

Makine Öğrenmesi

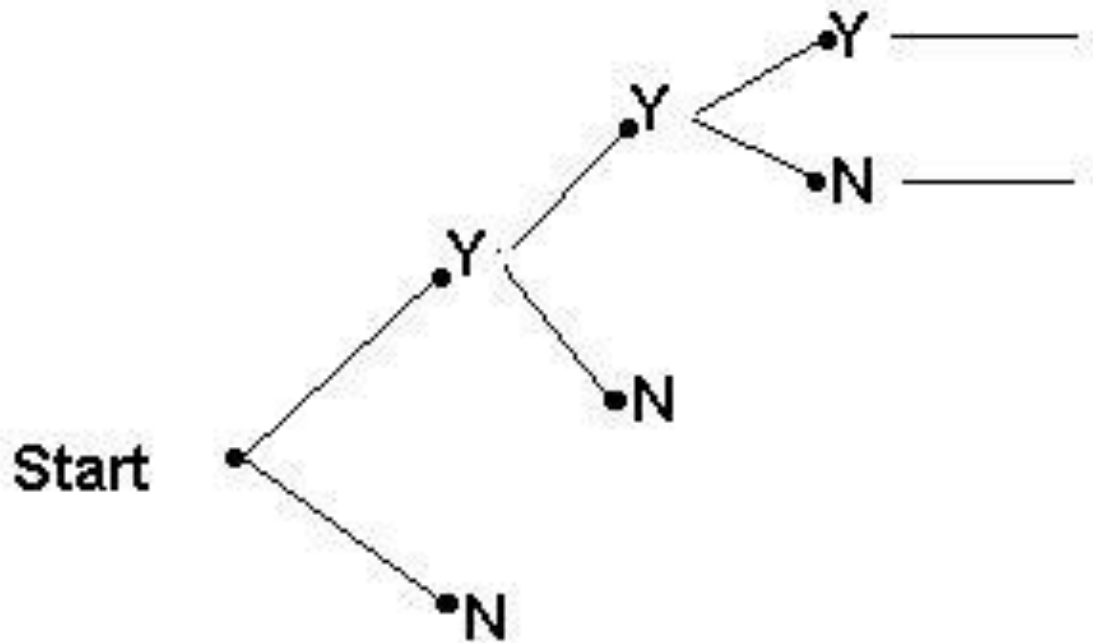


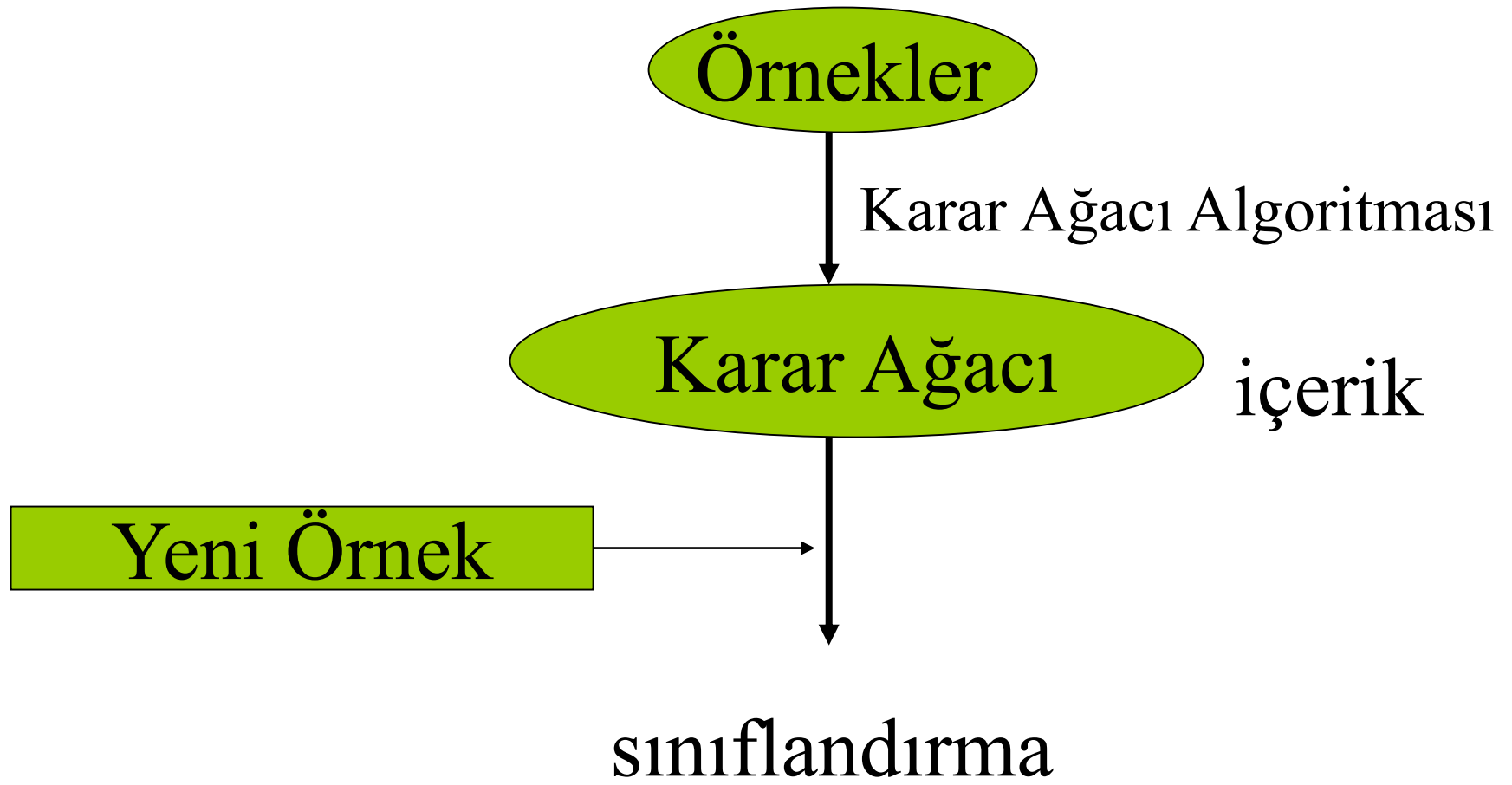
Karar Ağaçları

Doç. Dr. İlhan AYDIN

Karar Ağaçları (Decision trees)

Karar Ağaçları, Tree yapısında olup bir olayın sonuçlandırılmasında sorunun cevabına göre hareket ederler.





-
- Karar Ağacı algoritması 2 aşamadan oluşmaktadır:
 - Ağacı oluşturma
 - En başta bütün öğrenme kümesi ağaçtadır.
 - Ağacı budama
 - Öğrenme kümesindeki gürültülü verilerden oluşan ve test kümesinde hataya neden olan dallar silinir.

