelerinin belirlenmesi
- Örnek EM (Expectation Maximization) Algoritması

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



Model Tabanlı Demetleme Yöntemleri



- K nesneden oluşan bir veri kümesi $D=\{x_{x},x_{2},...,x_{k}\}$ her x_{i} (i=[1,...K]) nesnesi Θ parametre kümesiyle tanımlanan bir olasılık dağılımından oluşturulur.
- Olasılık dağılımının, $c_j \in C = \{c_j, c_2, \dots, c_g\}$ şeklinde G adet bileşeni vardır. Her \mathbf{e}_g $g \in [1, \dots, G]$ parametre kümesi g bileşeninin olasılık dağılımın belirleyen, \mathbf{e} kümesinin ayrışık bir alt kümesidir.
- Herhangi bir x_i nesnesi öncelikle, $p(c_j|\pmb{\theta})=r_{g'}(\varSigma_G r_j=1)$ olacak şekilde) bileşen katsayısına (ya da bileşenin seçilme olasılığına) göre bir bileşene
- Bu bileşen $p(\mathbf{x}_i/c_{o'}; \mathbf{\Theta}_o)$ olasılık dağılımına göre \mathbf{x}_i değişkenini oluşturur.
- Böylece bir x_i nesnesinin bu model için olasılığı bütün bileşenlerin olasılıklarının toplamıyla ifade edilebilir:

$$p(x_i | \boldsymbol{\Theta}) = \sum_{g=1}^{G} p(c_g | \boldsymbol{\Theta}) p(\mathbf{x}_i | c_g; \boldsymbol{\Theta}_g)$$
$$p(x_i | \boldsymbol{\Theta}) = \sum_{g=1}^{G} \tau_g p(\mathbf{x}_i | c_g; \boldsymbol{\Theta}_g)$$



Model Tabanlı Demetleme Problemi

- Model parametrelerinin belirlenmesi
 - Maximum Likelihood (ML) yaklaşımı

$$\ell_{ML}(\Theta_1, ..., \Theta_G; \tau_1, ..., \tau_G \mid D) = \prod_{i=1}^K \sum_{g=1}^G \tau_g \, p(x_i \mid c_g, \Theta_g)$$

Maximum Aposteriori (MAP) yaklaşımı

$$\ell_{\mathit{MAP}}(\Theta_1, ..., \Theta_G; \tau_1, ..., \tau_G \mid D) = \prod_{i=1}^K \sum_{g=1}^G \frac{\tau_g \, p(x_i \mid c_g, \Theta_g) \, p(\Theta)}{p(D)}$$

Uygulamada her ikisinin logaritması

$$L(\Theta_1, ..., \Theta_G; \tau_1, ..., \tau_G \mid D) = \sum_{i=1}^{K} \ln \sum_{j=1}^{G} \left(\tau_g p(x_i \mid c_g, \Theta_g) \right)$$

 $L(\Theta_{1},...,\Theta_{G};\tau_{1},...,\tau_{G}\mid D) = \sum_{i=1}^{K}\ln\sum_{g=1}^{G}\left(\tau_{g}p(x_{i}\mid c_{g},\Theta_{g})\right) + \ln p(\Theta)$



EM Algoritması

- Veri kümesi: $D=\{x_1, x_2, ..., x_k\}$
- Gizli değişkenler $H=\{z_1,z_2,...,z_K\}$ (her nesnenin hangi demete dahil olduğu bilgisi)
- Verinin eksik olduğu durumda, tam verinin beklenen değeri hesaplanır:

$$\begin{split} Q(\Theta, \Theta') &= E[L_c(D, H \mid \Theta) \mid D, \Theta') \\ &= \sum_{i=1}^{K} \sum_{g=i}^{G} p(c_g \mid x_i) [\ln p(x_i \mid c_g) + \ln \tau_g] \end{split}$$

- EM Algoritmasının adımları:
 - ullet Θ' için başlangıç değerleri atama
 - (E) Expectation: $Q(\Theta|\Theta')$ hesaplanması
 - (M) Maximization: $argmax Q(\Theta|\Theta')$

