



Yalova Üniversitesi  
Fen Bilimleri Enstitüsü  
Bilgisayar Mühendisliği Ana Bilim Dalı

Oruntü Tanıma Dersi

---

Sevdanur GENC - 135105008

---

**Oruntu Tanima - Decision Trees**

<b>Icindekiler</b>	<b>Sayfa</b>
1. Karar Agaclari	3
2. Entropy	3
2.1. Entropy - Ornek I	3
2.2. Entropy - Ornek II	3
3. Information Gain - Bilgi Kazanci	4
4. Gini Indexleme	4
4.1. Gini Indexleme - Ornek I	4
4.2. Gini Indexleme - Ornek II	4
5. CART - Calssification And Regression Tree	5
5.1. Siniflandirma Agaci (Classification Tree)	5
5.2. Regresyon Agaci (Regression Tree)	5
5.3. Regresyon Analizi	5
6. ID3 (Iterative Dichotomiser 3 - Tekrarlanan ikili yapi) Algoritmasi	6
7. C4.5 Algoritmasi	7
8. ID3 Algoritma Uygulamasi	8
9. C4.5 Algoritma Uygulamasi	15
10. Matlab Karar Agaci Uygulamalari	17
10.1. Uygulama I	17
10.2. Uygulama II	18
10.3. Uygulama III	19

## Karar Agaclari

Pek çok problemin çözümü için veriler üzerinde istatistik analizleri önemli bir yöntemdir. Bazı durumlarda istatistik kullanımı sınırlı olduğundan dolayı akıllı veri analizi yöntemlerini içeren algoritmalar ortaya çıkmıştır. Ancak bu yöntemlerin bazıları zayıftır. Bu zayıflığın söz konusu olmayan bir yaklaşım türü ise karar ağaçlarıdır. Karar ağaçları, hedef sonuçlarının yaklaşık değerlerini hesaplamak için kullanılan ve öğrenme verilerinin karar ağacı ile gösterildiği bir yöntemdir.

Karar ağaçları için kullanılacak karar ağacı algoritması iki aşamadan oluşmaktadır;

1. Ağacın Oluşturma ; Butun öğrenme verilerine sahip olunan bir kümedir.
2. Ağacın Budama ; Butun öğrenme verilerinin sahip olduğu kümeyle ait olan ve test kümesinde hataya sebep olan dalların ağaçtan budanması, silinmesidir.

Tüm bunlara baktığımızda, aslında kesin bir optimum sonuç vermeyecektir. Bunun için, gerekli optimum sonuca yaklaşılmada kullanılan kurallar bulunmaktadır.

Karar Ağacı Algoritmasının Sahip Olduğu Adımlar;

1. Karar ağacının sonuç olarak hangi kararı alacağı belirlenir.
2. Kurulacak sistemin Entropy'si hesaplanır.
3. Ağacın Root'u yani kökü belirlenir. Bunu belirleyebilmek için Information Gain hesaplanır. Yani bilgi kazancı, en yüksek olan ağacın en üstteki yerini alır.

**Entropy :** Rastgeleliğin, belirsizliğin ve beklenmeyen durumun ortaya çıkma olasılığını gösterir.

$S'$  i bir dataset olarak kabul edelim.  $S$  dataseti içerisinde bulunan örnekler aynı sınıfa ait ise entropy değeri 0, örnekler eşit dağılmışsa entropy 1 değerine yakın olacaktır. Örnekler sınıflar arasında rastgele dağılmışsa  $0 < \text{entropy} < 1$  değeri beklenmektedir.

$$\text{Entropy}(S) \equiv -p_i \log_2 p_i - p_i \log_2 p_i$$

### Entropy - Örnek I

Y	X1	X2	X3
Pi	3/6	2/6	1/6

Entropy belirsizliği hesaplanacak olursa;

$$E(S) = - (X1 \log_2 X1 + X2 \log_2 X2 + X3 \log_2 X3)$$

$$E(S) = - (3/6 \log_2 3/6 + 2/6 \log_2 2/6 + 1/6 \log_2 1/6)$$

$$E(S) = 1,4591$$

### Entropy - Örnek II

Y	X1	X2	X3	X4	X5	X6	X7	X8
Pi	Evet	Evet	Hayir	Hayir	Hayir	Hayir	Hayir	Hayir

Entropy belirsizliği hesaplanacak olursa;

$S_{uzayı} = \{\text{Evet}, \text{Evet}, \text{Hayir}, \text{Hayir}, \text{Hayir}, \text{Hayir}, \text{Hayir}, \text{Hayir}\}$

$P(\text{Evet}) = \text{Evet} / S$  ve  $P(\text{Hayir}) = \text{Hayir} / S$  şeklinde olasılıkları hesaplanır.

$$P(\text{Evet}) = 2/8 = 0.25 \text{ ve } P(\text{Hayir}) = 6/8 = 0.75$$

$$E(S) = - (P(\text{Evet}) \log_2 P(\text{Evet}) + P(\text{Hayir}) \log_2 P(\text{Hayir}))$$

$$E(S) = - (0.25 \log_2 0.25 + 0.75 \log_2 0.75)$$

$$E(S) = 0.97$$

**Information Gain (Bilgi Kazanci) :** Karar agaci yontemlerinde en ayirt edici ozelligi belirlemek amactir. Bunun icin de, her ozellik icin bilgi kazanci hesaplanır. Bu hesaplama, Entropy hesabimi kullanilmaktadir.

Her ozelligin bilgi kazancinda dogal olarak bolunmelere neden olacaktir. Entropy bu bolunmelerin olculerini azaltacagindan dolayi en iyi azaltmayi saglayan bolunme basarili olarak secilir.

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

**Gini Indexleme :** Dugum homojenliginin olcumunde kullanılan bir yontemdir. Kayitlarin butun siniflar arasinda esit olarak dagilmasiyla ilgilenir.

Bir t dugumundeki j sinifina ait bagil olasilik hesaplama ve buna bagli olarak gini tabanlı bolunme formulleri soyledir (ni = child kayitlari);

$$GINI(t) = 1 - \sum_j [p(j|t)]^2 \quad \text{ve} \quad GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i) \quad \text{formulleri kullanilir.}$$

### Gini Indexleme - Ornek I

Tek Oznitelikli C1 ve C2 isminde iki sinif olsun. Bu siniflardaki verilerin homojen dagilimlari sonucunda toplamda 6 ornek bulunmaktadır. Siniftaki orneklerin sayisi esit olana kadar dagilim gercekleşmektedir.

