

VERİ MADENCİLİĞİ

Temel Sınıflandırma Yöntemleri

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

1

Konular

- Sınıflandırma işlemi
 - Sınıflandırma tanımı
 - Sınıflandırma uygulamaları
- Sınıflandırma yöntemleri
 - Karar ağaçları
 - Yapay sinir ağları
 - Bayes sınıflandırıcılar
 - Bayes ağıları

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

2

Sınıflandırma

- Sınıflandırma (classification) problemi:
 - nesnelerden oluşan veri kümesi (**öğrenme kümesi**):
 $D=\{t_1, t_2, \dots, t_n\}$
 - her nesne niteliklerden oluşuyor, niteliklerden biri **sınıf** bilgisi
- Sınıf niteliğini belirlemek için diğer nitelikleri kullanarak bir **model** bulma
- Öğrenme kümesinde yer almayan nesneleri (**sinama kümesi**) mümkün olan en iyi şekilde doğru sınıflara atamak
- **sınıflandırma**=ayrık değişkenler için öngörüde (prediction) bulunma

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

3

Sınıflandırma

- Amaç: Bir niteliğin değerini diğer nitelikleri kullanarak belirlemek
 - verinin dağılımına göre bir model bulunur
 - bulunan model, başarımı belirlendikten sonra niteliğin gelecekteki ya da bilmeyen değerini tahmin etmek için kullanılabilir
 - model başarımı: doğru sınıflandırılmış sinama kümesi örneklerinin oranı
- Veri madenciliği uygulamasında:
 - ayrıntılı değerlerini tahmin etmek: sınıflandırma
 - sürekli nitelik değerlerini tahmin etmek: öngörü



- Sınıflandırma: hangi topun hangi sepete koyulabileceği
- Öngörü: Topun ağırlığı

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

4

Gözetimli & Gözetimsiz Sınıflandırma

- Gözetimli (Supervised) sınıflandırma = sınıflandırma
 - Sınıfların sayısı ve hangi nesnenin hangi sınıfı olduğu biliniyor.
- Gözetimsiz (Unsupervised) sınıflandırma = demetleme (clustering)
 - Hangi nesnenin hangi sınıfı olduğu bilinmiyor. Genelde sınıf sayısı bilinmiyor.



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

5

Konular

- Sınıflandırma işlemi
 - Sınıflandırma tanımı
 - Sınıflandırma uygulamaları
- Sınıflandırma yöntemleri
 - Karar ağaçları
 - Yapay sinir ağları
 - Bayes sınıflandırıcılar
 - Bayes ağıları

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

6

Sınıflandırma Uygulamaları

- Kredi başvurusu değerlendirme
- Kredi kartı harcamasının sahtekarlık olup olmadığına karar verme
- Hastalık teşhisisi
- Ses tanıma
- Karakter tanıma
- Gazete haberlerini konularına göre ayırma
- Kullanıcı davranışları belirleme



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

7

Sınıflandırma için Veri Hazırlama

- Veri dönüşümü:
 - Sürekli nitelik değeri ayrık hale getirilir
 - Normalizasyon ($[-1, \dots, 1]$, $[0, \dots, 1]$)
- Veri temizleme:
 - gürültüyü azaltma
 - gereksiz nitelikleri silme

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

8

Sınıflandırma İşlemi

- Sınıflandırma işlemi üç aşamadan oluşur:
 1. Model oluşturma
 2. Model değerlendirme
 3. Modeli kullanma

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

9

Sınıflandırma İşlemi: Model Oluşturma

1. Model Oluşturma:
 - Her nesnenin sınıf etiketi olarak tanımlanan niteliğinin belirdiği bir sınıfta olduğu varsayılar
 - Model oluşturmak için kullanılan nesnelerin oluşturduğu veri kümlesi öğrenme kümlesi olarak tanımlanır
 - Model farklı biçimlerde ifade edilebilir
 - IF – THEN – ELSE kuralları ile
 - Karar ağaçları ile
 - Matematiksel formüller ile

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

10

Sınıflandırma İşlemi: Model Değerlendirme

2. Model Değerlendirme:
 - Modelin başarımı (doğruluğu) sınama kümesi örnekleri kullanılarak belirlenir
 - Sınıf etiketi bilinen bir sınama kümesi örneği model kullanılarak belirlenen sınıf etiketileyi karşılaştırılır
 - Modelin doğruluğu, doğru sınıflandırılmış sınama kümesi örneklerinin toplam sınama kümesi örneklerine oranı olarak belirlenir
 - Sınama kümesi model öğrenirken kullanılmaz

