

Karar ağaçları

Decision Trees Nedir?

- Decision Tree genellikle sınıflandırma problemleri için kullanılan bir gözetimli öğrenme algoritmasıdır. Regresyon içinde kullanılabilir. Yüksek doğruluk, kararlılık ve yorumlanma kolaylığına sahiptir. Doğrusal olmayan ilişkileride oldukça iyi tespit edebilir.

Makine Öğrenimi Karar Ağaçları Ne İçin Kullanılır?

- Bu model çok yönlüdür, dolayısıyla karar ağacı algoritması makine öğreniminde birçok uygulama alanı bulmaktadır:
- Veri madenciliği
- Veri bilimi
- İstatistikler
- Bu tekniğin amacı, çeşitli girdi değişkenlerine dayalı olarak bir hedef değişkenin değerini tahmin eden bir model oluşturmaktır. Bu teknik şu gibi soruları cevaplamak için kullanılır: X yapılmalı mı? A'yı mı yoksa B'yi mi seçmeliyiz?

1. Karar analizi

- Karar verme sürecinde kullanılır. Karar ağacı makine öğrenimi modeli, karar sürecini ve sonucunu açıkça sunmak için veri görselleştirmede kullanılabilir. Bu, karar ağaçlarının en büyük avantajlarından biri, işlerinde makine öğrenimi ile ilgilenmeyen biri için bile anlaşılabilir ve kolay olmasıdır.

2. Sınıflandırma

- Sınıflandırma, bir sınıf değerini tahmin etmek veya açıklamak için kullanılan bir makine öğrenimi tekniğidir. Sınıflandırma algoritmaları ile bir veya daha fazla girdiye dayalı olarak bir olayın meydana gelme olasılığını değerlendirebilirsiniz. Karar ağacı makine öğrenimi modeli tam olarak bununla ilgilidir. Örneğin, e-posta'ları spam ve spam olmayan olarak ikiye ayırmak için kullanılabilir. Algoritmanız bir dizi soru sorar ve yanıtlara göre belirli bir e-posta'nın spam olup olmadığına karar verir.

3. Regresyon

- Denetimli makine öğrenimi tekniklerinden biridir. Regresyon, bir dizi önceki veriye dayanarak belirli bir değeri tahmin etmenize (veya açıklamaya) yardımcı olur. Bir çalışanın maaşı, hastalığın yayılması veya mülk değeri gibi konuları tahmin etmek için regresyonu kullanabilirsiniz. Bu durumda, buna regresyon ağacı denir.
- Gerçek yaşam koşullarında bu model insanlar tarafından yaygın olarak kullanılmaktadır. Arkadaşlarınızla buluşup buluşmayacağınızı ya da televizyon izleyip izlemeyeceğinizi her düşündüğünüzde, sadece zihninizde olsa bile karar ağacı tekniğini kullanırsınız. Makine öğreniminin çeşitli yönlerinde bu kadar yaygın olmasının nedeni budur. Bu, insanların öğrenme şeklini taklit etmeye çalışan bir teknolojidir ve insanlar farkında olmasalar bile bu tekniği kullanır.

Karar ağaçlarında terimler

- Karar ağaçlarının ilk hücrelerine **kök** (root veya root node) denir. Her bir gözlem kökteki koşula göre “Evet” veya “Hayır” olarak sınıflandırılır.
- Kök hücrelerinin altında **düğüm**ler (interval nodes veya nodes) bulunur. Her bir gözlem düğümler yardımıyla sınıflandırılır. Düğüm sayısı arttıkça modelin karmaşıklığı da artar.
- Karar ağacının en altında **yapraklar** (leaf nodes veya leaves) bulunur. Yapraklar, bize sonucu verir.