C1	C2	C1	C2	C1	C2	C1	C2
0	6	1	5	2	4	3	3

$$P(c1) = 0/6$$

$$P(c2) = 6/6$$

$$Gini = 1 - P(c1)^2 - P(c2)^2$$

$$Gini = 1 - 0^2 - 1^2$$

$$Gini = 0$$

$$P(c1) = 1/6$$

$$P(c2) = 5/6$$

$$Gini = 1 - P(c1)^2 - P(c2)^2$$

$$Gini = 1 - (1/6)^2 - (5/6)^2$$

$$Gini = 0.278$$

$$P(c1) = 2/6$$

$$P(c2) = 4/6$$

$$Gini = 1 - P(c1)^2 - P(c2)^2$$

$$Gini = 1 - (2/6)^2 - (4/6)^2$$

$$Gini = 0.444$$

$$P(c1) = 3/6$$

$$P(c2) = 3/6$$

$$Gini = 1 - P(c1)^2 - P(c2)^2$$

$$Gini = 1 - (3/6)^2 - (3/6)^2$$

$$Gini = 0.500$$

### Gini Indexleme - Ornek II

Kayitlarda ikiser oznitelik oldugu dusunulurse (Asagidaki formuller kullanilir);

$$GINI(t) = 1 - \sum_j [p(j|t)]^2$$

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

	N1	N2	Parent
C1	5	1	6
C2	2	4	6
Children	7	5	

$$Gini(N1) = 1 - P(c1)^2 - P(c2)^2$$

$$Gini(N1) = 1 - (5/7)^2 - (2/7)^2$$

$$Gini(N1) = 0.408$$

$$Gini(N2) = 1 - P(c1)^2 - P(c2)^2$$

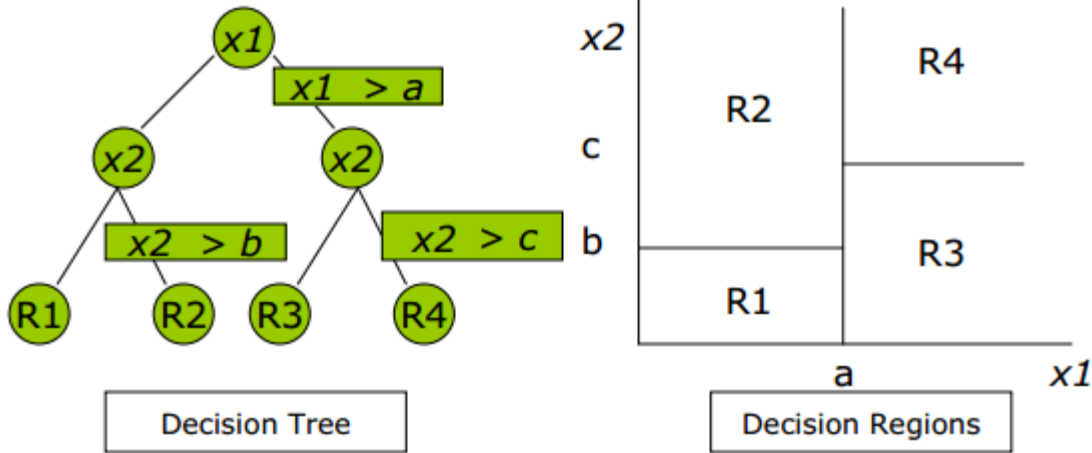
$$Gini(N2) = 1 - (1/5)^2 - (4/5)^2$$

$$Gini(N2) = 0.320$$

$$Gini(Children) = \sum (Children / Parent) * Gini(N)$$

$$Gini(Children) = 7/12 * 0.408 + 5/12 * 0.320$$

$$Gini(Children) = 0.371$$



### CART (Classification & Regression Trees)

Sınıflandırma ve regresyon ağaçları olarak bilinen bu modelin temelinde tek değişkenli ikili kararların bir hiyerarsisini içermektedir. Yaygın olarak kullanılan bu istatistiksel prosedür, verileri iki alt kümeye ayırmaktadır. Her bir alt küme içerisindeki veriler bir önceki alt kümeye ait verilerden biraz daha fazla homojen bir yapıya sahip olmaktadır. Birbirini devam eden bu süreçler en optimize edilmiş haldeki homojenlik kriterini veya durma koşullarını sağlayıncaya kadar kendini sürekli tekrar edecektir. Tüm bu süreçlerde en iyi seçeneklerin seçilmesine özen gösterilmektedir buna bağlı olarak bölünme kriterleri kullanılmaktadır. Kısaca özetleyecek olursak, CART ağaçları kesin bir heterojenliğe (impurity) sahiptirler ve bu heterojenlik iki değerli (binary) ağaçlar yardımıyla optimize edilerek homojen hale getirilmektedir. Hedef, aynı veya yakın sonuç çıktı değerlerinin olduğu alt gruplar yaratılmasıdır.

CART algoritmalarına örnek verecek olursak; Twoing ve Gini algoritmaları.

CART ağacını oluştururken en iyi dallara ayırma kriterini seçmek için Entropy'den faydalanılmaktadır. Bu kriteri en iyi şekilde sonuçlandırmak için ise kullanılan formül ;

$$\Psi\left(\frac{s}{t}\right) = 2 P_L P_R \sum_{j=1}^M \left| P\left(\frac{C_j}{t_L}\right) - P\left(\frac{C_j}{t_R}\right) \right|$$

$W(s/t)$  : Herhangi bir  $t$  düğümündeki  $s$  dalları

$t$  : Dalların yapılacağı düğümler

$C$  : Kriteri

$L$  : Ağacın sol yanı

$R$  : Ağacın sağ yanı

$P_L$  ve  $P_R$  : Eğitim seti içerisindeki bir verinin ağacın solunda ya da sağında olma olasılığı

$P(C_j / t_L)$  ve  $P(C_j / t_R)$  : Verilerin bulunduğu  $C_j$  sınıfındaki bir kaydın ağactaki yerinin solunda ya da sağında olma olasılığı.

Bu formüle dayanarak söylenebilecek kural ;

- Dalların en büyük kriteri göre gerçekleştiriliyorsa Twoing algoritması, en küçük kriterlere göre gerçekleştiriliyorsa Gini algoritmasının kullanılması tavsiye ediliyor. Gini algoritmasındaki amaç; her zaman her adımda en büyük veri kümelemesinin oluşturulmasıdır. Bu kümelemeler sonuçlandığında ilgilenilmeyen dallar budanabilir. Twoing algoritmasındaki amaç ise; her zaman ana düğüm ve yavru düğümlerin çoğunluğunun yarısı üzerinde çalışma hedefidir. Gini algoritmasına göre daha yavaş çalışacak ve veriler üzerinde daha dengeli bir tavir sergilemiş olacaktır.

CART Aaclari uzerinde calisirken minimum sayidaki n dugumu belirlenir. n dugumunun sayisini belirlerken genellikle veri setinin yuzde 10'u kadar bir deger secilir. Aksi bir degerin secilmesi algoritmayi ya hizlandirir yada yavaslatir ve bu test analiz sonuclarini yanlis degerlendirmis olur.

CART yaklasiminda, siniflar arasi ayrim maksimize edilirken, sinif icerisindeki varyasyonun minimize edilmesi bir kural olarak benimsenmistir. Hem kategorik hem de surekli bagimli degiskenlerin modellenmesi soz konusudur. Bagimli degiskenler eger kategorik ise yontem Siniflandirma Agaci (CT - Classification Tree), surekli ise Regresyon Agaci (Regression Tree) ismini almaktadir.