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

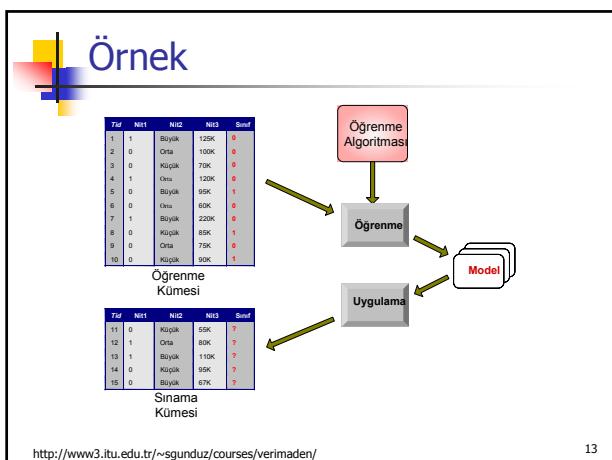
11

Sınıflandırma İşlemi: Modeli Kullanma

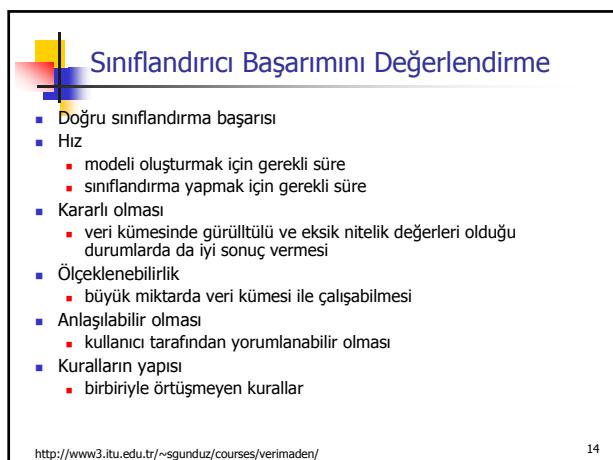
3. Modeli kullanma:
 - Model daha önce görülmemiş örnekleri sınıflandırmak için kullanılır
 - Örneklerin sınıf etiketlerini tahmin etme
 - Bir niteliğin değerini tahmin etme

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

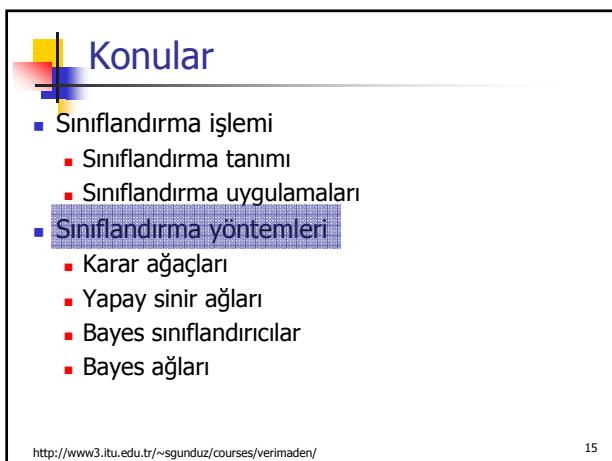
12



13



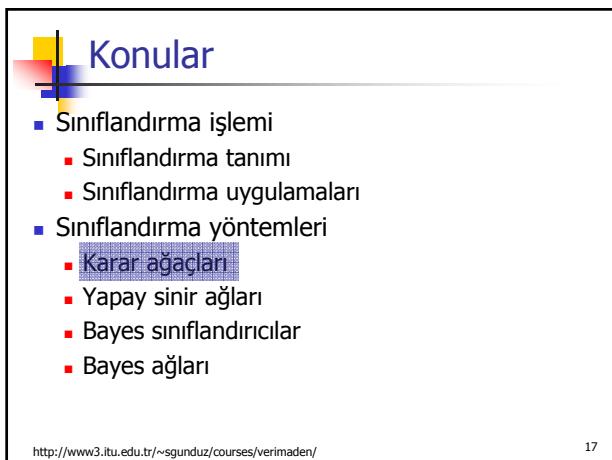
14



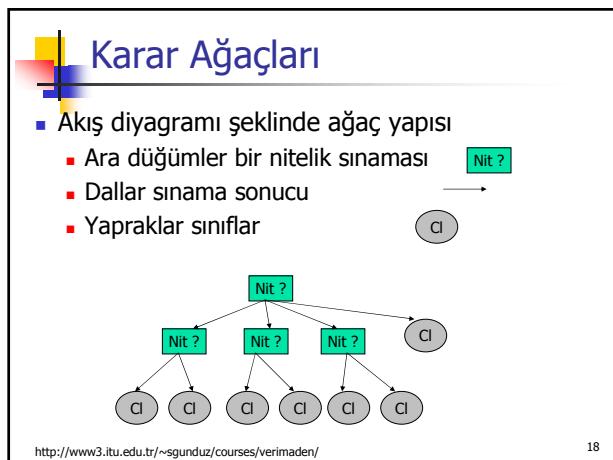
15



16



17



18

Örnek: Karar Ağacı

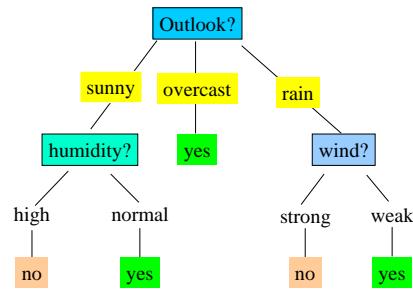
- J. Ross Quinlan'ın geliştirdiği ID3 modeline uyarlanmıştır:
 - hava tenis oynamaya uygun mu?

Day	Outlook	Temperature	Humidity	Wind	Play ball
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Strong	No
D6	Rain	Cool	Normal	Weak	Yes
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Hava durumu Verisi

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

Örnek: Karar Ağacı



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

20

Karar Ağacı Yöntemleri

- Karar ağacı oluşturma yöntemleri genel olarak iki aşamadan oluşur:
 - ağaç oluşturma
 - en başta bütün öğrenme kümesi örnekleri kökte seçilen niteliklere bağlı olarak örnek yinelemeli olarak bölünüyor
 - ağaç budama
 - öğrenme kümesindeki gürültülü verilerden oluşan ve sınıma kümesinde hataya neden olan dalları silme (sınıflandırma başarısını artırır)

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

21

Karar Ağacı Oluşturma

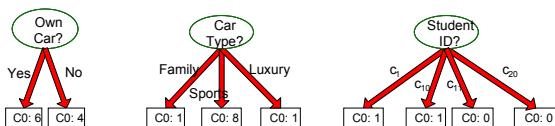
- Yinelemeli işlem
 - ağaç bütün verinin oluşturduğu tek bir düğümle başlıyor
- ci
 - eger örnekleri hepsi aynı sınıfa aitse düğüm yaprak olarak sonlanıyor ve sınıf etiketini alıyor
- Nit?
 - eger değilse örnekleri sınıflara en iyi bölecek olan nitelik seçiliyor
- işlem sona eriyor
 - örneklerin hepsi (çoğunluğu) aynı sınıfa ait
 - örnekleri bölecek nitelik kalmamış
 - kalan niteliklerin değerini taşıyan örnek yok

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

22

Örnekleri En İyi Bölgenin Nitelik Hangisi?