- **Entropy'e Dayalı Algoritmalar**

ID3, C 4.5

- **Sınıflandırma ve Regresyon Ağaçları (CART)**

Twoing, Gini

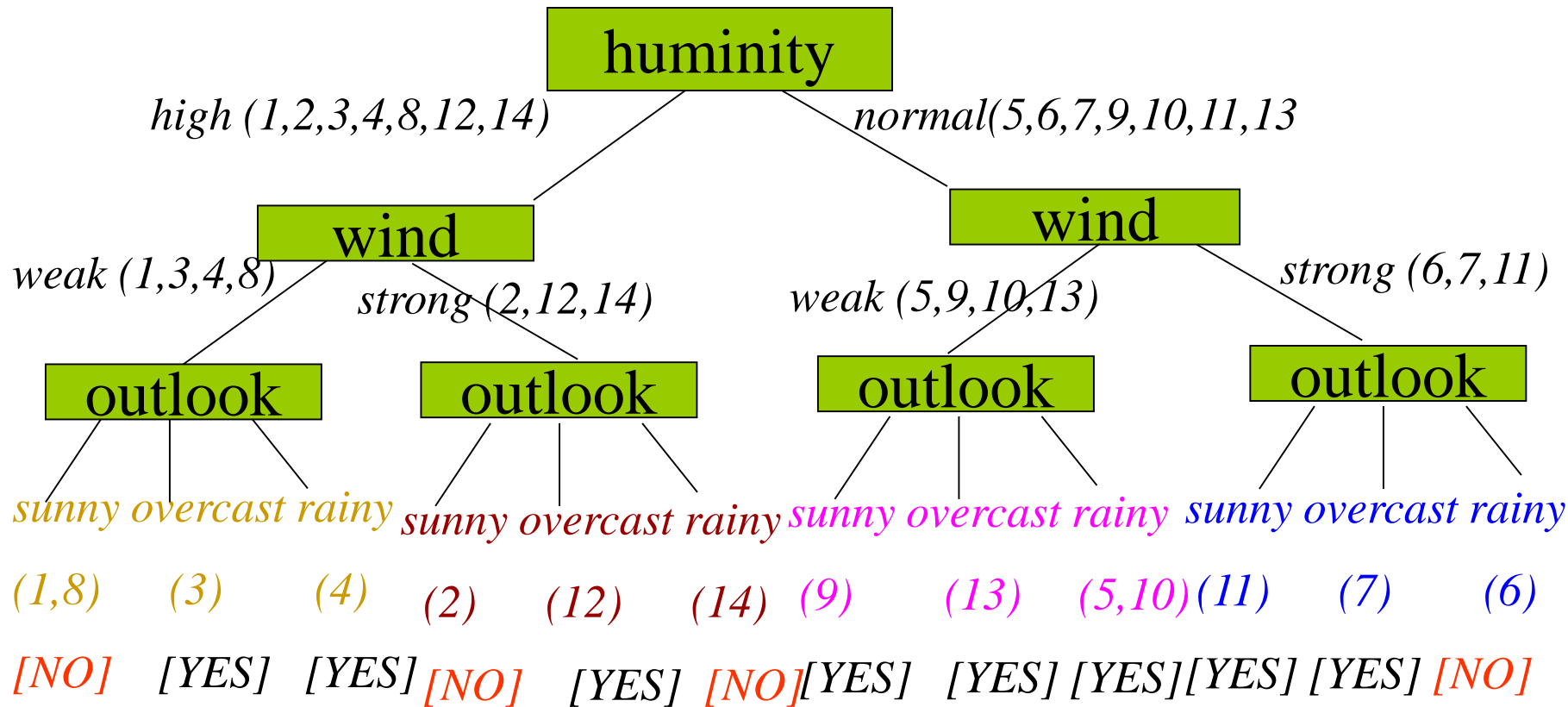
- **Bellek Tabanlı Sınıflandırma Algoritmaları**

K-EnYakın Komşu

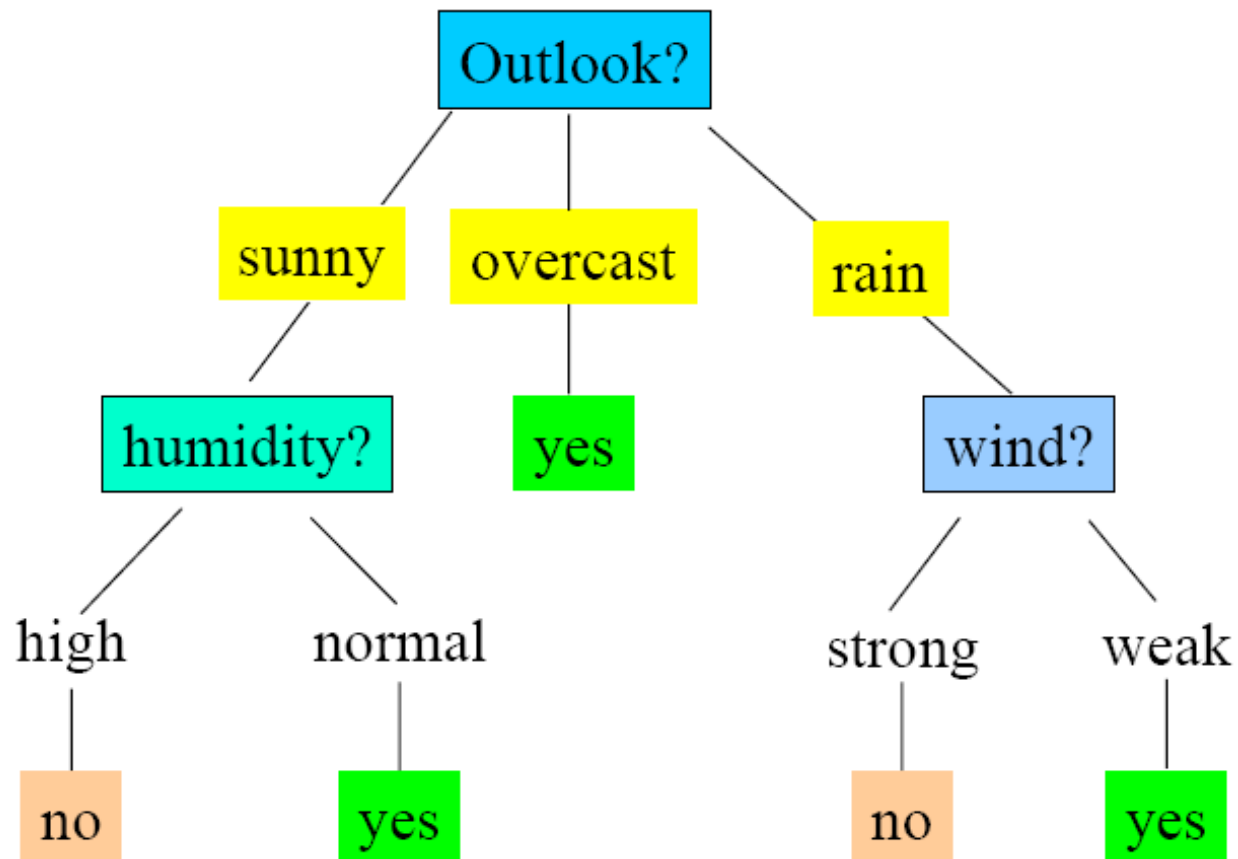
Örnek Karar Ağacı : Play Tennis ? (hava tenis oynamaya uygun mu?)