**Siniflandirma Agaci (CT)** : Siniflandirma agaci genelde turlerin dagilimi modellenmesi icin kullanilmaktadir. Bu sebepten bagimli degisken Var/Yok veya Evet/Hayir gibi ikili kategorileri icermektedir. Ikili bagimli degiskenlerinin homojenligine karar verirken Gini katisiklik olcumu kullanilir. Herhangi bir t dugumu icin g(t) fonksiyonu soyledir;

$$g(t) = \sum_{j \neq i} p(j|t)p(i|t)$$

Buradaki i ve j egitim setindeki hedef (bagimli) degiskenin kategorileridir. Egerki ikili kategorilerden olusan bir yontem kullaniliyorsa formul esitligi asagidaki gibi degisecektir;

$$g(t) = 2p(1|t)p(2|t)$$

Herhangi bir t dugumune gelen bir ornegin s olarak bilinmesi ile, hem sol taraf ayrimini (tl) hem de sag taraf (tg) ayrimini gerceklestirecektir.

$$\phi(s,t) = g(t) - p_L g(t_L) - p_R g(t_R)$$

Burada, t dugumundeki durumlarin oranini belirtirken sag taraftaki Pr, sol taraftaki Pl degerleri belirlenir.

**Regresyon Agaci (RT)** : Regresyon agac mantiginda siniflara yer yoktur. Buna bagli olarak Gini indeksleme de kullanilmaz. Agac olusturulurken ikiye ayrilan sonuclarda dugumlerin tahmini toplam varyansin minimize edilerek hesaplanmasi gerekiyor. Agac olusturulurken her bir dugum icin yapilmasi gereken minimizasyon yani azaltma islemi icin gerekli formul asagidadir;

$$\arg \min_{x_j \leq x_j^R, j=1, \dots, M} [P_L \text{Var}(Y_L) + P_R \text{Var}(Y_R)]$$

Burada yine, Pl ile sol dugum Pr ile sag dugum olasiliklari hesaplanmak istenmistir. Egitim setindeki degiskenlerin sayisini M ile ifade etmistir. Var(Yl) ve Var(Yr) karsilikli sag ve sol alt dugumlerin vektorlerini temsil etmektedir. Artiklari karelerinin azaltma algoritmasina gore asagidaki formulle;

$$i(t) = 1 - \sum_{k=1}^K p^2(k|t)$$

P(k|t) dugumunun t icerisinde bulundu gu sinifin k'nin kosullarina bagli ozelliklerini, K sinif sayisi ve k sinif sayisi indeksi ile t dugum indeksini belirtmistir.

### Regresyon Analizi :

Bir veya birden fazla kullanılan bagimsiz degiskenler ile bagimli degiskenlerin arasindaki iliskiyi kiyaslamak icin Regresyon analiz yontemi kullanilmaktadir. Iki yonteme sahiptir;

- I ) Tek degiskenli regresyon analiz modeli
- II ) Cok degiskenli regresyon analiz modeli

**I ) Tek degiskenli regresyon analiz modeli** : Bir bagimli degisken ve bir bagimsiz degisken arasindaki iliskiyi analiz eder. Bu iki iliski arasinda temsili olarak bir dogrusallik ifade vardir ve bu bir dogrunun denklemleri formulu ile ifade edilir.  $y = a + bx + e$  denklemleri kullanilabilir.

**II ) Çok degiskenli regresyon analiz modeli :** Bir bagimli degisken ve birden fazla bagimsiz degisken arasindaki iliskiyi analiz eder.

### **ID3 (Iterative Dichotomiser 3 - Tekrarlanan ikili yapı) :**

Ozunde Entropy hesaplamasi kullanan bir algoritmadir. Entropy, bir veri kumesindeki verilerin belirsizliginin sayisallastirilmesi demektir. Algoritmanin amaci, egitim kumesindeki verilerin agacin olusturulmasi esnasinda birbirine benzetilmesi gerekiyor, agac derinliginin minimum olmasi, karmasikliginda minimum olmasini saglarken kazancin maksimum olmasi gozle gorulur bir fark alacaktır. Entropy deger araligi  $0 < \text{entropy} < \log_2 n$  arasinda olmalidir. Entropy degeri  $\log_2 n$ 'e yaklastikca belirsizligin artmasi, 0'a yaklasmasiyla belirsizligin azalmasi olarak bilinecektir. Ilk hesaplanmasi gereken entropy, tum data setin hesaplanmasi ile olusur. Sonrasinda datasetin farkli nitelikleri icinde entropy hesaplanmaktadır. Tum bu islemleri bilgi edinim icin kullanilmaktadır. Bu kavramdaki kazanc ise, ilk hesaplanan entropy ile her bir alt kumenin olusumundaki entropilerin arasindaki fark hesaplandiktan sonra, farki buyuk olan karar agacinin sagligi acisinden en dogru dallanmayi yapmis olacaktır.

**Avantaj :** Olasilik kurallari icin egitim verileri kullanilir ve tum data set'teki veriler agac olusturulmasi icin analiz edilir, sonucta kisa agaclar olusturdugu icin en hizli yapiya sahip olmus olur. Bu da, Test sayilarinin azalmasi ve test verilerinin budanmasini saglamaktadır.

**Dezavantaj :** Dataset'ten aldigimiz egitim verilerimizin boyutu kucukse agac test edildiginde cikan sonucun basarisiz olma olasiligi cok yuksektir.

### **C4.5 Algoritmasi :**

ID3 algoritmasının gelistircisi, ID3 algoritması sonucunda siniflandirmalarda bazi eksiklikler ve sorunlar tespit etmistir. Bu sorunların giderilmesini C4.5 algoritması ile saglamistir. Kokeni tamamen ID3 algoritması olan C4.5 algoritmasına gelen ek ozellikler; bolunme-dagilma bilgisinin (split-info) edinilmesi, kayip degerleri olan ozelliklerin tespit edilmesi ve sayisal ozellikteki verilerin hesaba direk olarak katilmasi. Adim adim inceleyecek olursak;

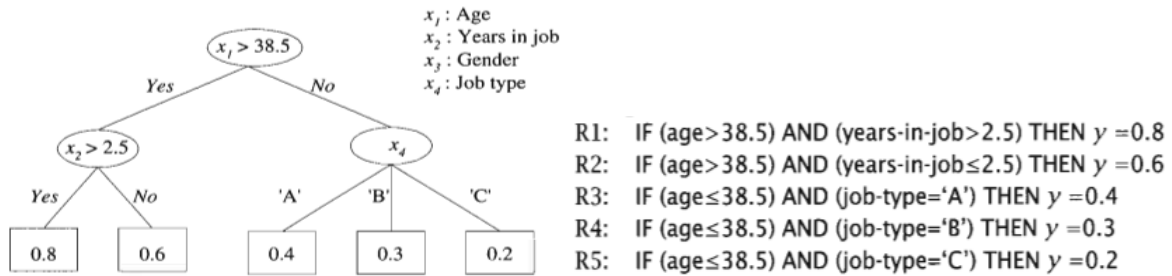
**1. Split Information - Bolunme dallanma bilgisi :** Bir ozniteligin olasiligi ne kadar yuksek olursa bilgi kazancida yuksek olacaktır, bu durum sonucun dogrulunu olumsuz etkileyecektir. Yani, bilgi kazancinin yuksek cikmasinin nedeni ozellik cesitliliğinin fazla olmasidir. Bu tarz gereksiz bilgilerin onlenmesi icin bolunme bilgisi kavrami algoritmaya katilmistir. Bu algoritma ile bilgi kazancini azaltarak gereksiz bazi cikarimlar yapilmasi engellenmistir.