- Bölmenden önce:
 - 10 örnek C0 sınıfında
 - 10 örnek C1 sınıfında



Hangisi daha iyi?

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

23

En iyi Bölme Nasıl Belirlenir?

- "Greedy" yaklaşım
 - çoğunlukla aynı sınıfa ait örneklerin bulunduğu (homojen) düğümler tercih edilir
- Düğümlün kalitesini ölçmek için bir yöntem

C0: 5
C1: 5

homojen değil
kalitesi düşük

C0: 9
C1: 1

homojen
kalitesi yüksek

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

24

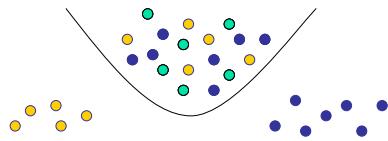
En İyi Bölün Nitelik Nasıl Belirlenir?

- İyilik Fonksiyonu (Goodness Function)
- Farklı algoritmalar farklı iyilik fonksiyonları kullanabilir:
 - bilgi kazancı (information gain): ID3, C4.5
 - bütün niteliklerin ayrık değerler aldığı varsayıiyor
 - sürekli değişkenlere uygulamak için değişiklik yapılabilir
 - gini index (IBM IntelligentMiner)
 - bütün niteliklerin sürekli değerler aldığı varsayıiyor
 - her nitelik için farklı bölmeye değerleri olduğu varsayıiyor
 - bölmeye değerlerini belirlemek için başka yöntemlere (demetleme gibi) ihtiyaç var
 - ayrık değişkenlere uygulamak için değişiklik yapılabilir

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

25

Bilgi Kazancı



sepetteki toplar farklı renklerde belirsizlik fazla topların hepsi aynı renkte ise daha belirsizlik yok

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

26

Bilgi / Entropi

- p_1, p_2, \dots, p_s toplamları 1 olan olasılıklar. Entropi (Entropy)

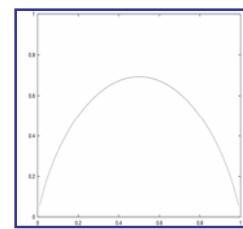
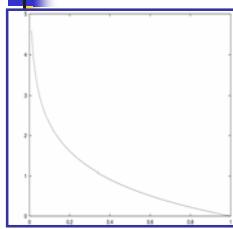
$$H(p_1, p_2, \dots, p_s) = -\sum_{i=1}^s p_i \log(p_i)$$

- Entropi rastgeleliği, belirsizliği ve beklenmeyen durumun ortaya çıkma olasılığını gösterir
- Sınıflandırmada
 - olayın olması beklenen bir durum
 - entropi=0

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

27

Entropi



- örnekler aynı sınıfa aitse entropi=0
- örnekler sınıflar arasında eşit dağılmışsa entropi=1
- örnekler sınıflar arasında rastgele dağılmışsa $0 < \text{entropi} < 1$

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

28

Örnek

- S veri kümesinde 14 örnek: C0 sınıfına ait 9, C1 sınıfına ait 5 örnek.
- Entropi

$$H(p_1, p_2, \dots, p_s) = -\sum_{i=1}^s p_i \log(p_i)$$

$$H(p_1, p_2) = - (9/14) \log_2 (9/14) - (5/14) \log_2 (5/14) = 0.940$$

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

29

Bilgi Kazancı (ID3 / C4.5)

- Bilgi kuramı kavramlarını kullanarak karar ağacı oluşturulur. Sınıflandırma sonucu için en az sayıda karşılaştırma yapmayı hedefler.
- Ağac bir niteliğe göre dallandığında entropi ne kadar düşer?
- A niteliğinin S veri kümesindeki bilgi kazancı

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|\text{S}_v|}{|S|} \text{Entropy}(\text{S}_v)$$

$\text{Values}(A)$, A niteliğinin alabileceğini değerler, S_v , $A=v$ olduğu durumda S'nin altkümesi.