Karar ağacı kullanılırken yapılan bazı varsayımlar aşağıdadır:

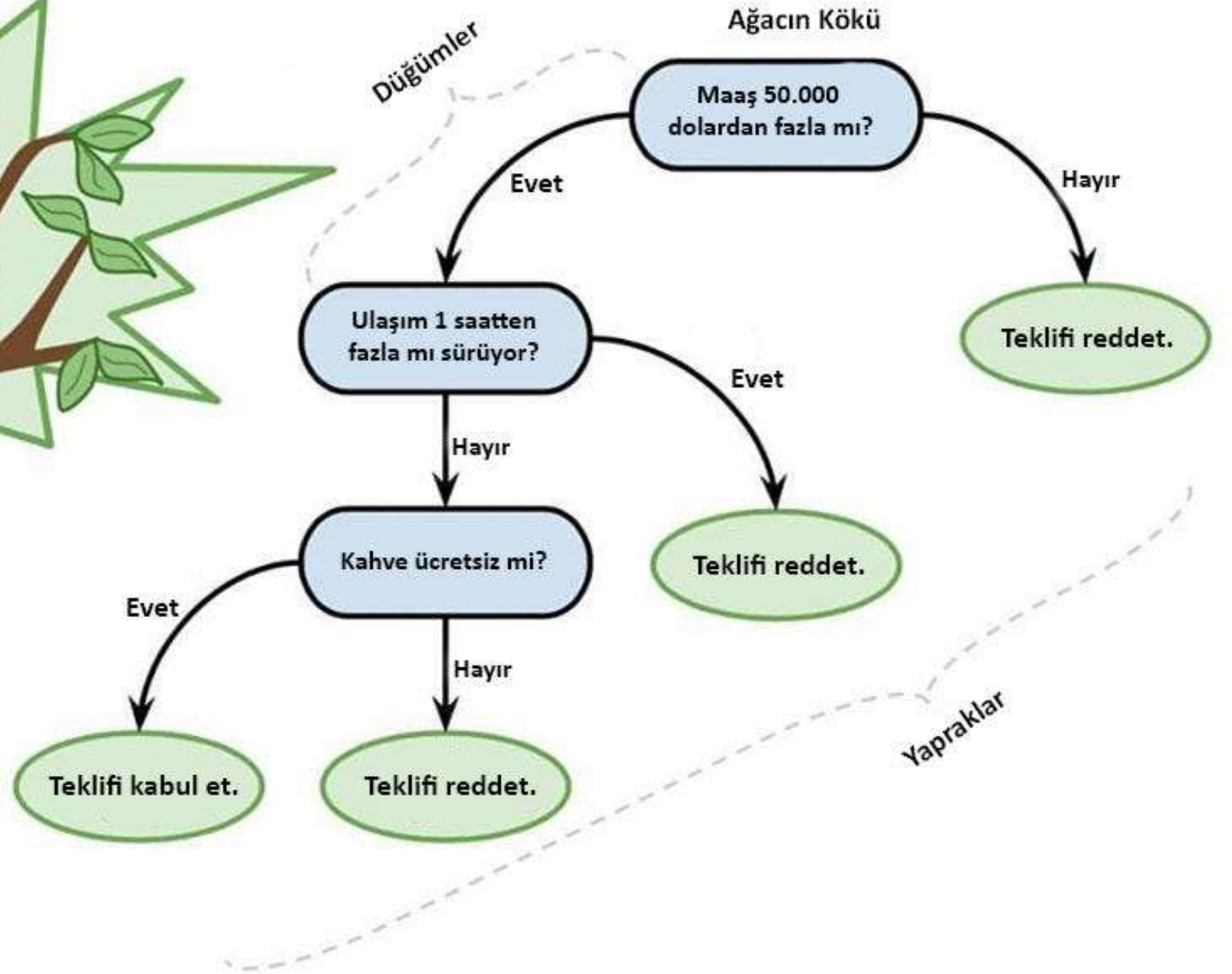
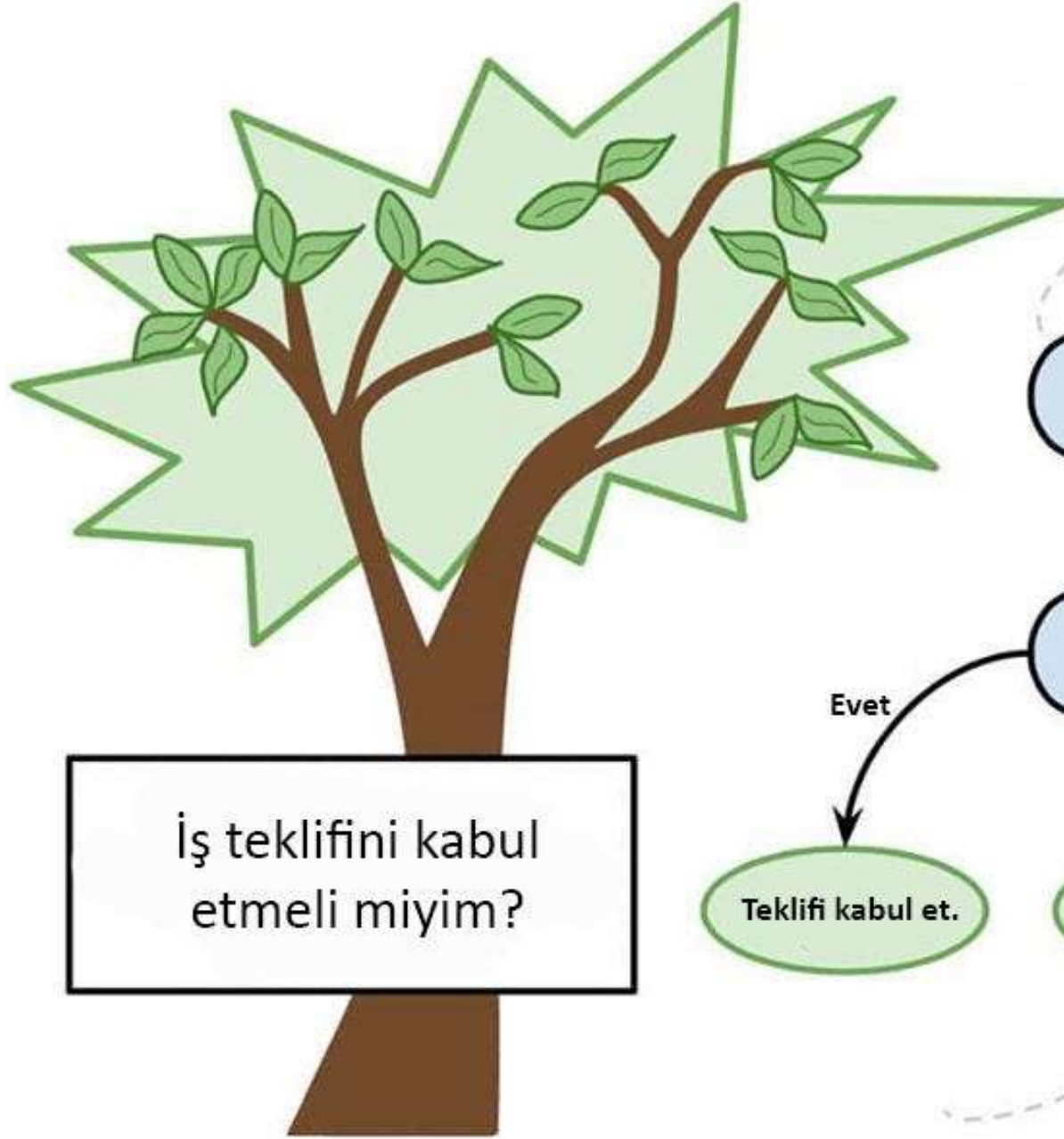
- Başlangıçta, tüm eğitim seti kök olarak kabul edilir.
- Özellik değerlerinin kategorik olması tercih edilir. Değerler sürekli ise, model oluşturmadan önce ayrıklaştırılırlar.
- Öznitelik değerleri temelinde, kayıtlar özyinelemeli olarak dağıtılır.
- Öznitelikleri kök veya dahili düğüm olarak sıralamak için istatistiksel yöntemler kullanılır.