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no

Özelliklerden biri kök seçilir (outlook, temp, humidity, wind)



veya



Karar Ağacı iyi bir çözümdür

ancak

optimum değildir

optimum bir karar ağacının oluşturulması
için **bir kuralın olması gerekir**

Bilgi kazancı ölçümü: Entropi

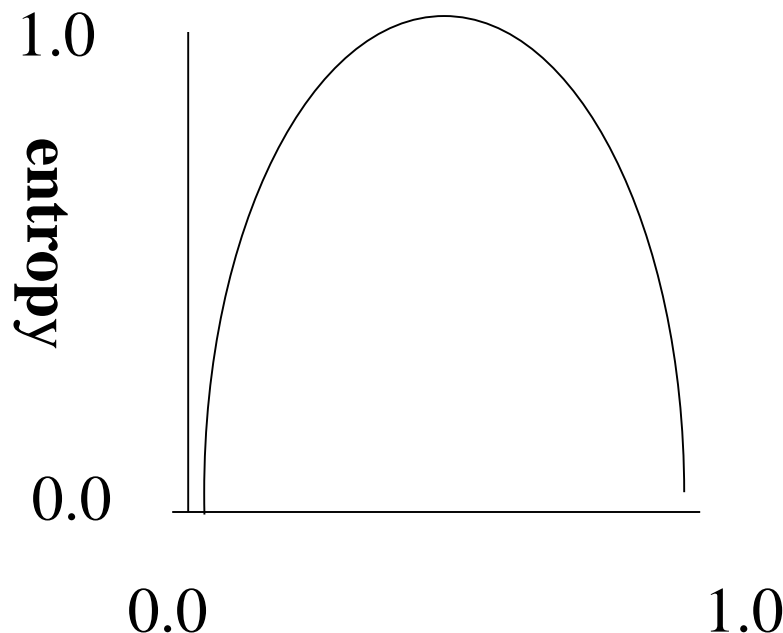
- Entropy rastgeleliğin, belirsizliğin ve beklenmeyen durumun ortaya çıkma olasılığını gösterir.
- Sınıflandırmada
 - örneklerin tümü aynı sınıfa ait ise $\text{entropy}=0$
 - örnekler sınıflar arasında eşit dağılmış ise $\text{entropi}=1$
 - örnekler sınıflar arasında rastgele dağılmış ise $0<\text{entropi}<1$

Information Gain (Bilgi Kazancı- maksimum kazanç)

$$\text{Entropy}(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Bütün örnekler aynı sınıfa ait ise $E(S)=0$ (homojen)

Bütün örnekler sınıflara eşit dağılmış ise $E(S)=1$ (heterojen)



Karar Ağacı oluşturma algoritması (ID 3)

Adım:1

- Karar ağacının hangi kararı alacağı belirlenir.
 - Örnek veri setinde tenis oynamaya gidilip gidilmeyeceğine (play tennis) karar verilecektir.

Sistemin Entropy si hesaplanır

Adım:2

$$E(s_1, s_2, \dots, s_m) = - \sum_{i=1}^m \frac{s_i}{S} \log_2 \frac{s_i}{S} \quad m \text{ sınıf sayısı}$$

14 tane örnek

9 tane YES

5 tane NO

$$\text{Entropy}(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

$$\text{Entropy}(S) \equiv -(9/14) \log_2 (9/14) - (5/14) \log_2 (5/14)$$

$$\text{Entropy}(S) \equiv 0.940$$

Adım :3

- Ağacın en üstünde yani kökte/root ta konumlanacak nitelik (özellik) belirlenir.

Neye göre belirlenir?

- Bilgi kazancı (information gain) en yüksek olan **ÖZELLİK** ağacın en üstünde konumlandırılır.

Adım : 4

- ▣ Bilgi kazancı (information gain) nasıl hesaplanır?

A özelliğinin, S örneği için kazancı (information gain)

$$\text{Gain}(S,A) \equiv \text{Entropy}(S) - \sum P(v) \text{Entropy}(S(v))$$

v: Values of A

$$P(v) \equiv |S(v)| / |S|$$

S:[9+,5-]

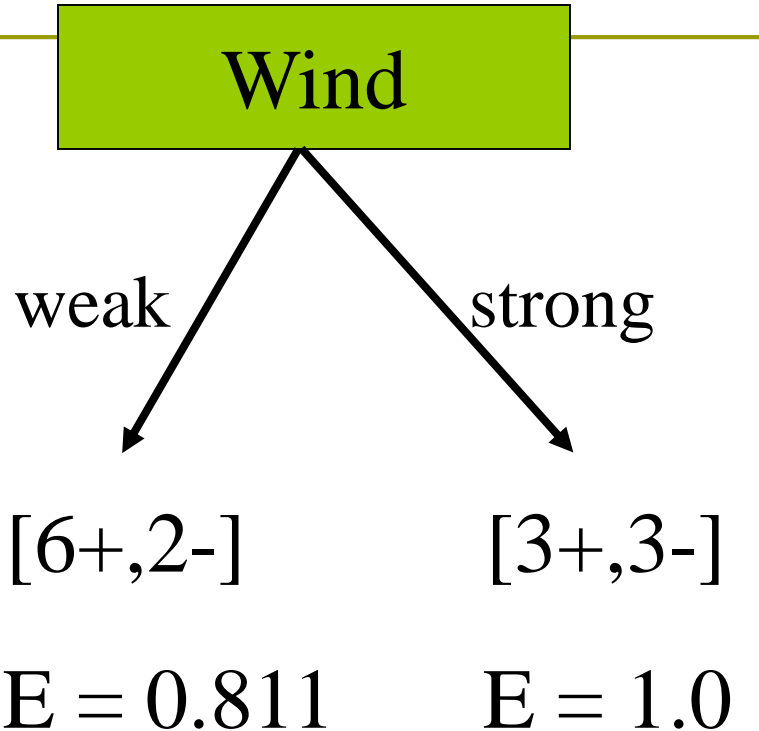
E = 0.940

Gain(S,wind) = ?

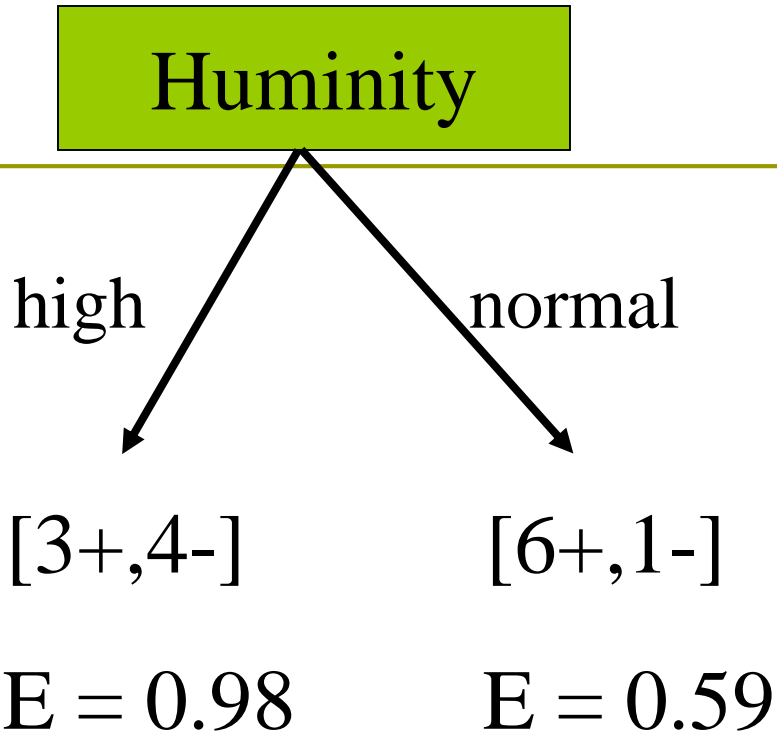
Gain(S,humidity) = ?

Gain(S,temperature) = ?

Gain(S,outlook) = ?