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

30

Örnek

- Bilgi kazancına göre nitelik seçme

toplam örnek sayısı $s=14$, iki sınıfı ayrılmış
 $s_1=9(\text{yes}), s_2=5(\text{no})$
 $\text{Entropy}(S) = -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) = 0.940$
 wind için: weak=8, strong=6
 weak: no=2, yes=6
 strong: no=3, yes=3
 $\text{Entropy}(S_{\text{weak}}) = -(6/8) \log_2(6/8) - (2/8) \log_2(2/8) = 0.811$
 $\text{Entropy}(S_{\text{strong}}) = -(3/6) \log_2(3/6) - (3/6) \log_2(3/6) = 1.00$
 $\text{Gain}(\text{wind}) = 0.940 - (8/14) * 0.811 - (6/14) * 1.00$

$\text{Gain}(\text{Outlook}) = 0.246$
$\text{Gain}(\text{Humidity}) = 0.151$
$\text{Gain}(\text{wind}) = 0.048$
$\text{Gain}(\text{Temperature}) = 0.029$

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

31

Örnek

```

graph TD
    A[Outlook?] -- sunny --> B{?}
    A -- overcast --> C{?}
    A -- rain --> D{?}
    B -- yes --> E[Sunny]
    B -- no --> F{?}
    C -- yes --> G[Overcast]
    C -- no --> H{?}
    D -- yes --> I[Rain]
    D -- no --> J{?}
    E --> K[Day: Sunny, Outlook: Overcast, Temperature: Cool, Humidity: Normal, Wind: Weak, Play ball: Yes]
    F --> L[Day: Overcast, Outlook: Overcast, Temperature: Cool, Humidity: Normal, Wind: Weak, Play ball: Yes]
    G --> M[Day: Rain, Outlook: Rain, Temperature: Cool, Humidity: Normal, Wind: Weak, Play ball: Yes]
    H --> N[Day: Rain, Outlook: Rain, Temperature: Cool, Humidity: Normal, Wind: Strong, Play ball: Yes]
    I --> O[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: No]
    J --> P[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: Yes]
    K --> Q[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    L --> R[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    M --> S[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    N --> T[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    O --> U[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: No]
    P --> V[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: Yes]
    Q --> W[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    R --> X[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    S --> Y[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    T --> Z[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    U --> AA[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: No]
    V --> BB[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: Yes]
    W --> CC[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    X --> DD[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    Y --> EE[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    Z --> FF[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    AA --> GG[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: No]
    BB --> HH[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: High, Wind: Strong, Play ball: Yes]
    CC --> II[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    DD --> JJ[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Strong, Play ball: Yes]
    EE --> KK[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
    FF --> LL[Day: Rain, Outlook: Rain, Temperature: Mild, Humidity: Normal, Wind: Weak, Play ball: Yes]
  
```

$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$, $\text{Entropy}(S_{\text{sunny}}) = 0.970$
 humidity için: high=3, normal=2
 high: no=3, yes=0
 normal: no=0, yes=2
 $\text{Entropy}(S_{\text{high}}) = 0$
 $\text{Entropy}(S_{\text{normal}}) = 0$
 $\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = 0.970 - (3/5)0.0 - (3/5)0.0 = 0.970$

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

32

Gini Index (IBM IntelligentMiner)

- Veri kümesi S içinde n sınıf varsa ve p_j C_j sınıfının olasılığı ise

$$gini(S) = 1 - \sum_{j=1}^n p_j^2$$

- Eğer veri kümesi S_1 ve S_2 altkümelere bölünüyorsa ve her altkümede sırasıyla N_1 ve N_2 örnek varsa:

$$gini_{split}(S) = \frac{N_1}{N} gini(S_1) + \frac{N_2}{N} gini(S_2)$$

- Gini Index değeri en küçük olan nitelik seçilir.

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

33

Örnek

$$GINI(S) = 1 - \sum_j [p_j]^2$$

C1	0
C2	6

$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$
 $\text{Gini} = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$

C1	1
C2	5

$P(C1) = 1/6 \quad P(C2) = 5/6$
 $\text{Gini} = 1 - (1/6)^2 - (5/6)^2 = 0.278$

C1	2
C2	4

$P(C1) = 2/6 \quad P(C2) = 4/6$
 $\text{Gini} = 1 - (2/6)^2 - (4/6)^2 = 0.444$

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

34

Sürekli Nitelikleri Bölme

- Bölmeleme:
 - Statik: En başta bölmelenir
 - Bölmeler eşit genişlik, eşit derinlik veya demetleme yöntemi ile bulunur.
 - Dinamik:
 - Sürekli nitelik A sıralanır. Birbirini izleyen ancak sınıf etiketi farklı olan nitelik değerleri bulunur. En fazla kazanç sağlayan bölge seçilir.

Temperature	40	48	60	72	80	90
Play tennis	No	No	Yes	Yes	Yes	No

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

35

Ağaç Oluşturmadan Temel Yaklaşımlar

- Bölme kriteri:
 - ağacın bir düğümünde karşılaştırma yapılacak niteliğin seçilmesi
 - farklı algoritmalar farklı iyilik fonksiyonları kullanabilir: bilgi kazancı, gini index,...
- Dallanma kriteri:
 - bir örneğin hangi dala ait olduğunu belirleme
 - ikiye dallanma (gini index), çoklu dallanma (bilgi kazancı)
- Durma kararı:
 - dallanma işleminin devam edip etmeyeceğine karar verme
- Etiketleme kuralı:
 - yaprak düğüm en çok örneği olan sınıfla etiketleniyor

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

36

Ağaç Oluşturma:

- Parçala ve çöz (divide and conquer)
 - kökten yapraklara
 - düğümü dallara ayırl
 - 'Greedy' algoritma
 - her adımda en iyi çözümü bul: her düğümde dallanmak için en iyi niteliği bul
 - her dal için algoritmayı uygula

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

37

Örnek Algoritma: ID3

- Bütün nitelikler ayrık
- Bir düğüm oluştur N:
 - Eğer örneklerin hepsi C sınıfına ait ise, N düşümü C etiketli yaprak
 - Eğer karşılaştırma yapılacak nitelik yoksa N düşümü en çok örnegi olan sınıf
- En büyük bilgi kazancı olan niteliği bölmek için seç
 - N'yi seçilen nitelik ile etiketle
 - niteliğin her A_i değeri için bir dal oluştur
 - S_i, örneklerin hepsinin A_i değeri aldığı dal
 - S_i boş —> bir yaprak oluşturup en çok örnegi olan sınıfla etiketle
 - S_i boş değil —> algoritmayı S_i düşümü üzerinde yinele
- Yaprak düğümlere kadar