**2. Ozelliklerin kayip degerlerinin tespit edilmesi :** Dataset üzerindeki veriler cesitli sebeplerden dolayi eskik olabilir. ID3 algoritması eksik olmayan bir veri kumesini dikkate alır. Eksik olan veriler yuzunden bazi bilgilerin bulunmasinda yanlisliklar cikmaktadır. Bu sebepten, 3 soruna karsi karsiya kalinmaktadır; bilgi kazanci ve bilgi oranlarının kayip oldugu veri kumesinin hesabi nasil yapılabilir, karar agaci olusturulurken oznitelik degeri olmayanlar alt dugumlere nasil yayilabilir ve bununla beraberinde test islemleri agacin dallarinda nasil yapılabilir? Tum bunlara cevap bulabilmek icin verilerin oznitelikleri ve bu niteliklere sahip bilgi kazanclari sayesinde cozumlere ulasilabilmektedir.

**3. Sayisal Ozellikteki verilerin hesaba direkt katilmasi :** Veri kumelerinde iki tip veri kullanilir; Nominal (kategorik) ve sayisal. ID3 algoritmasinda sadece nominal degerler kullanilirken, C4.5 algoritmasinda

sayisal verilerde yer verilmistir. Tabi bu kullanimda bir yonteme ihtiyac duyulmustur. Sayisal degerler arasinda uygun bir esik degerinin bulunmasi gerekiyor. Esik degeri bulunduktan sonra ikili bir bolunme ile veri kumeleri dagitilabiliyor. Yani, bu esik degerinden buyuk ve esik degerinden kucuk veriler olmak uzere ikiye ayrim yapiliyor.

Esik degeri belirlenirken, tum sayisal verileri kucukten buyuge bir sekilde siralariz.  $\{x_1, x_2, x_3, \dots, x_m\}$  kumesinde,  $m$  tane sayisal veri icerisinden  $x_i$  '.nci veriyi esik degeri olarak seceriz. Bu secimden sonraki siralama artik;  $\{x_1, x_2, x_3, \dots, x_i\}$  ve  $\{x_{i+1}, x_{i+2}, x_{i+3}, \dots, x_m\}$  seklinde iki grup haline donusecektir. Bunun anlami aslinda verilerimiz icerisinden  $m-1$  adet esik degeri secebilecegimizdir. Bu yuzden, olasi butun esik degerlerini  $(x_i + x_{i+1} / 2)$  seklinde formulize edebiliriz. Esik degerine  $e$  dersek,  $x_i < e$  sartini saglayan veriler kucuk  $x_i > e$  sartini saglayan veriler ise buyuk seklinde gruplandirilacaktir.



### ID3 UYGULAMASI

#### 1.Adım : Dataset

Week	Weather	Temperatures (Isi)	Wetness	Wind (Ruzgar)	Game
W1	Sunny	Warm	High	Slightly	No
W2	Sunny	Warm	High	Strong	No
W3	Cloudy	Warm	High	Slightly	Yes
W4	Rainy	Warmish	High	Slightly	Yes
W5	Rainy	Cold	Normal	Slightly	Yes
W6	Rainy	Cold	Normal	Strong	No
W7	Cloudy	Cold	Normal	Strong	Yes
W8	Sunny	Warmish	High	Slightly	No
W9	Sunny	Cold	Normal	Slightly	Yes
W10	Rainy	Warmish	Normal	Slightly	Yes
W11	Sunny	Warmish	Normal	Strong	Yes
W12	Cloudy	Warmish	High	Strong	Yes
W13	Cloudy	Warm	Normal	Slightly	Yes
W14	Rainy	Warmish	High	Strong	No

Game = {No, No, Yes, Yes, Yes, No, Yes, No, Yes, Yes, Yes, Yes, Yes, No}

$P(\text{Game, No}) = 5/14$

$P(\text{Game, Yes}) = 9/14$

$E(\text{Game}) = - ( P(\text{Game, No}) \log_2 P(\text{Game, No}) + P(\text{Game, Yes}) \log_2 P(\text{Game, Yes}) )$

$E(\text{Game}) = - ( 5/14 \log_2 5/14 + 9/14 \log_2 9/14 )$

$E(\text{Game}) = 0.940$



**1. Adim – Isi niteligi Entropy - Gain Cozumu**

$$[\text{Temperatures,Cold}] = 4$$

$$[\text{Temperatures,Warmish}] = 6$$

$$[\text{Temperatures,Warm}] = 4$$

$$E(\text{Game}) = 0.940$$

$$E(\text{Temp,Game}) = P(\text{Temp,Cold}) * E(\text{Temp,Cold}) + P(\text{Temp,Warmish}) * E(\text{Temp,Warmish}) + P(\text{Temp,Warm}) * E(\text{Temp,Warm})$$

$$E(\text{Temp,Game}) = 4/14 E(\text{Temp,Cold}) + 6/14 E(\text{Temp,Warmish}) + 4/14 E(\text{Temp,Warm})$$

$$E(\text{Temp,Cold}) = - (1/4 \log_2 1/4 + 3/4 \log_2 3/4)$$

$$E(\text{Temp,Cold}) = 0.811$$

$$E(\text{Temp,Warmish}) = - (2/6 \log_2 2/6 + 4/6 \log_2 4/6)$$

$$E(\text{Temp,Warmish}) = 0.918$$

$$E(\text{Temp,Warm}) = - (2/4 \log_2 2/4 + 2/4 \log_2 2/4)$$

$$E(\text{Temp,Warm}) = 1.00$$

$$E(\text{Temp,Game}) = 4/14 * 0.811 + 6/14 * 0.918 + 4/14 * 1.00$$

$$E(\text{Temp,Game}) = 0.911$$

$$\text{Gain}(\text{Temp,Game}) = E(\text{Game}) - E(\text{Temp,Game})$$

$$\text{Gain}(\text{Temp,Game}) = 0.940 - 0.911$$

$$\text{Gain}(\text{Temp,Game}) = 0.029$$

Week	Temperatures	Game
W1	Warm	No
W2	Warm	No
W3	Warm	Yes
W4	Warmish	Yes
W5	Cold	Yes
W6	Cold	No
W7	Cold	Yes
W8	Warmish	No
W9	Cold	Yes
W10	Warmish	Yes
W11	Warmish	Yes
W12	Warmish	Yes
W13	Warm	Yes
W14	Warmish	No