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



VERİ MADENCİLİĞİ İlişkilendirme Kuralları

Yrd. Doç. Dr. Şule Gündüz Öğüdücü http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



İlişkilendirme Kuralları Madenciliği

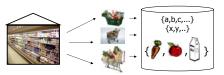
- İlişkilendirme kuralı madenciliği
 - Veri kümesi içindeki yaygın örüntülerin, nesneleri oluşturan öğeler arasındaki ilişkilerin bulunması
- İlişkilendirme kurallarını kullanma: veri içindeki kuralları belirleme
 - Hangi ürünler çoğunlukla birlikte satılıyor?
 - Kişisel bilgisayar satın alan bir kişinin bir sonraki satın alacağı ürün ne olabilir?
 - Yeni bir ilaca duyarlı olan DNA tipleri hangileridir?
 - Web dokümaları otomatik olarak sınıflandırılabilir mi?

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



İlişkilendirme Kuralları Bulma

Bir öğenin (veya öğeler kümesinin) varlığını harekette bulunan başka öğelerin varlıklarına dayanarak öngörme



hareketlerde alan öğeleri bulur

■ Kural şekli: "Gövde → Baş [destek, güven] "

satın alma(x, "ekmek") → satın alma (x, "süt") [%0.6, %65] öğrenci(x, "BLG"), kayıt(x, "VTYS") → not(x, "A") [%1, %75]

http://www3.itu.edu.tr/~squnduz/courses/verimaden/



Ilişkilendirme Kuralları Bulma

- Bütün öğelereden oluşan küme $I=\{i_1,i_2,...,i_d\}$ $I=\{$ ekmek, süt, bira, kola, yumurta, bez $\}$ Hareket $T_{\subseteq I}$, $T_{I}=\{f_{IJ},i_{J2}...,I_{IM}\}$ $TI=\{$ ekmek, süt $\}$ Hareketlerden oluşan veri kümesi $D=\{T_{JJ},T_{J2}...,T_{N}\}$



Market Alışveriş verisi

Hareket	Öğeler
T1	Ekmek, Süt
T2	Ekmek, Bez, Bira, Yumurta
Т3	Süt, Bez, Bira, Kola
T4	Ekmek, Süt, Bez, Bira
T5	Ekmek Süt Bez Kola

Yaygın öğeler:

Bez, bira Süt, ekmek, yumurta, kola Bira, ekmek, süt

Bulunan İlişkilendirme Kuralları

$$\begin{split} &\{\text{Bez}\} \rightarrow \{\text{Bira}\}, \\ &\{\text{Süt, Ekmek}\} \rightarrow \{\text{Yumurta, Kola}\}, \\ &\{\text{Bira, Ekmek}\} \rightarrow \{\text{Süt}\} \end{split}$$

http://www3.itu.edu.tr/~squnduz/courses/verimaden/



Yaygın Öğeler

- öğeler kümesi (Itemset)

 - Bir veya daha çok öğeden oluşan küme k-öğeler kümesi (k-itemset): k öğeden oluşan küme 3-öğeler kümesi: {Bez, Bira, Ekmek}
- Destek sayısı σ (Support count)
 - Bir öğeler kümesinin veri kümesinde görülme sıklığı
 - σ({Süt, Ekmek, Bez}) = 2
- Destek s (Support)
 - Bir öğeler kümesinin içinde bulunduğu hareketlerin toplam hareketlere oranı s({Süt, Ekmek, Bez}) = 2 /5
- Yaygın öğeler (Frequent itemset)
 - Destek değeri *minsup* eşik değerinden daha büyük ya da eşit olan öğeler kümesi

http://www3.itu.edu.tr/~sgunduz/courses/verimaden



İlişkilendirme Kuralları

- Veri kümesi $D=\{T_1,T_2,...,T_N\}$ en az, en küçük destek ve güven değerine sahip $X \to Y$ şeklinde kuralların bulunması $X \subset I, Y \subset I, X \cap Y = \emptyset$
- Kuralları değerlendirme ölçütleri destek (support) s: XYöğeler kümesinin bulunduğu hareket sayısının toplam hareket sayısına oranı

güven (confidence) c: X∪Y öğeler kümesinin bulunduğu hareket sayısının Xöğeler kümesi bulunan hareket sayısına oranı