Ayrıntılı bilgiler: http://dms.irb.hr/tutorial/tut_dtrees.php

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

38

Örnek Algoritma: C4.5

- Kökten yapraklara doğru ağaç oluşturma
- Bilgi kazancı yöntemini kullanıyor
- Bütün veri kümelerini bellekte tutuyor
 - Büyük veri kümeleri için uygun değil

<http://www.rulequest.com/Personal/c4.5r8.tar.gz>

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

39

Karar Ağacı Kullanarak Sınıflandırma

- Doğrudan
 - sınıflandırmak istenilen örnegin nitelikleri ağaç boyunca sınamır
 - ulaşılan yaprağın etiketi sınıf bilgisini verir
- Dolaylı
 - karar ağacı sınıflandırma kurallarına dönüştürülür
 - kökten yaprakların herbirene giden yollar için ayrı bir kural oluşturulur.
 - IF-THEN şeklinde kuralları insanlar daha kolay anlıyor
 - Örnek: IF Outlook="sunny" AND humidity="normal" THEN play tennis

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

40

Karar Ağacı Kullanarak Sınıflandırma

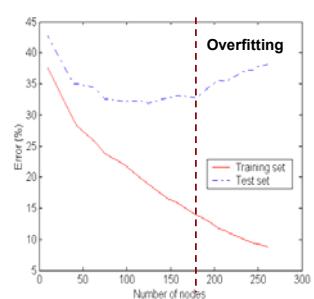
- Avantajları:
 - Karar ağacı oluşturmak zahmetsziz
 - Küçük ağaçları yorumlamak kolay
 - Anlaşılabili kurallar oluşturulabilir
 - Sürekli ve ayrı nitelik değerleri için kullanılabilir
- Dezavantajları:
 - Sürekli nitelik değerlerini tahmin etmekte çok başarılı değil
 - Sınıf sayısı fazla ve öğrenme kümlesi örnekleri sayısı az olduğunda model oluşturma çok başarılı değil
 - Zaman ve yer karmaşıklığı öğrenme kümlesi örnekleri sayısına (q), nitelik sayısına (h) ve oluşan ağaçın yapısına bağlı.
 - Ağaç oluşturma karmaşıklığı fazla, ağaç budama karmaşıklığı fazla
 - ağaç oluşturmak için zaman karmaşıklığı: $O(h \log q)$

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

41

Karar Ağaçlarında Aşırı Öğrenme

- Öğrenme kümelerindeki örneklerin azlığı veya gürültülü olması
- Aşırı öğrenmeye engelleyen iki yaklaşım
 - işlemi erken sona erdirme
 - işlemi sona eritmek için eşik değeri belirlemek gerekiyor
 - karar ağaçları oluştuktan sonra ağaç küçültme

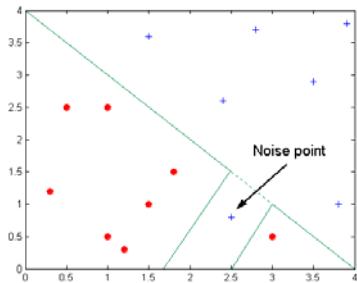


<http://www3.itu.edu.tr/~sgunduz/courses/verimaden/>

42



Aşırı Öğrenme: Gürültülü Örnekler



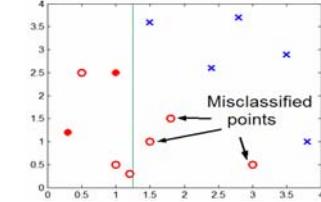
Gürültülü örnekler nedeniyle sınıfları ayıran düzlemin bozulması

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

43



Aşırı Öğrenme: Yetersiz Öğrenme Kümesi



- Öğrenme kümesindeki örnek sayısının yetersiz olması nedeniyle sınama kümesindeki örneklerin yanlış sınıflandırılması

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

44



Aşırı Öğrenme

- Gereğinden fazla karmaşık karar ağaçları aşırı öğrenmeye neden oluyor.
- Karar ağacının yeni örnekler üzerindeki başarısını tahmin etmek için öğrenme kümesi örnekleri yeterli olmuyor.
- Hatayı tahmin etmek için farklı yöntemler gereklidir.

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

45



Genel Hatayı Tahmin Etme

- Yerine koyma (Resubstitution) Hatası: öğrenme kümesi kullanılarak hesaplanan hata ($\sum e(t)$)
- Genel (Generalization) hata: sınama kümesi kullanılarak hesaplanan hata ($\sum e'(t)$)
- Genel hatayı tahmin etme yöntemleri:
 - İyimser yaklaşım: $e'(t) = e(t)$
 - Kötümser yaklaşım:
 - Her yaprak düğüm için: $e'(t) = (e(t)+0.5)$
 - Toplam hata: $e'(T) = e(T) + N \times 0.5$ (N : yaprak düğüm sayısı)
 - 30 yaprak düğümü olan bir karar ağıacı, 1000 öğrenme kümesi örneğinden 10 örneği yanlış sınıflandırırsa
 - Yerine koyma hatası: $10/1000 = \%1$
 - Genel hata: $(10 + 30 \times 0.5)/1000 = \% 2.5$
 - Ağacı budama: Genel hatayı tahmin etmek için geçerleme kümesi kullanılır