**1. Adim – Hava niteligi Entropy - Gain Cozumu**

$$[\text{Weather,Sunny}] = 5$$

$$[\text{Weather,Rainy}] = 5$$

$$[\text{Weather,Cloudy}] = 4$$

$$E(\text{Game}) = 0.940$$

$$E(\text{Weather,Game}) = P(\text{Weather,Sunny}) * E(\text{Weather,Sunny}) + P(\text{Weather,Cloudy}) * E(\text{Weather,Cloudy}) + P(\text{Weather,Rainy}) * E(\text{Weather,Rainy})$$

$$E(\text{Weather,Game}) = 5/14 E(\text{Weather,Sunny}) + 4/14 E(\text{Weather,Cloudy}) + 5/14 E(\text{Weather,Rainy})$$

$$E(\text{Weather,Sunny}) = - (3/5 \log_2 3/5 + 2/5 \log_2 2/5)$$

$$E(\text{Weather,Sunny}) = 0.971$$

$$E(\text{Weather,Rainy}) = - (2/5 \log_2 2/5 + 3/5 \log_2 3/5)$$

$$E(\text{Weather,Rainy}) = 0.971$$

$$E(\text{Weather,Cloudy}) = - (4/4 \log_2 4/4)$$

$$E(\text{Weather,Cloudy}) = 0$$

$$E(\text{Weather,Game}) = 5/14 * 0.971 + 5/14 * 0.971 + 4/14 * 0$$

$$E(\text{Weather,Game}) = 0.694$$

$$\text{Gain}(\text{Weather,Game}) = E(\text{Game}) - E(\text{Weather,Game})$$

$$\text{Gain}(\text{Weather,Game}) = 0.940 - 0.694$$

$$\text{Gain}(\text{Weather,Game}) = 0.247$$

Week	Weather	Game
W1	Sunny	No
W2	Sunny	No
W3	Cloudy	Yes
W4	Rainy	Yes
W5	Rainy	Yes
W6	Rainy	No
W7	Cloudy	Yes
W8	Sunny	No
W9	Sunny	Yes
W10	Rainy	Yes
W11	Sunny	Yes
W12	Cloudy	Yes
W13	Cloudy	Yes
W14	Rainy	No

**1. Adim – Nem niteligi Entropy - Gain Cozumu**

$$[\text{Wetness,High}] = 7$$

$$[\text{Wetness,Normal}] = 7$$

$$E(\text{Game}) = 0.940$$

$$E(\text{Wetness,Game}) = P(\text{Wetness,High}) * E(\text{Wetness,High}) + P(\text{Wetness,Normal}) * E(\text{Wetness,Normal})$$

$$E(\text{Wetness}, \text{Game}) = 7/14 E(\text{Wetness}, \text{High}) + 7/14 E(\text{Wetness}, \text{Normal})$$

$$E(\text{Wetness}, \text{Sunny}) = - (4/7 \log_2 4/7 + 3/7 \log_2 3/7)$$

$$E(\text{Wetness}, \text{Sunny}) = 0.985$$

$$E(\text{Wetness}, \text{Rainy}) = - (1/7 \log_2 1/7 + 6/7 \log_2 6/7)$$

$$E(\text{Wetness}, \text{Rainy}) = 0.592$$

$$E(\text{Wetness}, \text{Game}) = 7/14 * 0.985 + 7/14 * 0.592$$

$$E(\text{Wetness}, \text{Game}) = 0.789$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = E(\text{Game}) - E(\text{Wetness}, \text{Game})$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.940 - 0.789$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.151$$

Week	Wetness	Game
W1	High	No
W2	High	No
W3	High	Yes
W4	High	Yes
W5	Normal	Yes
W6	Normal	No
W7	Normal	Yes
W8	High	No
W9	Normal	Yes
W10	Normal	Yes
W11	Normal	Yes
W12	High	Yes
W13	Normal	Yes
W14	High	No

### 1. Adım – Ruzgar niteligi Entropy - Gain Cozumu

$$[\text{Wind}, \text{Slightly}] = 8$$

$$[\text{Wind}, \text{Strong}] = 6$$

$$E(\text{Game}) = 0.940$$

$$E(\text{Wind}, \text{Game}) = P(\text{Wind}, \text{Slightly}) * E(\text{Wind}, \text{Slightly}) + P(\text{Wind}, \text{Strong}) * E(\text{Wind}, \text{Strong})$$

$$E(\text{Wind}, \text{Game}) = 8/14 E(\text{Wind}, \text{Slightly}) + 6/14 E(\text{Wind}, \text{Strong})$$

$$E(\text{Wind}, \text{Slightly}) = - (2/8 \log_2 2/8 + 6/8 \log_2 6/8)$$

$$E(\text{Wind}, \text{Slightly}) = 0.811$$

$$E(\text{Wind}, \text{Strong}) = - (3/6 \log_2 3/6 + 3/6 \log_2 3/6)$$

$$E(\text{Wind}, \text{Strong}) = 1.00$$

$$E(\text{Wind}, \text{Game}) = 8/14 * 0.811 + 6/14 * 1.00$$

$$E(\text{Wind}, \text{Game}) = 0.892$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = E(\text{Game}) - E(\text{Wetness}, \text{Game})$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.940 - 0.892$$

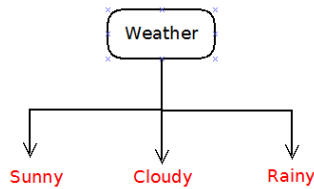
$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.048$$

Week	Wind	Game
W1	Slightly	No
W2	Strong	No
W3	Slightly	Yes
W4	Slightly	Yes
W5	Slightly	Yes
W6	Strong	No
W7	Strong	Yes
W8	Slightly	No
W9	Slightly	Yes
W10	Slightly	Yes
W11	Strong	Yes
W12	Strong	Yes
W13	Slightly	Yes
W14	Strong	No

### 1. Adım – Birinci Dallanma Sonucu

Oznitelikler	Bilgi Kazanci
<b>Weather</b>	<b>0.246</b>
<b>Temperature</b>	<b>0.029</b>
<b>Wetness</b>	<b>0.151</b>
<b>Wind</b>	<b>0.048</b>

Birinci Dallarima Sonucuna Ait Olusan Karar Agaci



**2. Adim**

Week	Weather	Temperatures	Wetness	Wind	Game
W1	Sunny	Warm	High	Slightly	No
W2	Sunny	Warm	High	Strong	No
W8	Sunny	Warmish	High	Slightly	No
W9	Sunny	Cold	Normal	Slightly	Yes
W11	Sunny	Warmish	Normal	Strong	Yes