Karar Ağacı Algoritması:

- 1) T öğrenme kümesini oluşturulur.
- 2) T kümesindeki örnekleri en iyi ayıran nitelikler belirlenir.
- 3) Seçilen nitelik ile ağacın düğümleri oluşturulur ve her bir düğümde alt düğümler veya ağacın yapraklarını oluşturulur. Alt düğümlere ait alt veri kümesinin örneklerini belirlenir

Karar Ağacı Algoritması:

- 4) 3. adımda oluşturulan her alt veri kümesi için ;
 - Örneklerin hepsi aynı sınıfa aitse
 - Örnekleri bölecek nitelik kalmamışsa
 - Kalan niteliklerin değerini taşıyan örnek yoksa işlemi sonlandır.
Diğer durumda alt veri kümesini ayırmak için 2. adımdan devam edilir.



Karar Ağaçlarının Oluşturulması

- Karar ağaçları, bir dizi kural kullanılarak veri setinden oluşturulur. Bu kurallar genellikle verinin belirli özelliklerine dayalı olarak bölünmesini içerir. Ağaç oluşturma sürecinde kullanılan yaygın yöntemler:
- **ID3 (Iterative Dichotomiser 3):** Entropi ve bilgi kazancı kullanarak veri setini böler.
- **C4.5:** ID3 algoritmasının geliştirilmiş versiyonudur ve kategorik verilerle daha iyi çalışır.
- **CART (Classification and Regression Trees):** Gini indeksi veya varyans azaltma yöntemlerini kullanarak veri setini böler.

Entropi

- **Entropi:** Bir veri kümesinin düzensizliğini veya belirsizliğini ölçer. Yüksek entropi, verinin daha karışık olduğunu gösterir.

$$\begin{aligned} \text{Entropy } H(X) &= - \left[\left(\frac{3}{8} \right) \log_2 \frac{3}{8} + \left(\frac{5}{8} \right) \log_2 \frac{5}{8} \right] \\ &= - [0.375 * (-1.415) + 0.625 * (-0.678)] \\ &= -(-0.53 - 0.424) \\ &= 0.954 \end{aligned}$$

Bilgi Kazancı

- **Bilgi Kazancı:** Bir özelliğin veriyi ne kadar iyi böldüğünü ölçer. Yüksek bilgi kazancı, daha iyi bir bölünme anlamına gelir.

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

Ezber (Overfitting: Aşırı Uyum):

- Tüm makine öğrenmesi yöntemlerinde verinin ana hatlarının modellenmesi esas alındığı için öğrenme modelinde ezberden (overfitting) kaçınılmalıdır.
- Tüm karar ağaçları önlem alınmazsa ezber yapar. Bu yüzden ağaç oluşturulurken veya oluşturulduktan sonra budama yapılmalıdır.

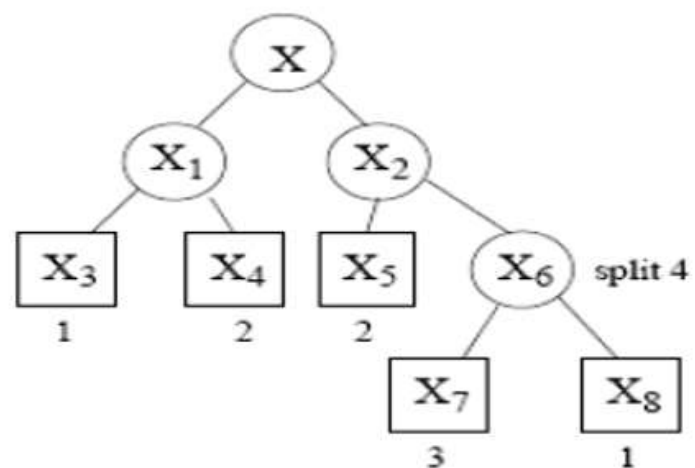
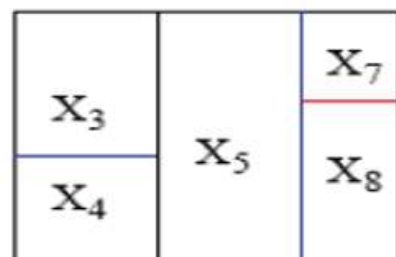
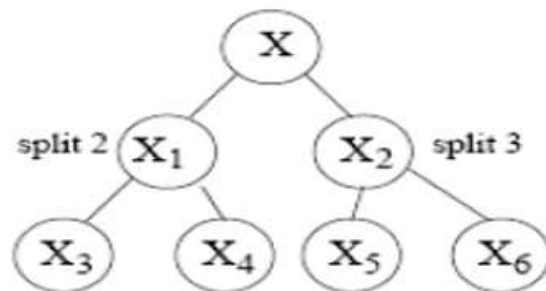
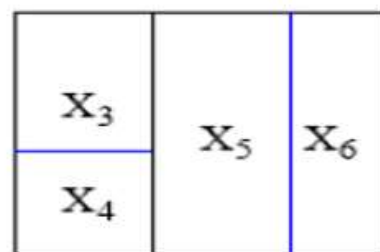
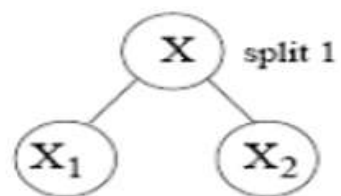
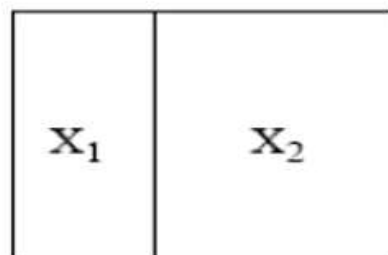
Ağaç Budama:

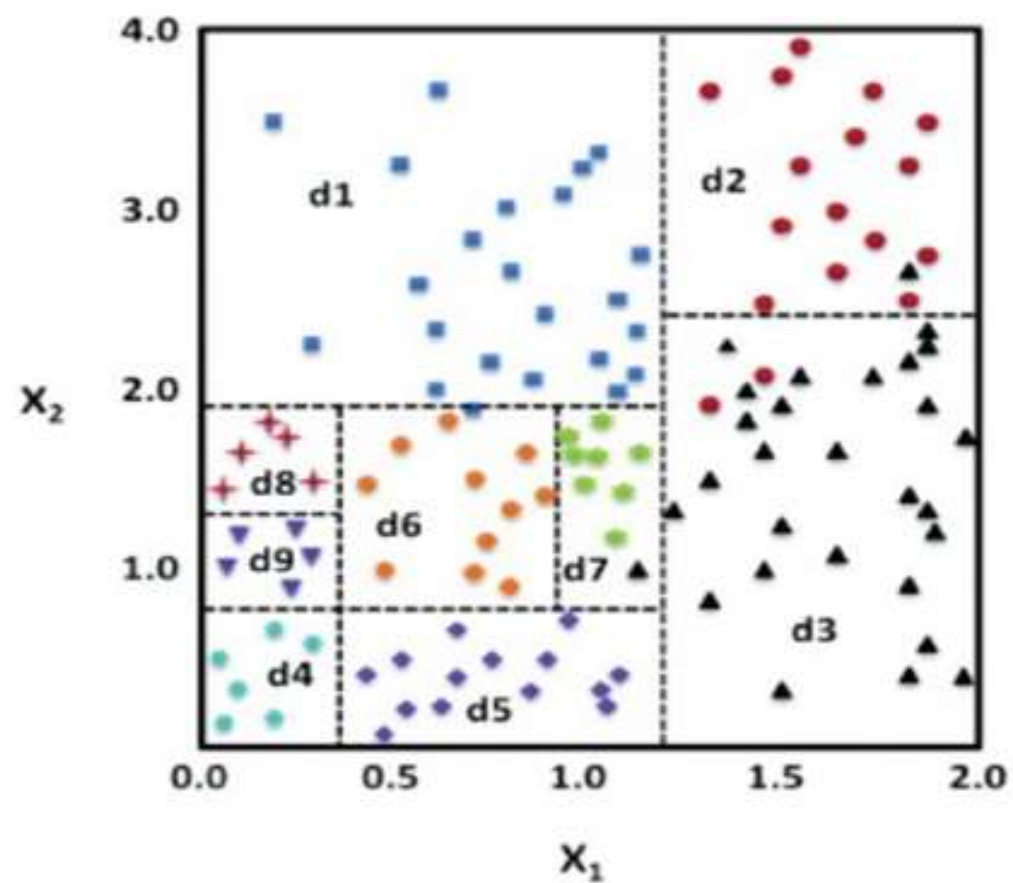
- Budama, sınıflandırmaya katkısı olmayan bölümlerin karar ağacından çıkarılması işlemidir. Bu sayede karar ağacı hem sade hem de anlaşılabilir hale gelir. İki çeşit budama yöntemi vardır; – Ön budama – Sonradan budama
Ön budama işlemi ağaç oluşturulurken yapılır. Bölünen nitelikler, değerleri belli bir esik değerinin (hata toleransının) üstünde değilse o noktada ağaç bölümleme işlemi durdurulur ve o an elde bulunan kümedeki baskın sınıf etiketi, yaprak olarak oluşturulur.
Sonradan Budama: Sonradan budama işlemi ağaç oluşturulduktan sonra devreye girer. Alt ağaçları silerek yaprak oluşturma, alt ağaçları yükseltme, dal kesme şeklinde yapılabilir.