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

46



Occam's Razor

- Genel hatası aynı olan iki modelden karmaşıklığı daha az olan seçilmeli
- Karmaşık modellerin veri içindeki gürültüyü öğrenme ihtimaleri daha fazla

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

47



Karar Ağacı Boyutunu Belirleme

- Veri kümesi öğrenme ve sınama kümesi olarak ayrılır
- Çapraz geçerleme kullanılır.
- Veri kümesinin tümü ağaç oluşturmak için kullanılır
 - İstatistiksel bir test ile (chi-square) düğüm eklemenin ya da ağaç küçültmenin katkısı sinanır
 - MDL (Minimum Description Length) yöntemi kullanılır: kodlama en aza indirildiğinde ağaçın büyümesi durdurulur

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

48

Konular

- Sınıflandırma işlemi
 - Sınıflandırma tanımı
 - Sınıflandırma uygulamaları
- Sınıflandırma yöntemleri
 - Karar ağaçları
 - **Yapay sinir ağları**
 - Bayes sınıflandırıcılar
 - Bayes ağları

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

49

Yapay Sinir Ağları ile Sınıflandırma

- İnsan beynindeki sinir hücrelerinin işlevini modelleyen bir yapı
- Birbiri ile bağlantılı katmanlardan oluşur.
 - katmanlar hücrelerden oluşur
- Katmanlar arasında iletim
- İleti katmanlar arasındaki bağın ağırlığına ve her hücrenin değerine bağlı olarak değişebilir

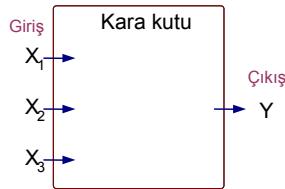


<http://www3.itu.edu.tr/~sgunduz/courses/verimaden/>

50

Örnek:

X ₁	X ₂	X ₃	Y
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	1
0	0	1	1
0	1	0	0
0	1	1	1
0	0	0	0

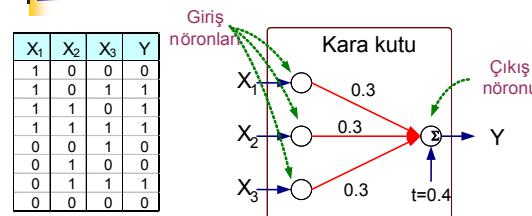


- En az iki giriş 1 ise çıkış 1, diğer durumlarda çıkış 0

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

51

Örnek



$$Y = I(0.3X_1 + 0.3X_2 + 0.3X_3 - 0.4 > 0)$$

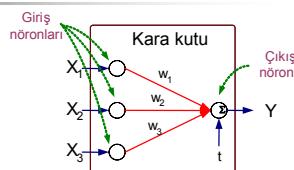
$$I(z) = \begin{cases} 1 & \text{eğer } z > 0 \\ 0 & \text{diğer} \end{cases}$$

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

52

Yapay Sinir Ağları

- Birbiri ile bağlantılı nöronlar ve ağırlıklar
- Çıkış nöronu kendisine gelen girişleri ağırlıkları olarak topluyor
- Çıkış nöronu bir eşik değeri ile karşılaştırılıyor



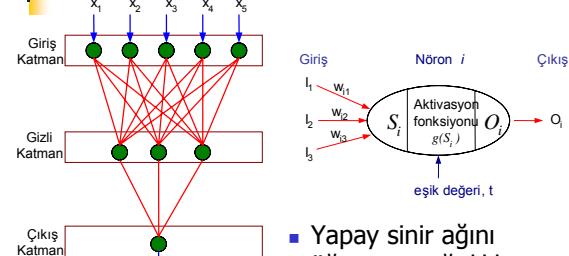
$$Y = I(\sum_i w_i X_i - t)$$

$$Y = \text{sign}(\sum_i w_i X_i - t)$$

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

53

Çok Katmanlı



- Yapay sinir ağını öğrenme: ağırlıkları öğrenme

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

54

Yapay Sinir Ağı ile Öğrenme

- Yapay sinir ağı oluşturma
 - giriş verisini modelleme
 - gizli katman sayısını, gizli katmanlardaki nöron sayısını belirleme
- Yapay sinir ağını eğitme
- Sinir ağını küçültme
- Sonucu yorumlama

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

55

Yapay Sinir Ağını Oluşturma

- Giriş nöron sayısı
 - Öğrenme kümelerindeki verilerin nitelik sayısı
- Gizli nöron sayısı
 - öğrenme sırasında ayarlanır
- Çıkış nöron sayısı
 - sınıf sayısı

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

56

Yapay Sinir Ağını Eğitme

- Amaç: Veri kümelerindeki örneklerin hepsini doğru sınıflandıracak ağırlıkları belirlemek
 - ağırlıklara rasgele değerler ata
 - öğrenme kümelerindeki giriş değerlerini teker teker sinir ağına uygula
 - çıkış hesapla
 - hata değerini hesapla $E = \sum_i [Y_i - f(w_i, X_i)]^2$
 - ağırlıkları hata fonksiyonunu enküçültucek şekilde düzelt

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

57

Yapay Sinir Ağını Küçültme

- Tam bağlı ağın anlaşılması çok güç
- n giriş nöron, h gizli nöron, m çıkış nöronu
 $h(m+n)$ ağırlık
- Küçültme: ağırlıklardan bazıları sınıflandırma sonucunu etkilemeyecek şekilde silinir