 $confidence(X \rightarrow Y) = \frac{\#(X \cup Y)}{\#Y}$

TID	Öğeler
T1	Ekmek, Süt
T2	Ekmek, Bez, Bira, Yumurta
T3	Süt, Bez, Bira, Kola
T4	Ekmek, Süt, Bez, Bira
T5	Ekmek Siit Bez Kola

Örnek $\{\text{Süt, Bez}\} \Rightarrow \{\text{Bira}\}$

support = $\frac{\sigma(sut, bira, bez)}{\sigma(sut, bira, bez)} = \frac{2}{\pi} = 0.4$ $confidence = \frac{\sigma(sut, bira, bez)}{\sigma(sut, bira, bez)} = \frac{2}{2}$ $\sigma(sut,bez)$

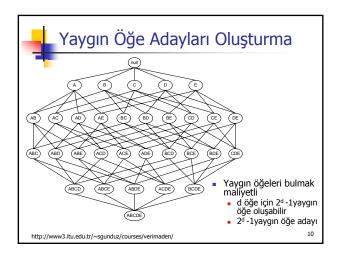


İlişkilendirme Kuralları Oluşturma

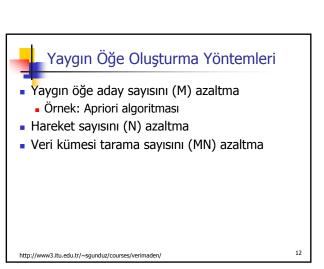
- İlişkilendirme kuralları madenciliğinde temel amaç D hareket kümesinden kurallar oluşturmak
 - kuralların destek değeri, belirlenen en küçük destek (minsup) değerinden büyük ya da eşit olmalı
 - kuralların güven değeri, belirlenen en küçük güven (minconf) değerinden büyük ya da eşit olmalı
- Brute-force yaklaşım
 - Olası bütün kuralları listele
 - Her kural için destek ve güven değeri hesapla
 - minsup ve minconf eşik değerlerinden küçük destek ve güven değerlerine sahip kuralları sil
 - hesaplama maliyeti yüksek













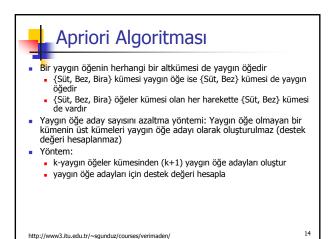
Apriori Algoritması

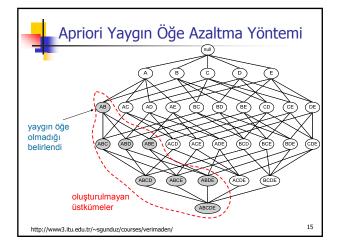
- Apriori yöntemi
 - Rakesh Agrawal, Ramakrishnan Srikant, Fast Algorithms for Mining Association Rules, Proc. 20th Int. Conf. Very Large Data Bases, VLDB'94
 - Heikki Mannila, Hannu Toivonen, Inkeri Verkamo, Efficient Algorithms for Discovering Association Rules. AAAI Workshop on Knowledge Discovery in Databases (KDD-94).
- Temel yaklaşım:

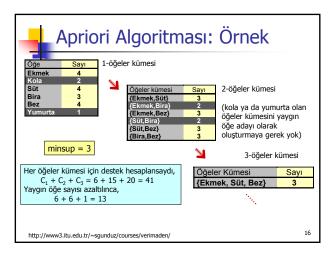
$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \ge s(Y)$$