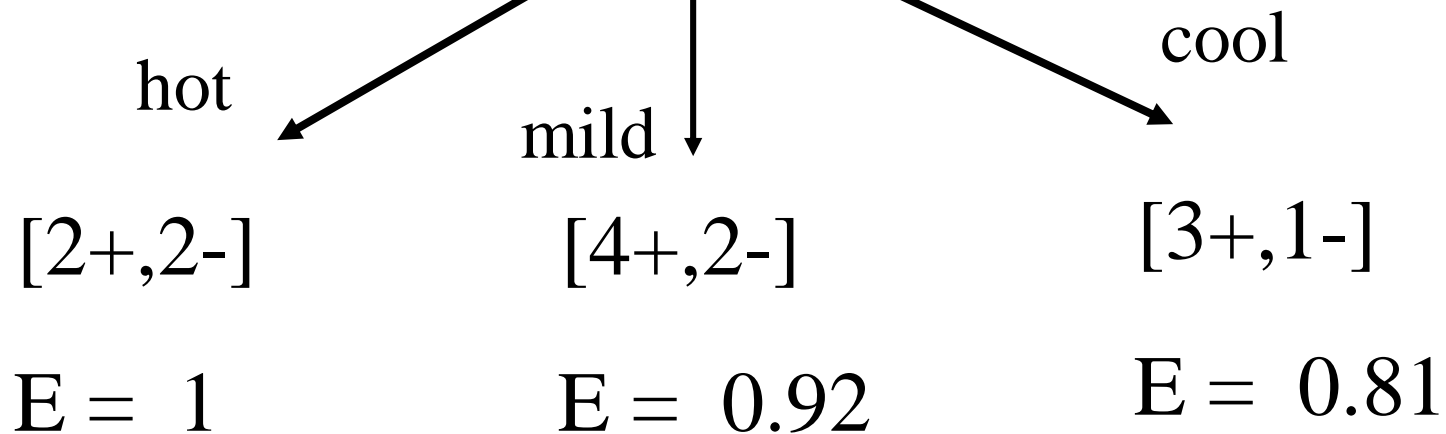


$$\begin{aligned}\text{Gain}(S,\text{wind}) &= 0.940 - [(8/14)[-(6/8)\log_2(6/8) - (2/8)\log_2(2/8)] \\ &\quad - [(6/14)[-(3/6)\log_2(3/6) - (3/6)\log_2(3/6)] \\ &= \mathbf{0.048}\end{aligned}$$



$$\begin{aligned}\text{Gain}(S, \text{humidity}) &= 0.940 - [(7/14)[-(3/7)\log_2(3/7) - (4/7)\log_2(4/7)]] \\ &\quad - [(7/14)[-(6/7)\log_2(6/7) - (1/7)\log_2(1/7)]] \\ &= 0.15\end{aligned}$$

Temperature



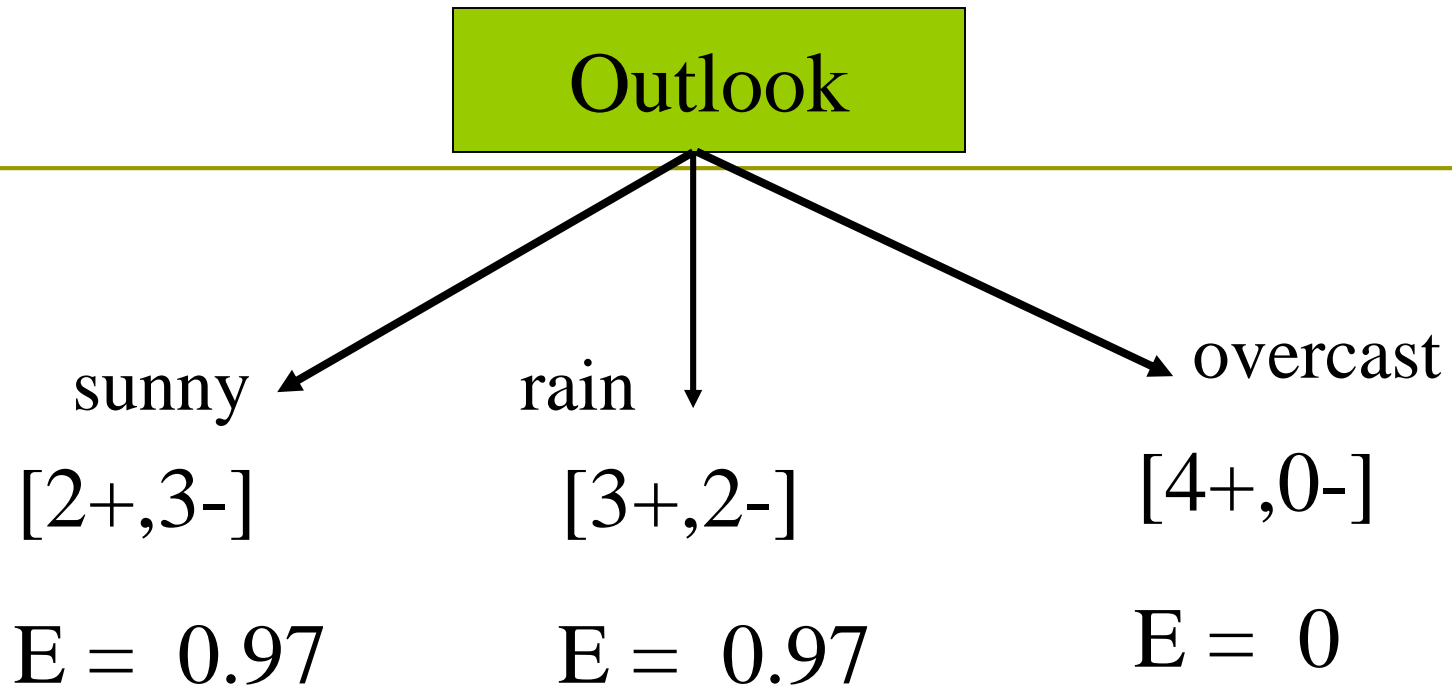
$$\text{Gain}(S, \text{temperature}) = 0.940$$

$$- [(4/14)[-(2/4)\log_2(2/4) - (2/4)\log_2(2/4)]]$$

$$- [(6/14)[-(4/6)\log_2(4/6) - (2/6)\log_2(2/6)]]$$

$$- [(4/14)[-(3/4)\log_2(3/4) - (1/4)\log_2(1/4)]]$$

$$= 0.027$$



Gain(S,temperature) = 0.940

$$\begin{aligned} & - [(5/14)[-(2/5)\log_2(2/5) - (3/5)\log_2(3/5)]] \\ & - [(5/14)[-(3/5)\log_2(3/5) - (2/5)\log_2(2/5)]] \\ & - [(4/14)[-(4/4)\log_2(4/4) - (0/4)\log_2(0/4)]] \\ & = \mathbf{0.246} \end{aligned}$$

$$\text{Gain}(S, \text{wind}) = 0.048$$

$$\text{Gain}(S, \text{humidity}) = 0.15$$

$$\text{Gain}(S, \text{temperature}) = 0.027$$

$$\text{Gain}(S, \text{outlook}) = 0.246$$

outlook

sunny
(1,2,8,9,11)
[+2,-3]
?

overcast
(3,7,12,13)
[+4,-0]
YES

rainy
(4,5,6,10,14)
[+3,-2]
?