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

58

Yapay Sinir Ağları

- Avantajları
 - doğru sınıflandırma oranı genelde yüksek
 - kararlı – öğrenme kümelerinde hata olduğu durumda da çalışıyor
 - çıkış ayrık, sürekli ya da ayrık veya sürekli değişkenlerden oluşan bir vektör olabilir
- Dezavantajları
 - öğrenme süresi uzun
 - öğrenilen fonksiyonun anlaşılmaması zor

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

59

Konular

- Sınıflandırma işlemi
 - Sınıflandırma tanımı
 - Sınıflandırma uygulamaları
- Sınıflandırma yöntemleri
 - Karar ağaçları
 - Yapay sinir ağları
 - Bayes sınıflandırıcıları
 - Bayes ağları

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

60

Bayes (İstatistiksel) Modelleme

- Bayes teoremini kullanan istatistiksel sınıflandırıcı
- Örneklerin hangi sınıfı hangi olasılıkla ait oldukları
- Naïve Bayes sınıflandırıcı
 - niteliklerin hepsi aynı derecede önemli
 - nitelikler birbirinden bağımsız
 - bir niteliğin değeri başka bir nitelik değeri hakkında bilgi içermiyor
 - sınıflandırma ve öğrenme problemleri

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

61

Bayes Teoremi

- X sınıflandırılacak örnek. Hipotez h , X örneğinin C sınıfına ait olduğu
- h hipotezinin sonrasal olasılığı (*posteriori probability*)

$$P(h|X) = \frac{P(X|h)P(h)}{P(X)}$$
- MAP (maximum posteriori) hipotez

$$h_{MAP} = \arg \max_{h \in H} P(h|D) = \arg \max_{h \in H} P(D|h)P(h)$$
- Çok sayıda olasılığı önceden kestirmek gerekiyor

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

62

Örnek

- $H + \bigcirc \xrightarrow{H} C$
 - $P(H) = P(\text{apple}) \quad P(X) = P(H + \bigcirc)$
 - $P(X|H) = P(\text{apple}) \text{ ise } H + \bigcirc$
- $$P(h|X) = \frac{P(X|h)P(h)}{P(X)}$$

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

63

Naïve Bayes Sınıflandırıcı

- $X = (X_1, X_2, \dots, X_n)$ örneğinin C sınıfında olma olasılığı ($P(C|X)$) nedir?
- $\frac{P(X|C_i)P(C_i)}{P(X)}$ değerini enbüyütme
 $\rightarrow P(X|C_i)P(C_i)$ değerini enbüyütme
- $P(C_i) = |S_i| / |S|, \quad S_i: C_i$ sınıfına ait örneklerin sayısı
- $P(X|C_i) = \prod_{k=1}^n P(x_k|C_i), \quad P(x_k|C_i) = s_{ik} / s_i$
- Hesaplama maliyetini azaltıyor, sadece sınıf dağılımları hesaplanıyor
- Naïve: nitelikler bağımsız

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

64

Hava Durumu Verisi için Olasılıklar

Outlook	Temperature		Humidity		Windy		Play	
	Yes	No	Yes	No	Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4
Overcast	4	0	Mild	4	2	Normal	6	1
Rainy	3	2	Cool	3	1			
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5
Rainy	3/9	2/5	Cool	3/9	1/5			

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

65

Hava Durumu Verisi için Olasılıklar

Outlook	Temperature		Humidity		Windy		Play	
	Yes	No	Yes	No	Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4
Overcast	4	0	Mild	4	2	Normal	6	1
Rainy	3	2	Cool	3	1			
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5
Rainy	3/9	2/5	Cool	3/9	1/5			

■ Yeni veri

$$P(C_i|X) = P(X|C_i) \times P(C_i) = \prod_{k=1}^n P(x_k|C_i) \times P(C_i)$$

$$\begin{aligned} P(\text{"yes"}|X) &= 2/9 \times 3/5 \times 3/9 \times 3/9 \times 4/5 \times 3/5 = 0.0053 \\ P(\text{"no"}|X) &= 3/5 \times 1/5 \times 4/9 \times 2/5 \times 6/9 \times 1/5 = 0.0206 \end{aligned}$$

Normalize edilmiş olasılıklar:

$$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = 0.205$$

$$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = 0.795$$

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

66

Sürekli Veriler için Olasılık

- Verinin normal dağılımdan geldiği varsayıiyor.

- Her sınıf-nitelik çifti için bir olasılık hesaplanıyor.

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma^2}}$$

- Gelir için sınıf=-1

- ortalama=110
- varyans=2975

$$P(Gelir = 120 | -1) = \frac{1}{\sqrt{2\pi(2975)}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Tid	Geri Ödeme	Medeni Durum	Gelir	Dolan dinci
1	Evet	Bekar	125K	-1
2	Hayır	Evi	100K	-1
3	Hayır	Bekar	70K	-1
4	Evet	Evi	120K	-1
5	Hayır	Bosanmış	95K	1
6	Hayır	Evi	60K	-1
7	Evet	Bosanmış	220K	-1
8	Hayır	Bekar	85K	1
9	Hayır	Evi	75K	-1
10	Hayır	Bekar	90K	1