Hava ozniteliginin gunesli degeri icin dallanma degerleri

Game = {No, No, No, Yes, Yes}

$P(\text{Game}, \text{No}) = 3/5$

$P(\text{Game}, \text{Yes}) = 2/5$

$E(\text{Game}) = - ( P(\text{Game}, \text{No}) \log_2 P(\text{Game}, \text{No}) + P(\text{Game}, \text{Yes}) \log_2 P(\text{Game}, \text{Yes}) )$

$E(\text{Game}) = - ( 3/5 \log_2 3/5 + 2/5 \log_2 2/5 )$

$E(\text{Game}) = 0.970$

**2. Adim – Isi niteligi Entropy - Gain Cozumu**

$[ \text{Temperatures}, \text{Cold} ] = 1$

$[ \text{Temperatures}, \text{Warmish} ] = 2$

$[ \text{Temperatures}, \text{Warm} ] = 2$

$E(\text{Game}) = 0.970$

$E(\text{Temp}, \text{Game}) = P(\text{Temp}, \text{Cold}) * E(\text{Temp}, \text{Cold}) + P(\text{Temp}, \text{Warmish}) * E(\text{Temp}, \text{Warmish}) + P(\text{Temp}, \text{Warm}) * E(\text{Temp}, \text{Warm})$

$E(\text{Temp}, \text{Game}) = 1/5 E(\text{Temp}, \text{Cold}) + 2/5 E(\text{Temp}, \text{Warmish}) + 2/5 E(\text{Temp}, \text{Warm})$

$E(\text{Temp}, \text{Cold}) = - ( 1/1 \log_2 1/1 )$

$E(\text{Temp}, \text{Cold}) = 0$

$E(\text{Temp}, \text{Warm}) = - ( 2/2 \log_2 2/2 )$

$E(\text{Temp}, \text{Warm}) = 0$

$E(\text{Temp}, \text{Warmish}) = - ( 1/2 \log_2 1/2 + 1/2 \log_2 1/2 )$

$E(\text{Temp}, \text{Warmish}) = 1.00$

$E(\text{Temp}, \text{Game}) = 1/5 * 0 + 2/5 * 0 + 2/5 * 1$

$E(\text{Temp}, \text{Game}) = 0.4$

$\text{Gain}(\text{Temp}, \text{Game}) = E(\text{Game}) - E(\text{Temp}, \text{Game})$

$\text{Gain}(\text{Temp}, \text{Game}) = 0.970 - 0.4$

$\text{Gain}(\text{Temp}, \text{Game}) = 0.570$

Week	Temperatures	Game
W1	Warm	No
W2	Warm	No
W8	Warmish	No
W9	Cold	Yes
W11	Warmish	Yes

**2. Adim – Nem niteligi Entropy - Gain Cozumu**

$[ \text{Wetness}, \text{High} ] = 3$

$[ \text{Wetness}, \text{Normal} ] = 2$

$E(\text{Game}) = 0.970$

$E(\text{Wetness}, \text{Game}) = P(\text{Wetness}, \text{High}) * E(\text{Wetness}, \text{High}) + P(\text{Wetness}, \text{Normal}) * E(\text{Wetness}, \text{Normal})$

$E(\text{Wetness}, \text{Game}) = 3/5 E(\text{Wetness}, \text{High}) + 2/5 E(\text{Wetness}, \text{Normal})$

$E(\text{Wetness}, \text{Sunny}) = - ( 3/3 \log_2 3/3 )$

$E(\text{Wetness}, \text{Sunny}) = 0$

$E(\text{Wetness}, \text{Rainy}) = - ( 2/2 \log_2 2/2 )$

$E(\text{Wetness}, \text{Rainy}) = 0$

$$E(\text{Wetness}, \text{Game}) = 3/5 * 0 + 2/5 * 0$$

$$E(\text{Wetness}, \text{Game}) = 0$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = E(\text{Game}) - E(\text{Wetness}, \text{Game})$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.970 - 0$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.970$$

Week	Wetness	Game
W1	High	No
W2	High	No
W8	High	No
W9	Normal	Yes
W11	Normal	Yes

## 2. Adım – Ruzgar niteligi Entropy - Gain Cozumu

$$[\text{Wind}, \text{Slightly}] = 3$$

$$[\text{Wind}, \text{Strong}] = 2$$

$$E(\text{Game}) = 0.970$$

$$E(\text{Wind}, \text{Game}) = P(\text{Wind}, \text{Slightly}) * E(\text{Wind}, \text{Slightly}) + P(\text{Wind}, \text{Strong}) * E(\text{Wind}, \text{Strong})$$

$$E(\text{Wind}, \text{Game}) = 3/5 E(\text{Wind}, \text{Slightly}) + 2/5 E(\text{Wind}, \text{Strong})$$

$$E(\text{Wind}, \text{Slightly}) = - (2/3 \log_2 2/3 + 1/3 \log_2 1/3)$$

$$E(\text{Wind}, \text{Slightly}) = 0.918$$

$$E(\text{Wind}, \text{Strong}) = - (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$$

$$E(\text{Wind}, \text{Strong}) = 1.00$$

$$E(\text{Wind}, \text{Game}) = 3/5 * 0.918 + 2/5 * 1.00$$

$$E(\text{Wind}, \text{Game}) = 0.951$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = E(\text{Game}) - E(\text{Wetness}, \text{Game})$$

$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.970 - 0.951$$

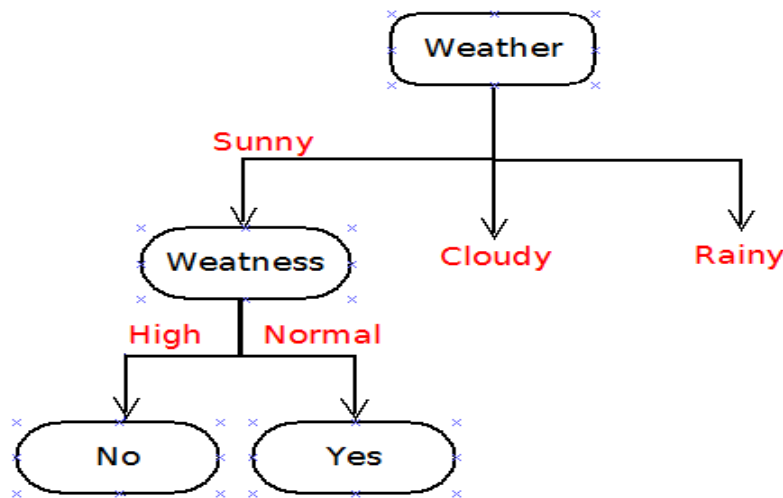
$$\text{Gain}(\text{Wetness}, \text{Game}) = 0.019$$

Week	Wind	Game
W1	Slightly	No
W2	Strong	No
W8	Slightly	No
W9	Slightly	Yes
W11	Strong	Yes