Budama işleminde ne zaman duracağınıza karar vermenin birkaç yolu vardır:

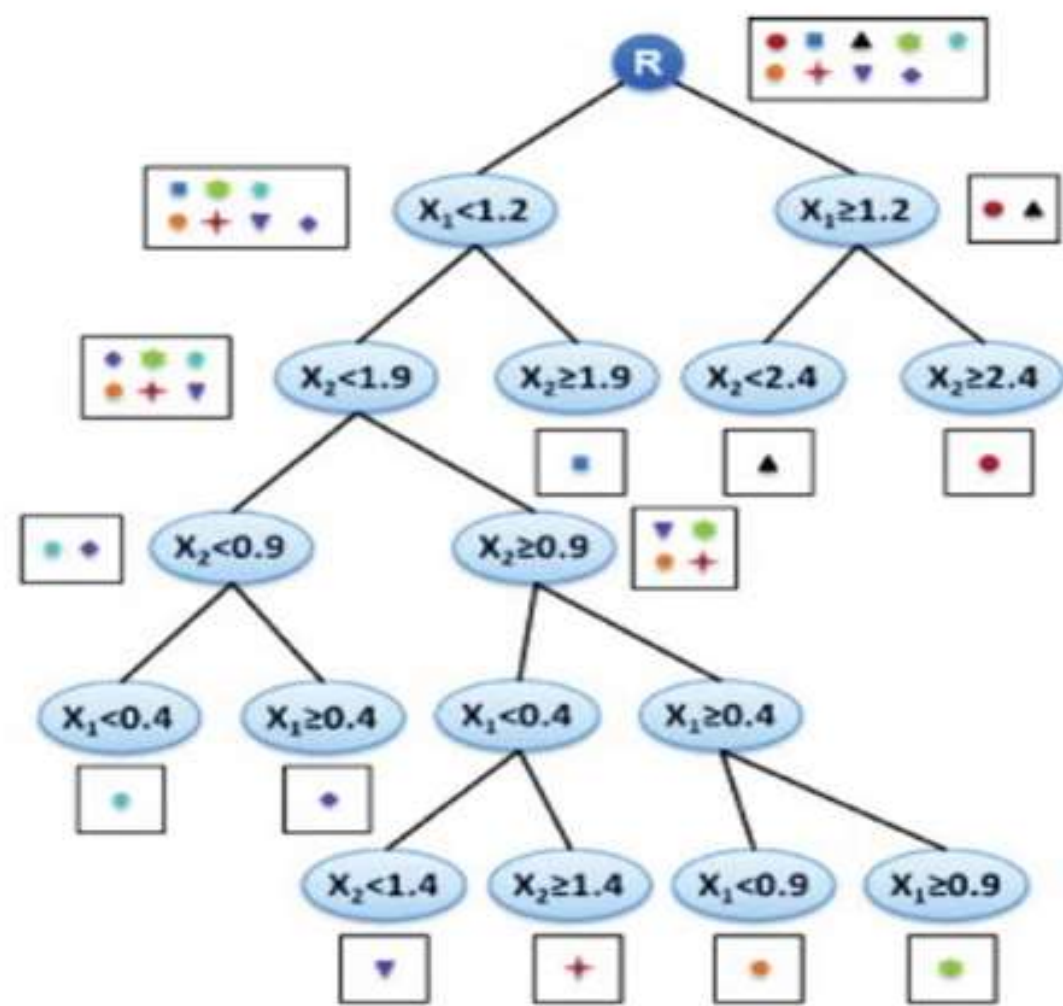
- Tüm terminal düğümleri saf olana kadar devam edilir.
- Her bir terminal düğümündeki veri sayısı belirli bir eşikten, örneğin 5'ten, hatta 1'den büyük olmayana kadar devam edilir.
- Ağaç yeterince büyük olduğu sürece, ilk ağacın boyutu kritik değildir.
- Buradaki anahtar, ilk ağacı yeniden budamadan önce yeterince büyük yapmaktır!

Sınıflandırma Ağaçları:





(a)



(b)

Adım 1: Hedefin entropisini hesaplanır.

$$\begin{aligned}\text{Entropy}(\text{PlayGolf}) &= \text{Entropy}(5,9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= -(0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94\end{aligned}$$

Adım 2: Veri kümesi daha sonra farklı niteliklere bölünür. Her dal için entropi hesaplanır. Daha sonra, bölme için toplam entropi elde etmek için orantılı olarak eklenir. Ortaya çıkan entropi, bölünmeden önceki entropiden çıkarılır. Sonuç, Bilgi Kazanımı veya entropideki azalmadır.