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

67

Örnek

X=(Geri Ödeme="Hayır", Medeni Durum="Evi", Gelir=120k)

$$\begin{aligned} P(\text{Geri Ödeme}=\text{"Evet"}| -1) &= 3/7 \\ P(\text{Geri Ödeme}=\text{"Hayır"}| -1) &= 4/7 \\ P(\text{Geri Ödeme}=\text{"Evet"}| 1) &= 0 \\ P(\text{Geri Ödeme}=\text{"Hayır"}| 1) &= 1 \\ P(\text{Medeni Durum}=\text{"Evi"}| -1) &= 4/7 \\ P(\text{Medeni Durum}=\text{"Bekar"}| -1) &= 2/7 \\ P(\text{Medeni Durum}=\text{"Bekar"}| 1) &= 2/3 \\ P(\text{Medeni Durum}=\text{"Boşanmış"}| 1) &= 1/3 \end{aligned}$$

$$\begin{aligned} \text{Gelir:} \\ \text{Sınıf} &= -1 \\ \text{ortalama} &= 110 \\ \text{varyans} &= 2975 \\ \text{Sınıf} &= 1 \\ \text{ortalama} &= 90 \\ \text{varyans} &= 25 \end{aligned}$$

Tid	Geri Ödeme	Medeni Durum	Gelir	Dolan dinci
1	Evet	Bekar	125K	-1
2	Hayır	Evi	100K	-1
3	Hayır	Bekar	70K	-1
4	Evet	Evi	120K	-1
5	Hayır	Bosanmış	95K	1
6	Hayır	Evi	60K	-1
7	Evet	Bosanmış	220K	-1
8	Hayır	Bekar	85K	1
9	Hayır	Evi	75K	-1
10	Hayır	Bekar	90K	1

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

68

Örnek

X=(Geri Ödeme="Hayır", Medeni Durum="Evi", Gelir=120k)

$$\begin{aligned} P(\text{Geri Ödeme}=\text{"Evet"}| -1) &= 3/7 \\ P(\text{Geri Ödeme}=\text{"Hayır"}| -1) &= 4/7 \\ P(\text{Geri Ödeme}=\text{"Evet"}| 1) &= 0 \\ P(\text{Geri Ödeme}=\text{"Hayır"}| 1) &= 1 \\ P(\text{Medeni Durum}=\text{"Evi"}| -1) &= 4/7 \\ P(\text{Medeni Durum}=\text{"Bekar"}| -1) &= 2/7 \\ P(\text{Medeni Durum}=\text{"Bekar"}| 1) &= 2/3 \\ P(\text{Medeni Durum}=\text{"Boşanmış"}| 1) &= 1/3 \end{aligned}$$

$$\begin{aligned} \text{Gelir:} \\ \text{Sınıf} &= -1 \\ \text{ortalama} &= 110 \\ \text{varyans} &= 2975 \\ \text{Sınıf} &= 1 \\ \text{ortalama} &= 90 \\ \text{varyans} &= 25 \end{aligned}$$

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

69

Olasılığın Sıfır Olması

- Her sınıfta bir niteliğin her değeri olmazsa

- koşullu olasılıklardan biri 0

- o sınıf'a ait olma olasılığı 0

Olasılıklar

$$\text{Original: } P(A_i | C) = \frac{N_{ic}}{N_c}$$

c: sınıf sayısı

Toplamları 1 olmak zorunda

$$\text{Laplace: } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

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

70

Bayes Sınıflandırıcılar

Avantajları:

- gerçeklemesi kolay
- çoğu durumda iyi sonuçlar

Dezavantajları

- varsayımlı: sınıf bilgisi verildiğinde nitelikler bağımsız
- gerçek hayatı değişkenler birbirine bağımlı
- değişkenler arası ilişki modellenemiyor

Çözüm:

- Bayes ağları

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

71

Konular

Sınıflandırma işlemi

- Sınıflandırma tanımı
- Sınıflandırma uygulamaları

Sınıflandırma yöntemleri

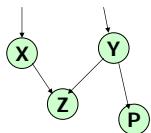
- Karar ağaçları
- Yapay sinir ağları
- Bayes sınıflandırıcılar
- Bayes ağları

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

72

Bayes Ağları

- Niteliklerin altkümesinin birbiri ile bağımsız olduğunu varsayıyor
- Yönlü çevrimsiz çizge (directed acyclic graph) ve koşullu olasılık tablolarından oluşur
- Her değişken A için bir tablo var
 - nitelığın ebeveynlerine olan koşullu olasılıkları

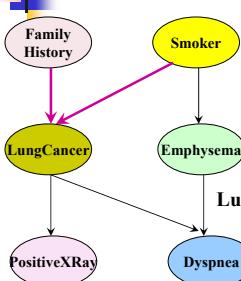


- düğümler: rasgele değişkenler
- ayritlar: olasılıklı bağılilik
- X ve Y, Z değişkeninin ebeveyni
- Y, P değişkeninin ebeveyni
- Z ve P arasında bağ yok

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

73

Örnek



(FH, S) (FH, ~S)(~FH, S) (~FH, ~S)

	LC	0.8	0.5	0.7	0.1
	~LC	0.2	0.5	0.3	0.9

LungCancer için koşullu olasılık tablosu

Bayes Ağları

$$P(LC="yes" | FH="yes", S="yes") = 0.8$$
$$P(LC="no" | FH="no", S="no") = 0.9$$

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

74

Bayes Ağlarının Eğitilmesi

- Ağ yapısı ve tüm değişkenler biliniyorsa koşullu olasılıklar hesaplanır
- Ağ yapısı belli ancak bazı değişkenler eksik ise yinelemeli öğrenme uygulanır
 - gradient descent algoritması
- D. Heckerman. [A Tutorial on Learning with Bayesian Networks](#). In *Learning in Graphical Models*, M. Jordan, ed., MIT Press, Cambridge, MA, 1999. Also appears as Technical Report MSR-TR-95-06, Microsoft Research, March, 1995.

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

75