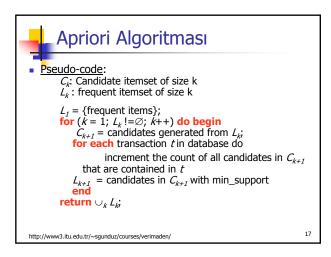
- Bir öğeler kümesinin destek değeri altkümesinin destek değerinden büyük olamaz
- anti-monotone özellik

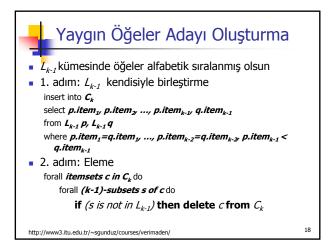
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/













Örnek: Yaygın Öğeler Adayı Oluşturma

- Yaygın öğe adayları olusturma:
 - L, kendisiyle birleştirilir (self join)
 - eleme
- Örnek:
 - L₃={abc, abd, acd, ace, bcd}
 - Kendisiyle birleştirme: L₃*L₃
 - abc ve abd öğeler kümesinden abcd
 - acd ve ace öğeler kümesinden acde
 - - ade L₂ kümesinin bir elemanı olmadığından acde yaygın öğelere dahil
 - C_a={ abcd}

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

19

21



Yaygın Öğelerden İlişkilendirme Kuralları Oluşturma

- Sadece güçlü ilişkilendirme kuralları oluşuyor
- Yaygın öğeler minsup değerini sağlıyor
- Güçlü ilişkilendirme kuralları *minconf* değerini sağlıyor.
- Güven $(A \rightarrow B) = Prob(B|A) = Destek(A \cup B)$ Destek(A)

- Yöntem:
 - Her yaygın öğeler kümesi f'in altkümelerini oluştur
 - Her altküme s için, $s \rightarrow (f-s)$ ilişkilendirme kuralı oluştur

 $destek(f) / destek(s) \ge minconf$ ise

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

20



Apriori Algoritmasını Geliştirme

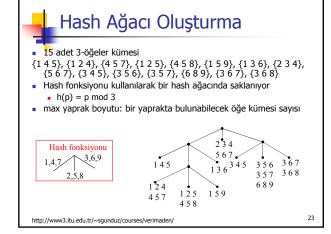
- Veritabanı tarama sayısını azaltma
 - k-yaygın öğeler kümesi için veritabanı k kez taraniyor
- Hash yöntemi ile veritabanı tarama sayısını azaltma
- Veritabanındaki hareket sayısını ve hareketlerdeki öğe sayısını azaltma
 - yaygın olmayan öğelerin veritabanında yer almasına gerek yok
- Arama uzayını parçalama

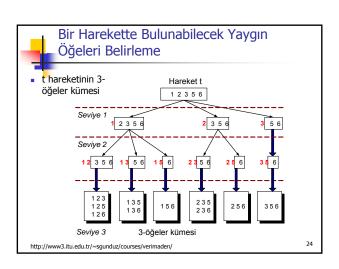
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

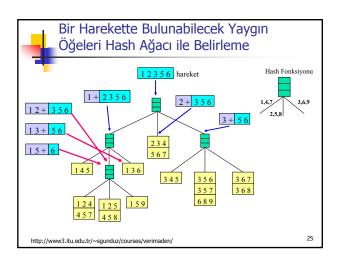
Eniyilime: Hash Ağacı

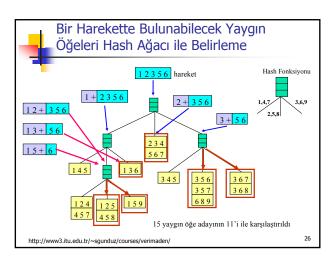
- Her yaygın öğe adayının destek değerinin belirlenmesi için veri kümesi taranır
 - Yaygın öğe adayı sayısı çok büyük olabilir
- Karşılaştırma sayısını azaltmak için yaygın öğe adayları hash ağacında saklanır
 - Her hareket bütün yaygın öğelerle karşılaştırılmak yerine hash ağacının ilgili bölümündeki yaygın öğeler adayları ile karşılaştırır

