

**BİM448-VERİ MADENCİLİĞİ VE BİLGİ KEŞFİ  
YILIÇI SINAVI**

Yrd.Doç.Dr. Songül ALBAYRAK

5 Mayıs 2009

Sınav Süresi: 90 dakika

Öğrencinin Adı ve Soyadı:

Öğrenci No:

**SORULAR**

1-[35puan] Aşağıda verilen 8 örnek ve 3 özellikten oluşan küçük veriseti için pozitif ve negatif sınıfları belirlemek için bir karar ağacı oluşturunuz. Karar ağacını oluşturmak için entropy ve bilgi kazancını hesaplayınız.

T Veriseti

Hair	Weight	Lotion	Result
blonde	light	no	sunburned (positive)
blonde	average	yes	none (negative)
brown	average	yes	none
blonde	average	no	sunburned
red	heavy	no	sunburned
brown	heavy	no	none
brown	heavy	no	none
blonde	light	yes	none

$$\text{Entropy}(T) = -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = 0,95 \quad \{3P, 5N\}$$

$$\begin{aligned} \text{Entropy}_{\text{hair}}(T) &= \frac{4}{8} \left( -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) + \frac{3}{8} \left( -\frac{3}{3} \log_2 \frac{3}{3} - 0 \right) + \frac{1}{8} \left( -1 \log_2 1 \right) = 0,5 \\ \text{Blonde} &= 4 \quad \{2+, 2-\} \\ \text{Brown} &= 3 \quad \{3-\} \\ \text{Red} &= 1 \quad \{1+\} \end{aligned}$$

$$\begin{aligned} \text{Entropy}_{\text{weight}}(T) &= \frac{2}{8} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{3}{8} \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{3}{8} \left( -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) = 0,94 \\ \text{Light} &= 2 \quad \{1+, 1-\} \\ \text{Average} &= 3 \quad \{2+, 1-\} \\ \text{Heavy} &= 3 \quad \{1+, 2-\} \end{aligned}$$

$$\begin{aligned} \text{Entropy}_{\text{lotion}}(T) &= \frac{5}{8} \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) + \frac{3}{8} \left( -\frac{3}{3} \log_2 \frac{3}{3} - 0 \right) = 0,61 \\ \text{No} &= 5 \quad \{3+, 2-\} \\ \text{Yes} &= 3 \quad \{3-\} \end{aligned}$$

Özellik

Kazanç

Hair

$$\text{Gain}_{\text{hair}} = \text{Entropy}(\tau) - \text{Entropy}(\tau)_{\text{hair}} = 0,95 - 0,5 = 0,45$$

Weight

$$\text{Gain}_{\text{weight}} = 0,95 - 0,94 = 0,01$$

Lotion

$$\text{Gain}_{\text{lotion}} = 0,95 - 0,61 = 0,34$$

Sağ taraf maksimum kazanç verdiği için karar ağacının köküne yerleştirilir.

Hair

Blonde

Brown

Red

Weight	Lotion	Result
light	no	+
average	yes	-
average	no	+
light	yes	-

Weight	Lotion	Result
average	yes	-
heavy	no	-
heavy	no	-

3 negative

Weight	Lotion	Result
heavy	no	+

1 positive

Lotion

Yes

No

Weight	Result
Average	-
Light	-

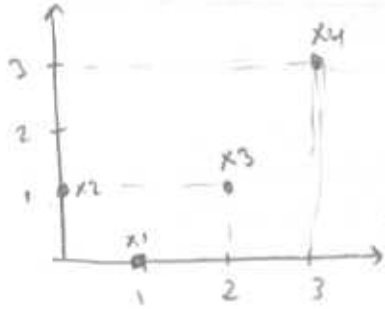
2 negative

Weight	Result
Light	+
Average	+

2 positive

2-[30 puan] Verisetinde iki boyutlu uzayda verilmiş 4 adet örnek bulunmaktadır.  $X_1 = \{1,0\}$ ,  $X_2 = \{0,1\}$ ,  $X_3 = \{2,1\}$ ,  $X_4 = \{3,3\}$

K-ortalamalı kümeleme yöntemine göre verisetini iki kümeye bölmek için başlangıç küme merkezleri  $X_1$  ve  $X_4$  olarak seçilmiş olsun. Algoritmanın ilk iki adımı için küme elemanlarını ve yeni küme merkezlerini hesaplayınız. Hesaplama Euclidean mesafesi kullanınız.



$$C_1 = (1,0) \text{ ve } C_2 = (3,3) \text{ ise}$$

$$X_1 \text{ için } d(X_1, C_1) = 0$$

$$d(X_1, C_2) = \sqrt{10}$$

$$X_2 \text{ için } d(X_2, C_1) = \sqrt{2}$$

$$d(X_2, C_2) = \sqrt{10}$$

$$X_3 \text{ için } d(X_3, C_1) = \sqrt{2}$$

$$d(X_3, C_2) = \sqrt{5}$$

$$X_4 \text{ için } d(X_4, C_1) = \sqrt{13}$$

$$d(X_4, C_2) = 0$$

$C_1$  kümesinin üyeleri  $\{X_1, X_2, X_3\} \Rightarrow$  yeni küme merkezi

$C_2$  kümesinin üyeleri  $\{X_4\} \Rightarrow$  yeni küme merkezi

$$C_1 = (1, \frac{2}{3})$$

$$C_2 = (3,3)$$

$$C_1 = (1, \frac{2}{3})$$

$$C_2 = (3,3)$$

$X_1$  için

$$d(X_1, C_1) = 2/3$$

$$d(X_1, C_2) = \sqrt{13}$$

$X_2$  için

$$d(X_2, C_1) = \sqrt{10/9}$$

$$d(X_2, C_2) = \sqrt{10}$$

$X_3$  için

$$d(X_3, C_1) = \sqrt{10/9}$$

$$d(X_3, C_2) = \sqrt{5}$$

$X_4$  için

$$d(X_4, C_1) = \sqrt{85/9}$$

$$d(X_4, C_2) = 0$$

$C_1$  kümesinin üyeleri  $\{X_1, X_2, X_3\} \Rightarrow C_1 = (1, 2/3)$

$C_2$  kümesinin üyeleri  $\{X_4\} \Rightarrow C_2 = (3,3)$

— İlk iterasyonda küme merkezleri değişmediği için iterasyon sonlanır.

3-[35 puan] Aşağıdaki tabloda verilen 10 örnek ve 3 özellikten oluşan verisetini eğitim veriseti olarak kullanarak Naive-Bayes sınıflayıcı modeli geliştirilecektir. Bu verisetinde 10 otomobilin rengi, modeli ve üretilidiği yer(yerli veya ithal) özelliklerine göre çalınma riski(var veya yok şeklinde) belirtilmiştir. Hangi sınıfa ait olduğu bilinmeyen (Red Domestic SUV) otomobilin çalınma riskini Naive Bayes sınıflayıcı kullanarak belirleyiniz.

Example No:	Color	Type	Origin	Stolen?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

5 YES  
5 NO

$$\begin{aligned}
 P(R, D, SUV | Yes) &= P(R | Yes) P(D | Yes) P(SUV | Yes) \times P(Yes) \\
 &= \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{8} \cdot \frac{5}{10} = \frac{6}{250}
 \end{aligned}$$

$$\begin{aligned}
 P(R, D, SUV | No) &= P(R | No) P(D | No) P(SUV | No) \times P(No) \\
 &= \frac{2}{5} \cdot \frac{3}{5} \cdot \frac{3}{5} \cdot \frac{5}{10} = \frac{18}{250}
 \end{aligned}$$

NO > YES

STOLEN = NO