## 2. Adım – İkinci Dallanma Sonucu

Oznitelikler	Bilgi Kazancı
Temperature	0.570
Wetness	0.970
Wind	0.019

İkinci Dallanma Sonucuna Ait Olusan Karar Agaci



**3. Adim**

Week	Weather	Temperatures	Wetness	Wind	Game
W3	Cloudy	Warm	High	Slightly	Yes
W7	Cloudy	Cold	Normal	Strong	Yes
W12	Cloudy	Warmish	High	Strong	Yes
W13	Cloudy	Warm	Normal	Slightly	Yes

Hava ozniteliginin bulutlu degeri icin dallanma degerleri

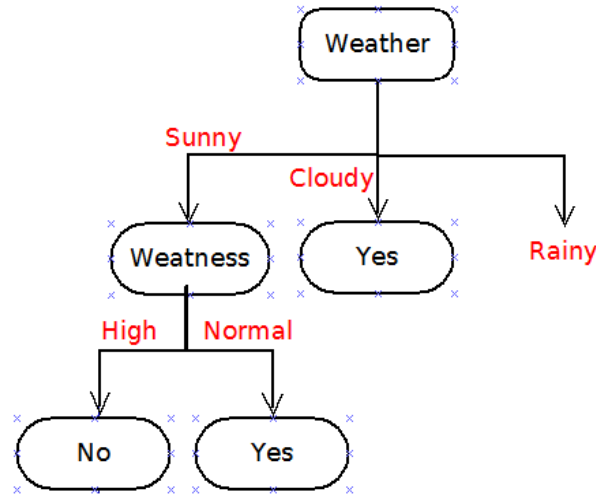
Game = {Yes, Yes, Yes, Yes}

$P(\text{Game}, \text{Yes}) = 4/4$

$E(\text{Game}) = - (P(\text{Game}, \text{Yes}) \log_2 P(\text{Game}, \text{Yes}))$

$E(\text{Game}) = - (4/4 \log_2 4/4)$

$E(\text{Game}) = 1$

**3. Adim – Ucuncu Dallanma Sonucu****4. Adim**

Week	Weather	Temperatures	Wetness	Wind	Game
W4	Rainy	Warmish	High	Slightly	Yes
W5	Rainy	Cold	Normal	Slightly	Yes
W6	Rainy	Cold	Normal	Strong	No
W10	Rainy	Warmish	Normal	Slightly	Yes
W14	Rainy	Warmish	High	Strong	No

Hava ozniteliginin yagmurlu degeri icin dallanma degerleri

Game = {Yes, Yes, No, Yes, No}

$P(\text{Game}, \text{No}) = 2/5$

$P(\text{Game}, \text{Yes}) = 3/5$

$E(\text{Game}) = - (P(\text{Game}, \text{No}) \log_2 P(\text{Game}, \text{No}) + P(\text{Game}, \text{Yes}) \log_2 P(\text{Game}, \text{Yes}))$

$E(\text{Game}) = - (2/5 \log_2 2/5 + 3/5 \log_2 3/5)$

$E(\text{Game}) = 0.970$

**4. Adım – Isi niteligi Entropy - Gain Cozumu**

[Temperatures,Cold] = 2

[Temperatures,Warmish] = 3

$E(\text{Game}) = 0.970$

$E(\text{Temp,Game}) = P(\text{Temp,Cold}) * E(\text{Temp,Cold}) + P(\text{Temp,Warmish}) * E(\text{Temp,Warmish})$

$E(\text{Temp,Game}) = 2/5 E(\text{Temp,Cold}) + 3/5 E(\text{Temp,Warmish})$

$E(\text{Temp,Cold}) = - (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$

$E(\text{Temp,Cold}) = 1$

$E(\text{Temp,Warmish}) = - (2/3 \log_2 2/3 + 1/3 \log_2 1/3)$

$E(\text{Temp,Warmish}) = 0.918$

$E(\text{Temp,Game}) = 2/5 * 1 + 3/5 * 0.918$

$E(\text{Temp,Game}) = 0.951$

$\text{Gain}(\text{Temp,Game}) = E(\text{Game}) - E(\text{Temp,Game})$

$\text{Gain}(\text{Temp,Game}) = 0.970 - 0.951$

$\text{Gain}(\text{Temp,Game}) = 0.019$

Week	Temperatures	Game
W4	Warmish	Yes
W5	Cold	Yes
W6	Cold	No
W10	Warmish	Yes
W14	Warmish	No

**4. Adım – Ruzgar niteligi Entropy - Gain Cozumu**

[Wind,Slightly] = 3

[Wind,Strong] = 2

$E(\text{Game}) = 0.970$

$E(\text{Wind,Game}) = P(\text{Wind,Slightly}) * E(\text{Wind,Slightly}) + P(\text{Wind,Strong}) * E(\text{Wind,Strong})$

$E(\text{Wind,Game}) = 3/5 E(\text{Wind,Slightly}) + 2/5 E(\text{Wind,Strong})$

$E(\text{Wind,Slightly}) = - (3/3 \log_2 3/3)$

$E(\text{Wind,Slightly}) = 0$

$E(\text{Wind,Strong}) = - (2/2 \log_2 2/2)$

$E(\text{Wind,Strong}) = 0$

$E(\text{Wind,Game}) = 3/5 * 0 + 2/5 * 0$

$E(\text{Wind,Game}) = 0$

$\text{Gain}(\text{Wetness,Game}) = E(\text{Game}) - E(\text{Wetness,Game})$

$\text{Gain}(\text{Wetness,Game}) = 0.970 - 0$

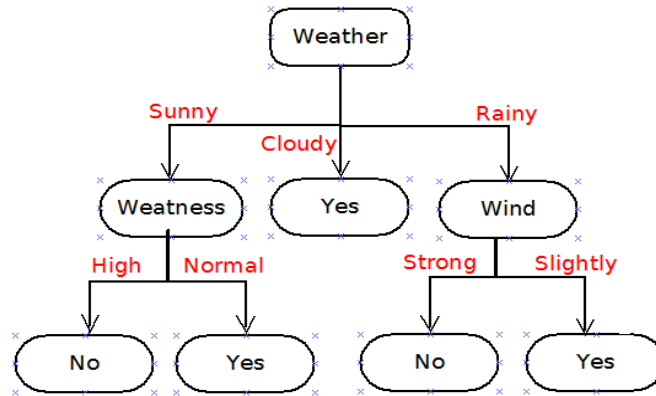
$\text{Gain}(\text{Wetness,Game}) = 0.970$

Week	Wind	Game
W4	Slightly	Yes
W5	Slightly	Yes
W6	Strong	No
W10	Slightly	Yes
W14	Strong	No

#### 4. Adım – Dorduncu Dallanma Sonucu

Oznitelikler	Bilgi Kazancı
Temperature	0.019
Wind	0.970

Dorduncu Dallanma Sonucuna Ait Olusan Karar Agaci



#### C4.5 UYGULAMASI

Dataset

No	Job	Age	Married	Decision
1	Doctor	70	Yes	Class1
2	Doctor	90	Yes	Class2
3	Doctor	85	No	Class2
4	Doctor	95	No	Class2
5	Doctor	70	No	Class1
6	Teacher	90	Yes	Class1
7	Teacher	78	No	Class1
8	Teacher	65	Yes	Class1
9	Teacher	75	No	Class1
10	Engineer	80	Yes	Class2
11	Engineer	70	Yes	Class2
12	Engineer	80	No	Class1
13	Engineer	70	No	Class1
14	Engineer	96	No	Class1

