

# VERİ MADENCİLİĞİ İlişkilendirme Kuralları

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



#### İlişkilendirme Kuralları Madenciliği

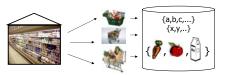
- İlişkilendirme kuralı madenciliği
  - Veri kümesi içindeki yaygın örüntülerin, nesneleri oluşturan öğeler arasındaki ilişkilerin bulunması
- İlişkilendirme kurallarını kullanma: veri içindeki kuralları belirleme
  - Hangi ürünler çoğunlukla birlikte satılıyor?
  - Kişisel bilgisayar satın alan bir 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	Elemak Siit Doz Kola

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

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

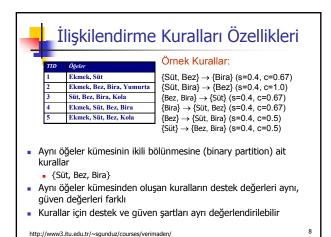
http://www3.itu.edu.tr/~sgunduz/courses/verimaden



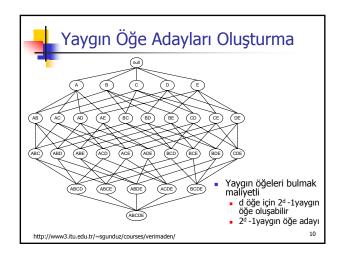
## İlişkilendirme Kuralları Oluşturma

- 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

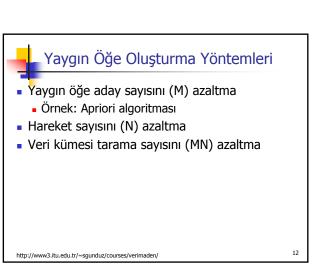
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/













## Apriori Algoritması

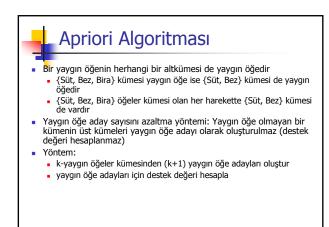
- Apriori yöntemi
  - Rakesh Agrawal, Ramakrishnan Srikant, Fast Algorithms for Mining Association Rules, Proc. 20th Int. Conf. Very Large Data Bases, VLDB'94
  - Heikki Mannila, Hannu Toivonen, Inkeri Verkamo, Efficient Algorithms for Discovering Association Rules. AAAI Workshop on Knowledge Discovery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

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

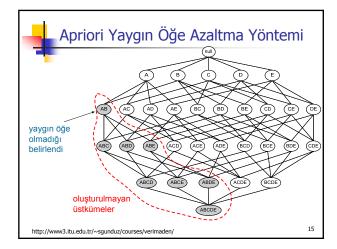
- Bir öğeler kümesinin destek değeri altkümesinin destek değerinden büyük olamaz
- anti-monotone özellik

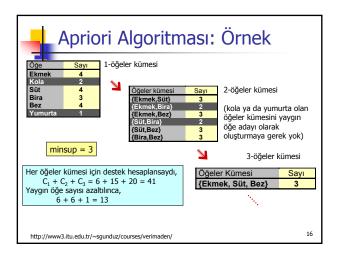
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

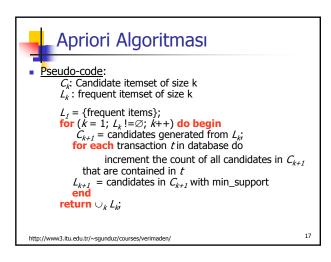
13

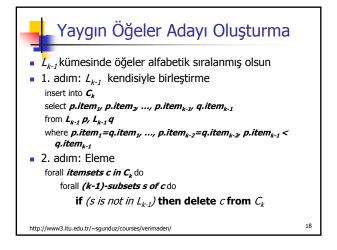


http://www3.itu.edu.tr/~sgunduz/courses/verimaden/











#### Örnek: Yaygın Öğeler Adayı Oluşturma

- Yaygın öğe adayları oluşturma:
  - L, kendisiyle birleştirilir (self join)
  - eleme
- Örnek:
  - L<sub>3</sub>={abc, abd, acd, ace, bcd}
  - Kendisiyle birleştirme: L<sub>3</sub>\*L<sub>3</sub>
    - abc ve abd öğeler kümesinden abcd
    - acd ve ace öğeler kümesinden acde
  - Fleme
    - ade L<sub>3</sub> kümesinin bir elemanı olmadığından acde yaygın öğelere dahil edilmez
  - C<sub>4</sub>={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 → (f-s) ilişkilendirme kuralı oluştur eğer:

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

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

20



### Apriori Algoritmasını Geliştirme

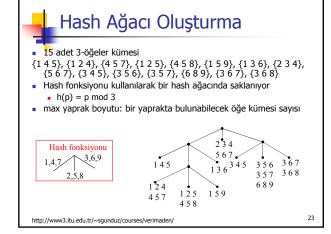
- Veritabanı tarama sayısını azaltma
  - k-yaygın öğeler kümesi için veritabanı k kez taranıyor
- Hash yöntemi ile veritabanı tarama sayısını azaltma
- Veritabanındaki hareket sayısını ve hareketlerdeki öğe sayısını azaltma
  - yaygın olmayan öğelerin veritabanında yer almasına gerek yok
- Arama uzayını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