$$S_{\text{sunny}} = [+2, -3]$$

$$E(\text{sunny}) = -(2/5)\log_2(2/5) - (3/5)\log_2(3/5) = 0.97$$

$$\text{Gain}(S_{\text{sunny}, \text{humidity}}) = ?$$

$$\text{Gain}(S_{\text{sunny}, \text{temp}}) = ?$$

$$\text{Gain}(S_{\text{sunny}, \text{wind}}) = ?$$

$$\text{Gain}(S_{\text{sunny,humidity}}) = 0.97$$

$$-(2/5)[-(2/2)\log_2(2/2)-(0/2)\log_2(0/2)]$$

$$-(3/5)[-(0/3)\log_2(0/3)-(3/3)\log_2(3/3)]$$

$$= 0.97$$

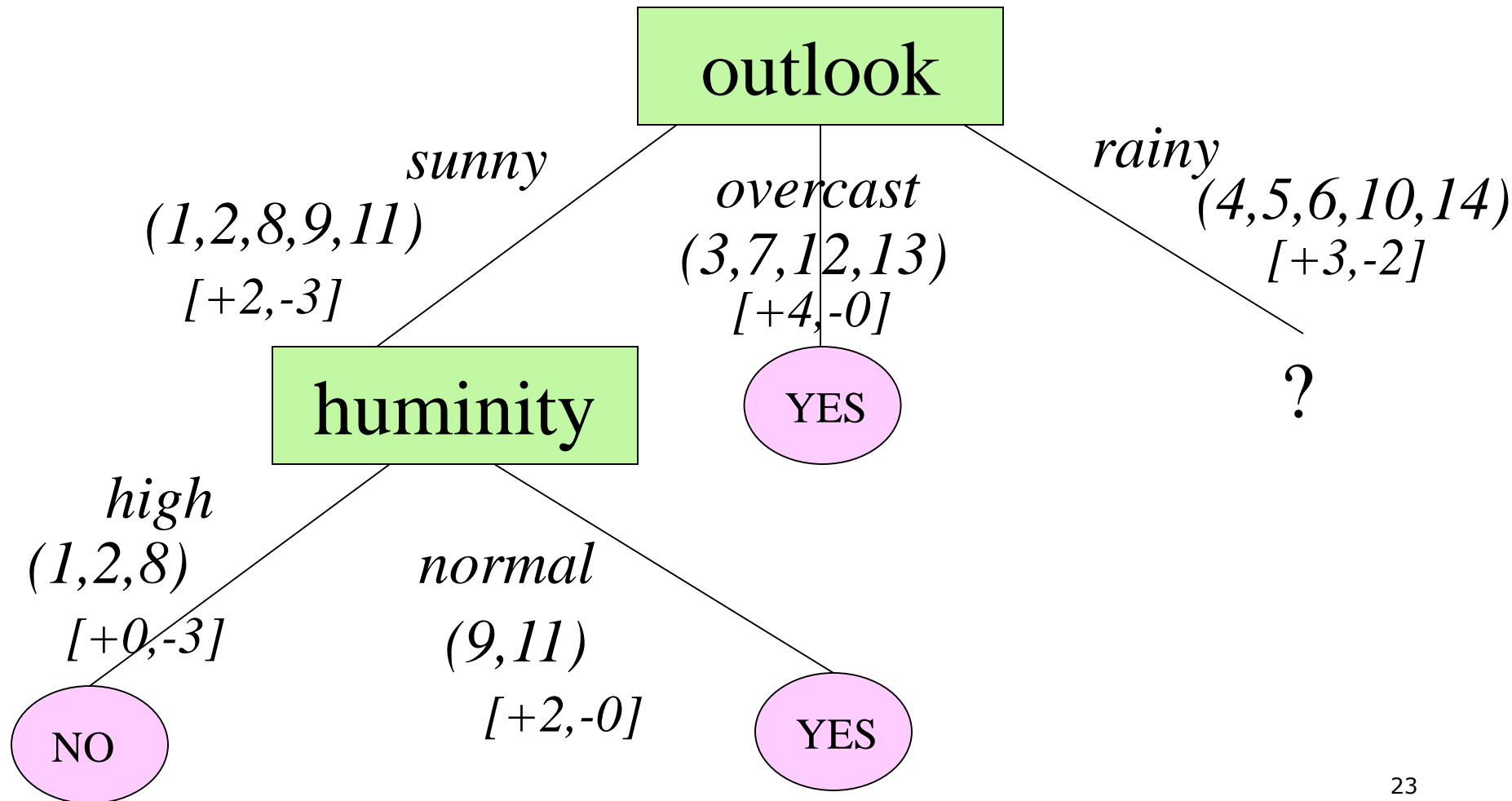
$$\text{Gain}(S_{\text{sunny,wind}}) = 0.97$$