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/











Karmaşıklığı Etkileyen Etmenler

- Veri kümesindeki boyut sayısı (öğe sayısı)
 - her öğenin destek değerini saklamak için daha fazla saklama alanına ihtiyaç var
- Veri kümesinin büyüklüğü
 - Veri kümesindeki tarandığı için veri kümesindeki hareket sayısının fazla olması algoritmanın çalışma süresini uzatır
- Hareketlerin ortalama büyüklüğü (öğe sayısı)
 - Yoğun veri kümelerinde hareketlerdeki öğe sayısı fazla olur
 - Yaygın öğelerde daha fazla öğe olur
- En küçük destek değerini belirleme
 - En küçük destek değerini küçültme daha fazla sayıda yaygın öğe oluşmasına neden olur
 - Yaygın öğe adayı sayısının ve yaygın öğelerin boyunun artmasına neden olur

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/



27

Yaygın Öğeleri Belirlemede Sorunlar

- Sorunlar
 - Veri kümesinin birçok kez taranması

 - ven kumesinin pirçok kez taranması Yaygın öğeler aday sayısının fazlalığı Yaygın öğeleri bulmak için $I_2I_2...I_{100}$ veri kümesini tarama sayısı: 100• Aday sayısı: $(_{100}^2)$ + $(_{100}^2)$ + ... + $(_{100}^1)$ 0°) = 2^{100} ·1 = $1.27*10^{30}$! Yaygın öğeler adayları için destek değerinin hesaplanması

 - Veri kümesi tarama sayısını azaltma (S. Brin R. Motwani, J. Ullman, and S. Tsur. *Dynamic itemset counting and implication rules for market basket* data. In SIGMOD'97)
 - Vaygin oge aday sayisini azaltma (J. Park, M. Chen, and P. Yu. An effective hash-based algorithm for mining association rules. In SIGMOD'95)
 - Destek değerinin hesaplanmasını kolaylaştırma (M. Zaki et al. New algorithms for fast discovery of association rules. In KDD'97)
 - Yaygın öğeler adayı oluşturmadan yaygın öğeler bulunabilir. (J. Han, J. Pei, and Y. Yin, *Mining Frequent Patterns without Candidate Generation*. In SIGMOD'00.

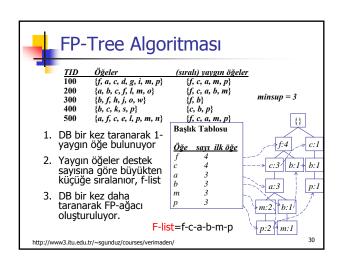
http://www3.itu.edu.tr/~squnduz/courses/verimaden/

28



Aday Oluşturmadan Yaygın Öğeleri Belirleme

- Kısa yaygın öğelere yeni öğeler eklenerek daha uzun yaygın öğeler elde etme
- Örnek:
 - "abc" bir yaygın öğe
 - Veri kümesinde içinde "abc" öğeleri bulunan hareketler (DB|abc)
 - DB|abc içinde d yaygın öğe olarak bulunursa: "abcd" yaygın öğe olarak belirlenir





FP-Ağacının Özelliği

- Bütünlük
 - Yaygın öğeleri bulmak için gerekli tüm bilgiyi barındırır
- Sıkıştırılmış
 - Yaygın olmayan öğeler FP-ağacında bulunmaz
 - Destek sayısı daha büyük olan öğeler köke daha yakın
 - Asıl veri kümesinden daha büyük değil

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

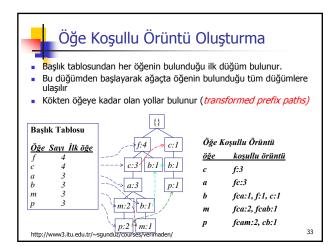
Örüntüleri ve Veri Kümesini Bölme

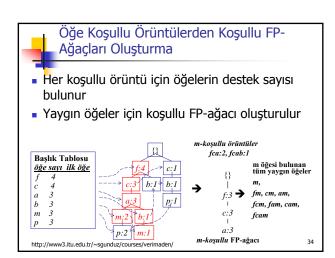
- Yaygın öğeler f-listesine göre altkümelere bölünür
 - F-list=f-c-a-b-m-p
 - p öğesi bulunan örüntüler
 - m öğesi bulunan ancak p öğesi bulunmayan örüntüler
 - ..