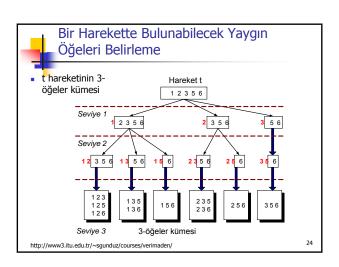
http://www3.itu.edu.tr/~sgunduz/courses/verimaden/

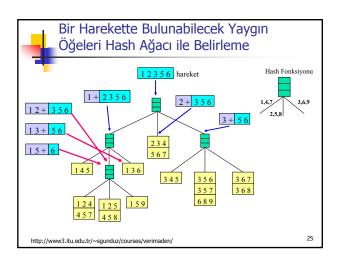
#### Eniyilime: Hash Ağacı

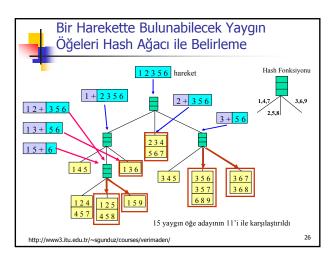
- Her yaygın öğe adayının destek değerinin belirlenmesi için veri kümesi taranır
  - Yaygın öğe adayı sayısı çok büyük olabilir
- Karşılaştırma sayısını azaltmak için yaygın öğe adayları hash ağacında saklanır
  - Her hareket bütün yaygın öğelerle karşılaştırılmak yerine hash ağacının ilgili bölümündeki 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:  $(_{100}^2)$  + ... +  $(_{11}^1\circ_0^0)$  =  $2^{100}$ ·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

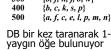


#### Aday Oluşturmadan Yaygın Öğeleri Belirleme

- 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

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

**FP-Tree Algoritması** Öğeler (sıralı) yaygın öğeler {f, a, c, d, g, i, m, p} {a, b, c, f, l, m, o} {b, f, h, j, o, w} {b, c, k, s, p} {f, c, a, m, p} {f, c, a, b, m} {f, b} 200



Yaygın öğeler destek sayısına göre büyükten küçüğe sıralanıor, f-list

3. DB bir kez daha taranarak FP-ağacı oluşturuluyor.

minsup = 3 $\{c, b, p\}$  $\{f, c, a, m, p\}$ 8 Baslık Tablosu f:4 sayı ilk öğe c:3 b:1 b:1 a:3 p:1 m:2 b:1 F-list=f-c-a-b-m-p p:2 m:1



## 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/

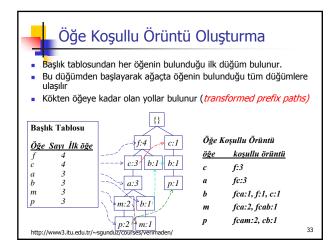
#### Örüntüleri ve Veri Kümesini Bölme

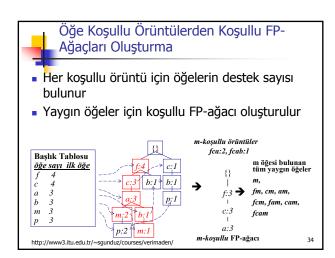
- Yaygın öğeler f-listesine göre altkümelere bölünür
  - F-list=f-c-a-b-m-p
  - p öğesi bulunan örüntüler
  - m öğesi bulunan ancak p öğesi bulunmayan örüntüler
  - •

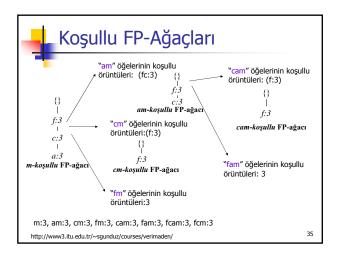
31

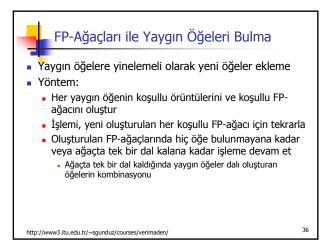
- c öğesi bulunan ancak a, b, m, p öğesi bulunmayan örüntüler
- f öğesi bulunan örüntüler

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 ightarrow ve  $\varnothing 
    ightarrow$  L kuralları geçerli kurallar değildir

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

37

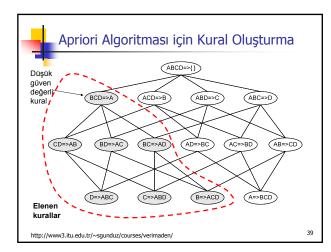


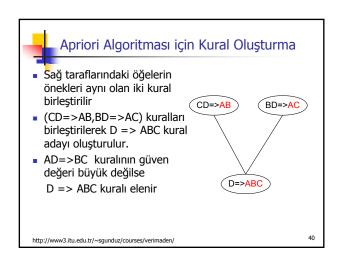
## İlişkilendirme Kuralları Oluşturma

- 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/

41

# 4

## 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/