$$-(3/5)[-(1/3)\log_2(1/3)-(2/3)\log_2(2/3)]$$

$$-(2/5)[-(1/2)\log_2(1/2)-(1/2)\log_2(1/2)]$$

$$= 0.019$$

$$\text{Gain}(S_{\text{sunny,temp}}) = 0.57$$



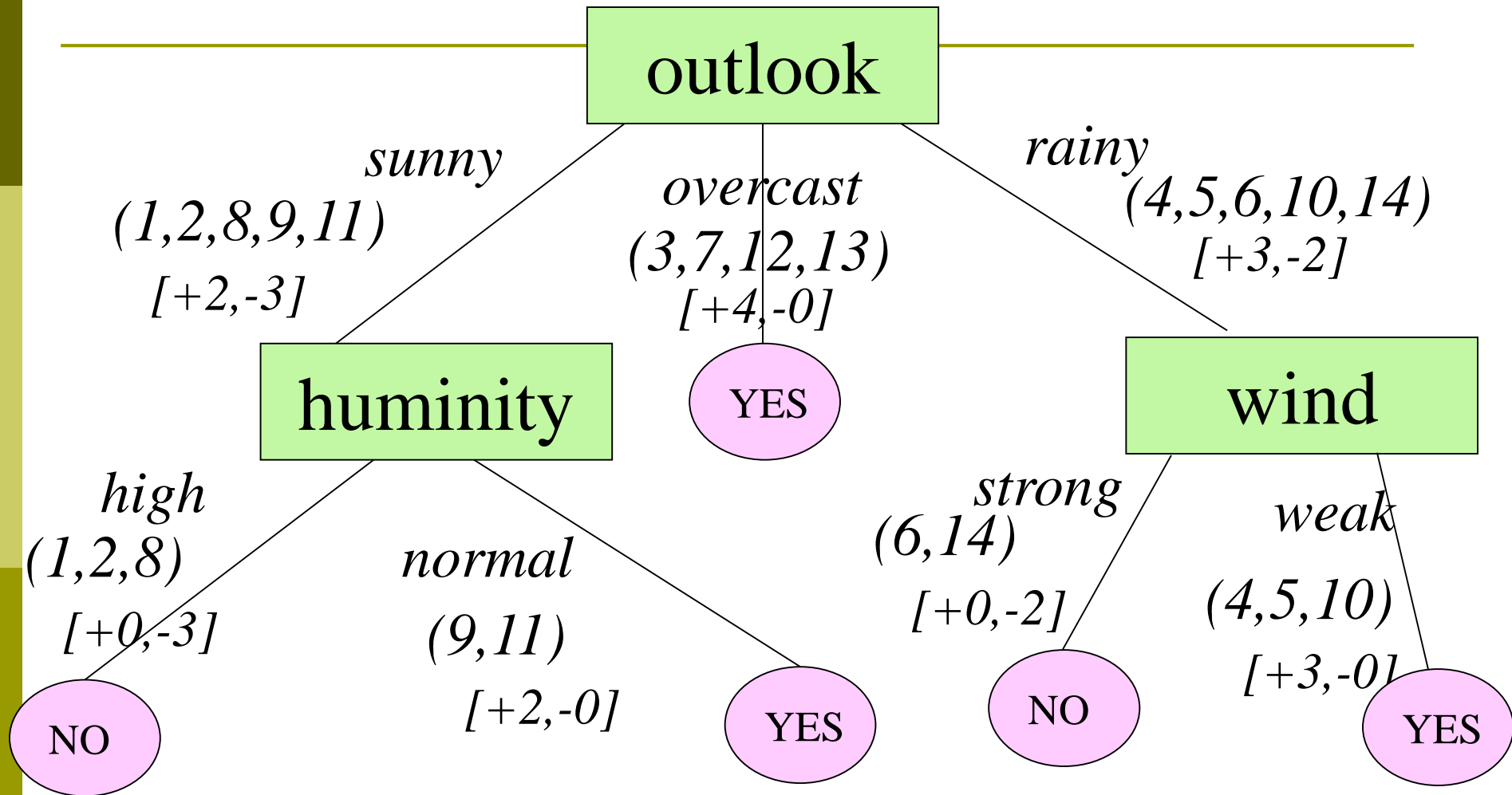
Aynı işlem

$$\text{Gain}(S_{\text{rainy,humidity}}) = ?$$

$$\text{Gain}(S_{\text{rainy,temp}}) = ?$$

$$\text{Gain}(S_{\text{rainy,wind}}) = ?$$

bulmak için yapılır.



Karar Ağacı kullanarak sınıflandırma

□ Avantajları:

- Karar ağacı oluşturmak zahmetsizdir
- Küçük ağaçları yorumlamak kolaydır
- Anlaşılabilir kurallar oluşturulabilir
- Sürekli ve ayrık nitelik değerleri için kullanılabilir

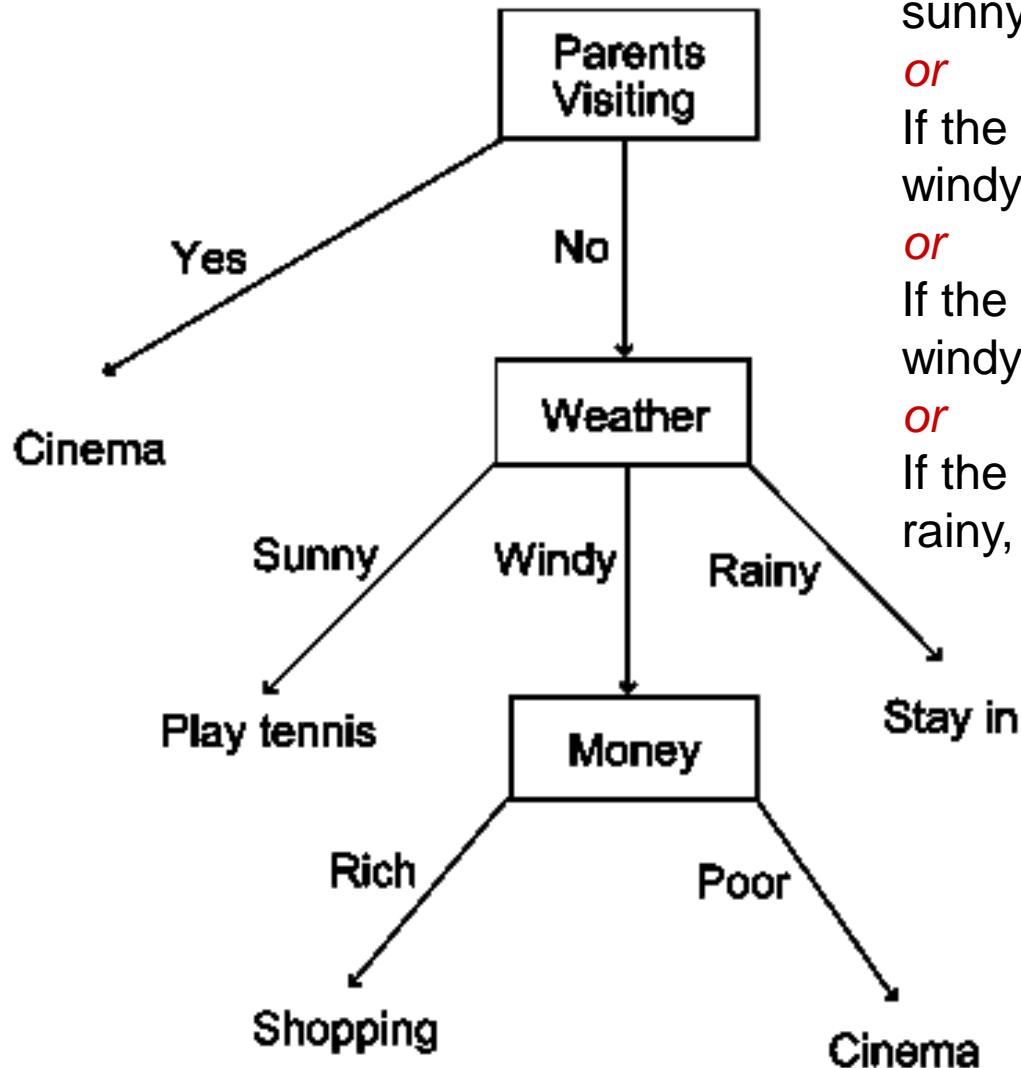
Karar Ağacı kullanarak sınıflandırma

■ Dezavantajları:

- Sürekli nitelik değerlerini tahmin etmekte çok başarılı değildir
- Sınıf sayısı fazla ve öğrenme kümesi örnekleri sayısı az olduğunda model oluşturma çok başarılı değildir
- Zaman ve yer karmaşıklığı öğrenme kümesi örnekleri sayısına, nitelik sayısına ve oluşan ağacın yapısına bağlıdır
- Hem ağaç oluşturma karmaşıklığı hem de ağaç budama karmaşıklığı fazladır

Örnek: Karar Ağacı oluşturma

weekend	weather	parent	money	decision
w1	sunny	yes	rich	cinema
w2	sunny	no	rich	tennis
w3	windy	yes	rich	cinema
w4	rainy	yes	poor	cinema
w5	rainy	no	rich	stay in
w6	rainy	yes	poor	cinema
w7	windy	no	poor	cinema
w8	windy	no	rich	shopping
w9	windy	yes	rich	cinema
w10	sunny	no	rich	tennis



If the parents are visiting, then go to the cinema

or

If the parents are not visiting *and* it is sunny, then play tennis

or

If the parents are not visiting *and* it is windy *and* you're rich, then go shopping

or

If the parents are not visiting *and* it is windy *and* you're poor, then go to cinema