Age = {65, 70, 75, 80, 85, 90, 95, 96}

Esik Degeri =>  $t = (x_i + x_{i+1}) / 2$

$t = (80 + 85) / 2$

$t = 83$

Age <= 83 ve Age > 83

No	Job	Age	Married	Decision
1	Doctor	Age <= 83	Yes	Class1
2	Doctor	Age > 83	Yes	Class2
3	Doctor	Age > 83	No	Class2
4	Doctor	Age > 83	No	Class2
5	Doctor	Age <= 83	No	Class1
6	Teacher	Age > 83	Yes	Class1
7	Teacher	Age <= 83	No	Class1
8	Teacher	Age <= 83	Yes	Class1
9	Teacher	Age <= 83	No	Class1
10	Engineer	Age <= 83	Yes	Class2
11	Engineer	Age <= 83	Yes	Class2
12	Engineer	Age <= 83	No	Class1
13	Engineer	Age <= 83	No	Class1
14	Engineer	Age > 83	No	Class1

$$[\text{Decision}, \text{Class1}] = 9 / 14$$

$$[\text{Decision}, \text{Class2}] = 5 / 14$$

$$E(\text{Decision}) = - ( P(\text{Decision}, \text{Class1}) \log_2 P(\text{Decision}, \text{Class1}) + P(\text{Decision}, \text{Class2}) \log_2 P(\text{Decision}, \text{Class2}) )$$

$$E(\text{Decision}) = - ( 5/14 \log_2 5/14 + 9/14 \log_2 9/14 )$$

$$E(\text{Decision}) = 0.940$$

$$E(\text{Decision}) = 0.940$$

$$[\text{Age} \leq 83, \text{Class1}] = 7$$

$$[\text{Age} \leq 83, \text{Class2}] = 2$$

$$E(\text{Age} \leq 83, \text{Decision}) = - ( 7/9 \log_2 7/9 + 2/9 \log_2 2/9 )$$

$$E(\text{Age} \leq 83, \text{Decision}) = 0.764$$

$$[\text{Age} > 83, \text{Class1}] = 2$$

$$[\text{Age} > 83, \text{Class2}] = 3$$

$$E(\text{Age} > 83, \text{Decision}) = - ( 2/5 \log_2 2/5 + 3/5 \log_2 3/5 )$$

$$E(\text{Age} > 83, \text{Decision}) = 0.970$$

$$E(\text{Age}, \text{Decision}) = P(\text{Age} \leq 83, \text{Decision}) * E(\text{Age} \leq 83, \text{Decision}) + P(\text{Age} > 83, \text{Decision}) *$$

$$E(\text{Age} > 83, \text{Decision})$$

$$E(\text{Age}, \text{Decision}) = 9/14 E(\text{Age} \leq 83, \text{Decision}) + 5/14 E(\text{Age} > 83, \text{Decision})$$

$$E(\text{Age}, \text{Decision}) = 9/14 * 0.764 + 5/14 * 0.970$$

$$\text{Gain}(\text{Age}, \text{Decision}) = E(\text{Decision}) - E(\text{Age}, \text{Decision})$$

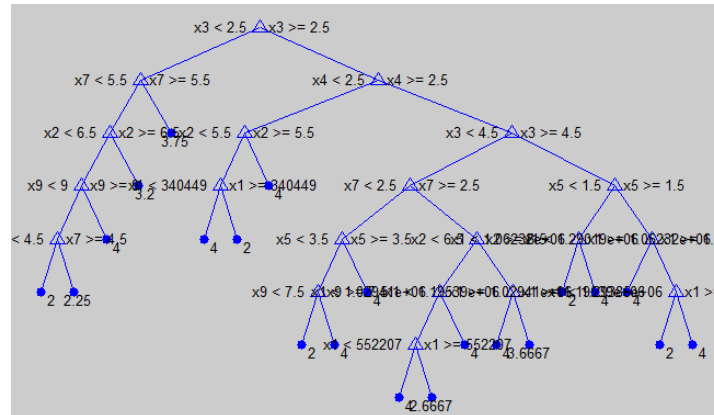
$$\text{Gain}(\text{Age}, \text{Decision}) = 0.940 - 0.837$$

$$\text{Gain}(\text{Age}, \text{Decision}) = 0.103$$





## Regression Tree



**Uygulama 2 :** Breast Cancer Wisconsin (Original) Data Set kullanılarak dogruluk, hata oranlari ve confusion matrix degerleri hesaplaniyor.

### Matlab Kod

```
clear all;
close all;
clc;

dataset = load('breast-cancer-wisconsin.data');
dataEgitim = dataset(1:600,1:10);
dataTest = dataset(601:683,1:10);
classEgitim = dataset(1:600,11);
classTest = dataset(601:683,11);

tree = ClassificationTree.fit(dataEgitim, classEgitim)
t = classregtree(dataEgitim, classEgitim);

cvv = crossval(tree);
error = kfoldLoss(cvv)
dogruluk = 1 - error

c1 = tree.predict(dataTest);
cMat = confusionmat(classTest, c1)

error = 0.0517
dogruluk = 0.9483
cMat = 67 2
      0 14
```

**Uygulama 3 :**  $y=f(x_1, x_2, x_3)$ 

$y$	$x_1$	$x_2$	$x_3$
-	0	0	0
-	1	0	0
+	0	0	1
+	1	0	1
+	0	1	0
+	1	1	0
-	0	1	1
-	1	1	1

## Matlab Kod

```

clear all;
close all;
clc;

x1 = [0 1 0 1 0 1 0 1]';
x2 = [0 0 0 0 1 1 1 1]';
x3 = [0 0 1 1 0 0 1 1]';
inData = [x1, x2, x3];
outData = ['- ', '- ', '+ ', '+ ', '+ ', '+ ', '- ', '- '];

mytree = treefit(inData, outData, 'method', 'classification', 'splitmin', 2,
'prune', 'on', 'splitcriterion', 'gdi')
treedisp(mytree);

```

## Decision tree for classification

```

1  if x1<0.5 then node 2 elseif x1>=0.5 then node 3 else -
2  if x2<0.5 then node 4 elseif x2>=0.5 then node 5 else -
3  if x2<0.5 then node 6 elseif x2>=0.5 then node 7 else -
4  if x3<0.5 then node 8 elseif x3>=0.5 then node 9 else -
5  if x3<0.5 then node 10 elseif x3>=0.5 then node 11 else -
6  if x3<0.5 then node 12 elseif x3>=0.5 then node 13 else -
7  if x3<0.5 then node 14 elseif x3>=0.5 then node 15 else -
8  class = -
9  class = +
10 class = +
11 class = -
12 class = -
13 class = +
14 class = +
15 class = -

```