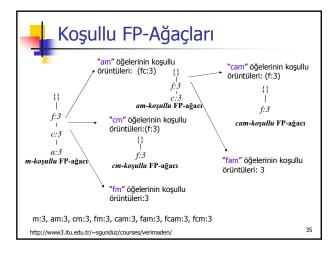
31

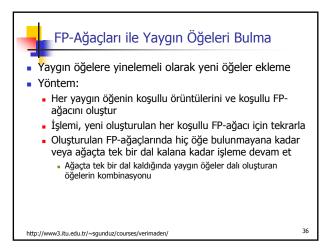
- c öğesi bulunan ancak a, b, m, p öğesi bulunmayan örüntüler
- f öğesi bulunan örüntüler

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/











İlişkilendirme Kuralları Oluşturma

- L yaygın öğelerden f⊂L altkümelerinin bulunması
 f→L-f kurallarının en küçük güven değeri koşulunu sağlaması gerekir
- Eğer {A,B,C,D} yaygın öğeler ise olası ilişkilendirme kuralları

- |L| = k için $2^k 2$ ilişkilendirme kuralı adayı vardır
 - L $\rightarrow \emptyset$ ve $\emptyset \rightarrow$ L kuralları geçerli kurallar değildir

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

37

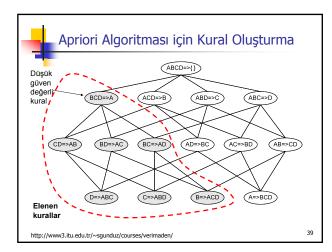


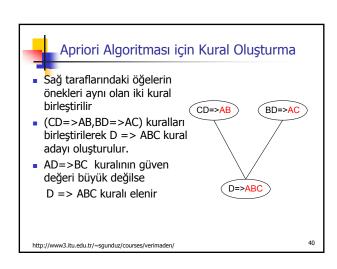
İlişkilendirme Kuralları Oluşturma

- İlişkilendirme kurallarının güven değerlerinin antimonotone özelliği yok
 - c(ABC \rightarrow D) değeri c(AB \rightarrow D) değerinden küçük ya da büyük olabilir
- Aynı yaygın öğeler kümesinden L = {A,B,C,D} oluşan ilişkilendirme kurallarının güven değerleri için antimonotone özelliği var
 - $c(ABC \rightarrow D) \ge c(AB \rightarrow CD) \ge c(A \rightarrow BCD)$
 - İlişkilendirme kuralının solunda bulunan öğe sayısı büyük olan kuralların güven değerleri de büyüktür.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

38







Destek Değerinin Etkisi

- minsup değeri büyük belirlenirse veri kümesinden bazı örüntüler elde edilmeyebilir:
 - veri kümesinde daha az bulunan
 - önemli bilgi taşıyan
- minsup değeri küçük belirlenirse
 - yöntem karmaşıklaşır
 - çok fazla sayıda yaygın öğeler kümesi elde edilir
- Tek bir destek değeri herzaman yeterli olmayabilir.

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

Bulunan Kuralların Önemi

- Tarafsız Ölçüt:
 - Örüntüler veri kümesinden elde edilen istatistiklere göre sıralanır
 - güven, destek, Jaccard, Gini, ...
- Taraflı Ölçüt
 - Örüntüler kullanıcının değerlendirmesine göre sıralanır
 - Bulunan örüntü kullanıcının beklentisi dışındaysa ilginçtir (Silberschatz & Tuzhilin)

http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