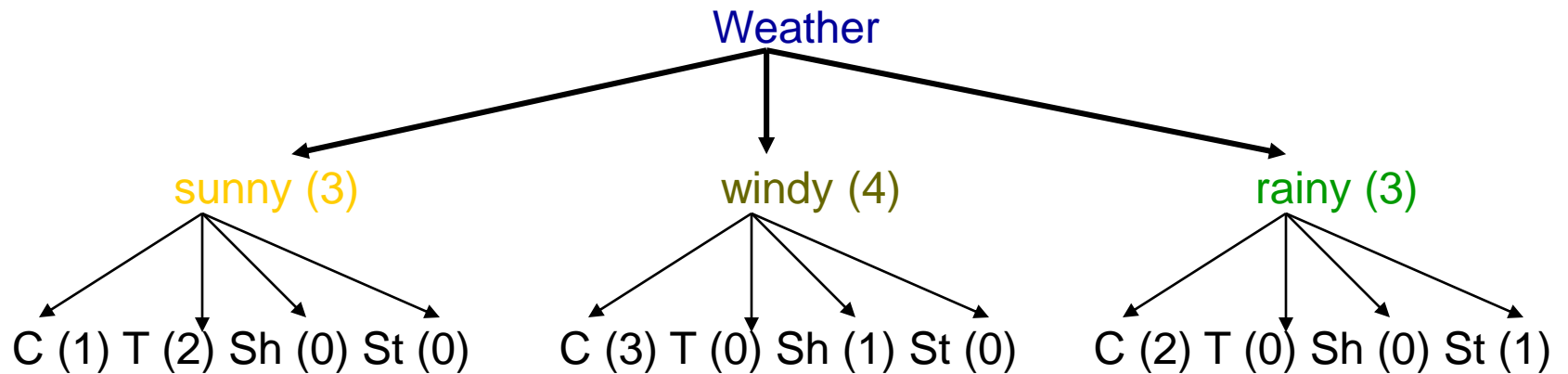
or

If the parents are not visiting *and* it is rainy, then stay in.

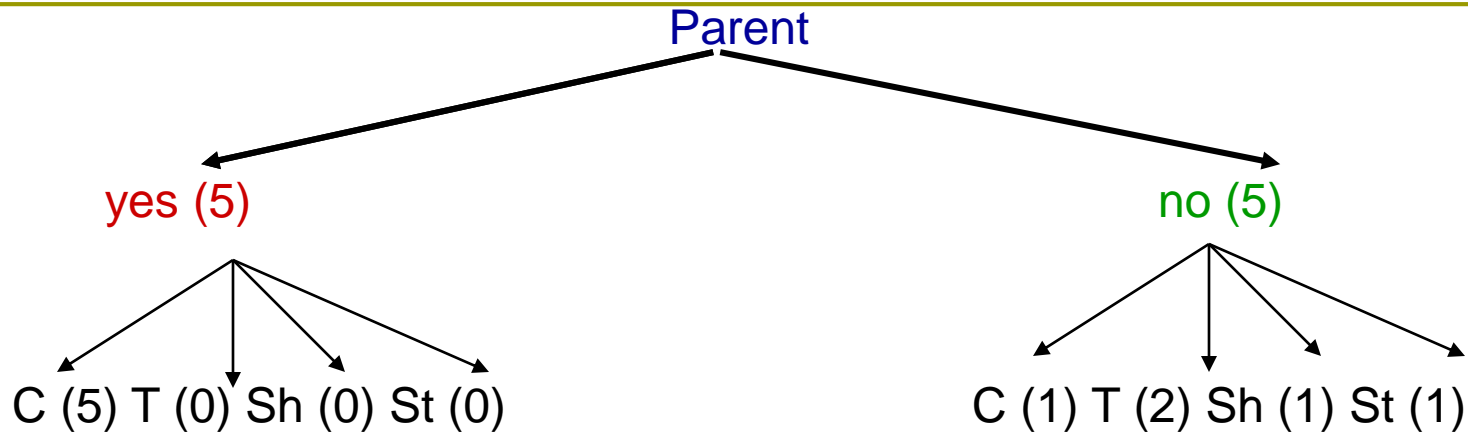
Önce sistemin Entropy si hesaplanır.

$$\begin{aligned}\text{Entropy}(S) &= -p_{\text{cinema}} \log_2(p_{\text{cinema}}) - p_{\text{tennis}} \log_2(p_{\text{tennis}}) - p_{\text{shopping}} \log_2(p_{\text{shopping}}) - p_{\text{stay_in}} \log_2(p_{\text{stay_in}}) \\ &= -(6/10) * \log_2(6/10) - (2/10) * \log_2(2/10) - (1/10) * \log_2(1/10) - (1/10) * \log_2(1/10) \\ &= -(6/10) * -0.737 - (2/10) * -2.322 - (1/10) * -3.322 - (1/10) * -3.322 \\ &= 0.4422 + 0.4644 + 0.3322 + 0.3322 = \mathbf{1.571}\end{aligned}$$

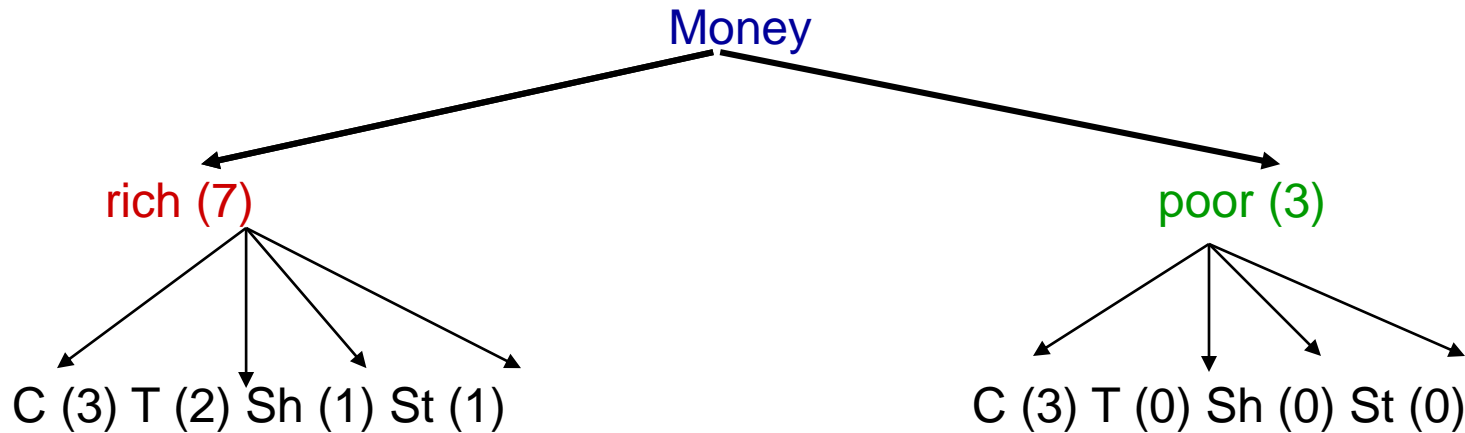
Köke yerleştirilecek özellik belirlenir (weather, parents, money ?)



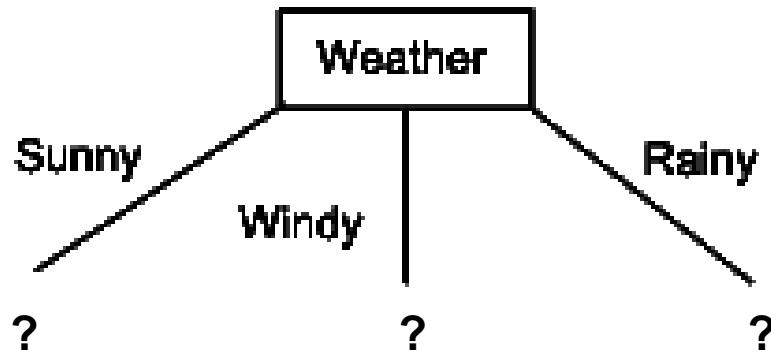
$$\begin{aligned}\text{Gain}(S, \text{weather}) &= 1.571 - [(3/10)[-(1/3)\log_2(1/3) - (2/3)\log_2(2/3)]] \\ &\quad - [(4/10)[-(3/4)\log_2(3/4) - (1/4)\log_2(1/4)]] \\ &\quad - [(3/10)[-(1/3)\log_2(1/3) - (2/3)\log_2(2/3)]] \\ &= \mathbf{0.70}\end{aligned}$$



$$\begin{aligned}\text{Gain}(S, \text{parent}) &= 1.571 - [(5/10)[- (1/5)\log_2(1/5) - (2/5)\log_2(2/5) - (1/5)\log_2(1/5) - (1/5)\log_2(1/5)]] \\ &\quad - [(5/10)[- (5/5)\log_2(5/5) - 0]] \\ &= 0.61\end{aligned}$$



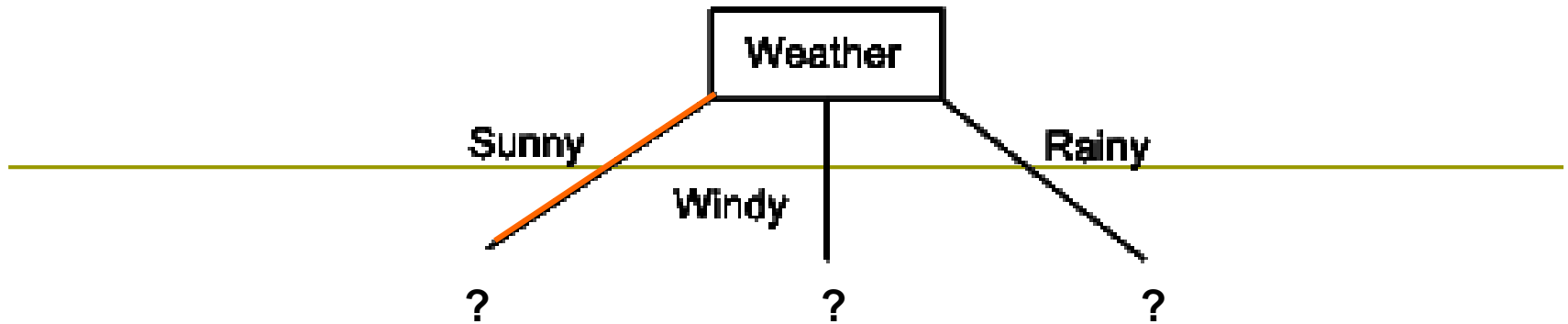
$$\begin{aligned}
 \text{Gain}(\text{S}, \text{money}) &= 1.571 - [(7/10)[- (3/7)\log_2(3/7) - (2/7)\log_2(2/7) - (1/7)\log_2(1/7) - (1/7)\log_2(1/7)]] \\
 &\quad - [(3/10)[- (3/3)\log_2(3/3) - 0]] \\
 &= \mathbf{0.2816}
 \end{aligned}$$



$S_{\text{sunny}} = \{W1, W2, W10\}$. W1, W2 ve W10, sırasıyla Cinema, Tennis ve Tennis.

$S_{\text{windy}} = \{W3, W7, W8, W9\}$. W3, W7, W8 ve W9 sırasıyla
Cinema, Cinema, Shopping ve Cinema.

$S_{\text{rainy}} = \{W4, W5, W6\}$. W4, W5 ve W6 sırasıyla Cinema, Stay in ve Cinema.



Sunny'nin altına gelecek özelliği belirlemek için $\text{Gain}(S_{\text{sunny, parents}})$ ve $\text{Gain}(S_{\text{sunny, money}})$ hesaplanmalıdır.

Önce, $\text{Entropy}(S_{\text{sunny}})$ değeri bilinmelidir.

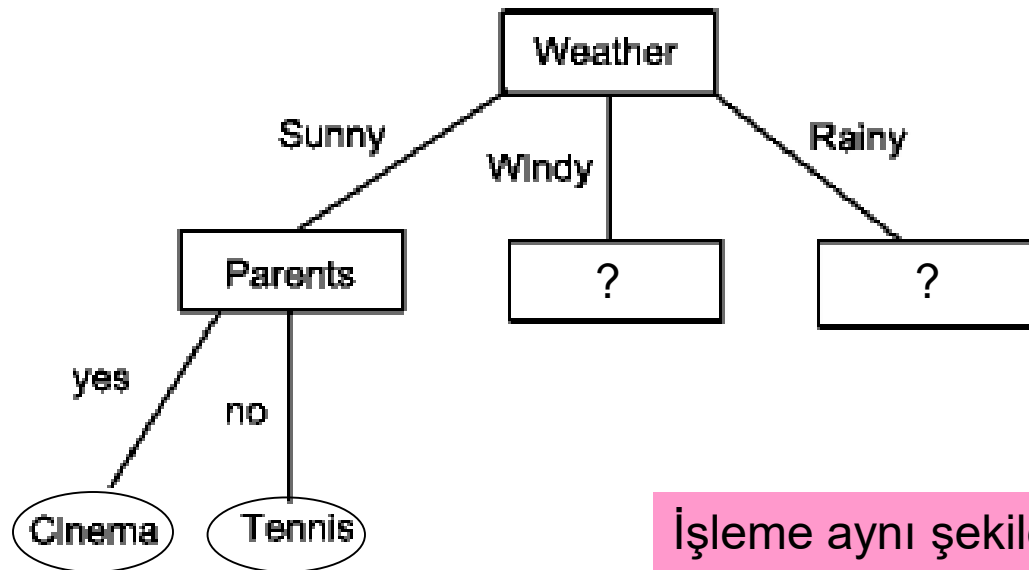
Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W10	Sunny	No	Rich	Tennis

$$E(\text{sunny}) = -(1/3)\log_2(1/3) - (2/3)\log_2(2/3) = \mathbf{0.918}$$

Sunny'nin altına gelecek özelliği belirlemek için
 $\text{Gain}(S_{\text{sunny, parents}})$ ve $\text{Gain}(S_{\text{sunny, money}})$?

$$\text{Gain}(S_{\text{sunny, parents}}) = 0.918 - (1/3)[-(1/1)\log_2(1/1)] - (2/3)[-(2/2)\log_2(2/2)] = 0.918$$

$$\text{Gain}(S_{\text{sunny, money}}) = 0.918 - (3/3)[-(1/3)\log_2(1/3) - (2/3)\log_2(2/3)] = 0$$



İşleme aynı şekilde devam edilir

C4.5 Algoritması

- Quinlan'ın ID3 algoritması yine aynı kişi tarafından genişletilerek C4.5 adını almıştır.
- ID3 algoritmasında bir özellik sayısal değerlere sahip ise sonuç alınamamaktadır. Bu yüzden C4.5 algoritması geliştirilmiştir.
- Sayısal değerler ile çalışılırken bir eşik değeri belirlenir. Bu eşik değeri bulunurken özelliğin değerleri sıralanır ve $[v_i, v_{i+1}]$ aralığının orta noktası alınır, ve bu değer t eşik değeri olarak belirlenir. Ve özellik değeri bu t eşik değerinden büyük veya küçük eşit olmak üzere ikiye ayrılır.

Bilinmeyen Nitelik Değeri

- C4.5 kayıp verilere sahip örneklerde bir düzeltme faktörü kullanır.
- $H(X)$ ve $H(X,T)$ değerleri hesaplanırken, bilinen niteliklere sahip örnekler alınmıştır.
- F faktörü kullanılarak kazanç ölçütü düzeltilir.

$F = \text{Veri tabanında değeri bilinen niteliğe sahip örneklerin sayısı} / \text{Veri tabanındaki tüm örneklerin sayısı}$

$\text{Kazanç}(X) = F(H(T) - H(X,T))$