		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3
Gain = 0.247			

		Play Golf	
		Yes	No
Temp.	Hot	2	2
	Mild	4	2
	Cool	3	1
Gain = 0.029			

		Play Golf	
		Yes	No
Humidity	High	3	4
	Normal	6	1
Gain = 0.152			

		Play Golf	
		Yes	No
Windy	False	6	2
	True	3	3
Gain = 0.048			

$$\text{Gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

$$\begin{aligned}\text{G}(\text{PlayGolf}, \text{Outlook}) &= \text{E}(\text{PlayGolf}) - \text{E}(\text{PlayGolf}, \text{Outlook}) \\ &= 0.940 - 0.693 = 0.247\end{aligned}$$

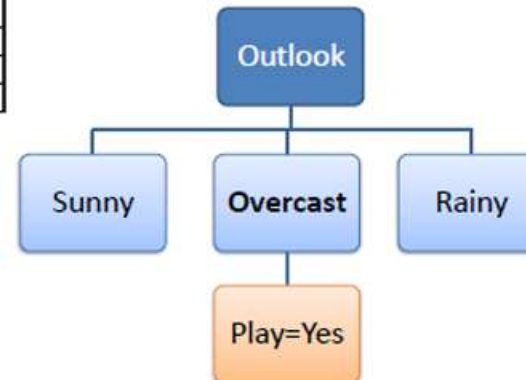
Adım 3: Karar düğümü olarak en büyük bilgi kazancına sahip öznitelik seçilir, veri seti dallarına bölünür ve aynı işlemi her dalda tekrarlanır.

		Play Golf	
		Yes	No
★	Sunny	3	2
	Overcast	4	0
	Rainy	2	3
Gain = 0.247			

		Outlook	Temp	Humidity	Windy	Play Golf
Outlook	Sunny	Sunny	Mild	High	FALSE	Yes
		Sunny	Cool	Normal	FALSE	Yes
		Sunny	Cool	Normal	TRUE	No
		Sunny	Mild	Normal	FALSE	Yes
		Sunny	Mild	High	TRUE	No
	Overcast	Overcast	Hot	High	FALSE	Yes
		Overcast	Cool	Normal	TRUE	Yes
		Overcast	Mild	High	TRUE	Yes
		Overcast	Hot	Normal	FALSE	Yes
	Rainy	Rainy	Hot	High	FALSE	No
		Rainy	Hot	High	TRUE	No
		Rainy	Mild	High	FALSE	No
		Rainy	Cool	Normal	FALSE	Yes
		Rainy	Mild	Normal	TRUE	Yes

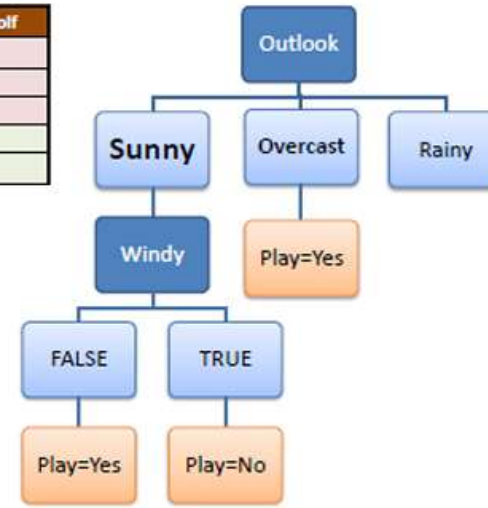
Step 4a: Entropisi 0 olan bir dal, bir yaprak düğümdür.

Temp	Humidity	Windy	Play Golf
Hot	High	FALSE	Yes
Cool	Normal	TRUE	Yes
Mild	High	TRUE	Yes
Hot	Normal	FALSE	Yes



Adım 4b: Entropisi 0'dan büyük olan bir dalın daha fazla bölünmesi gerekir.

Temp	Humidity	Windy	Play Golf
Mild	High	FALSE	Yes
Cool	Normal	FALSE	Yes
Mild	Normal	FALSE	Yes
Cool	Normal	TRUE	No
Mild	High	TRUE	No



Adım 5: ID3 algoritması, tüm veriler sınıflandırılana kadar yaprak olmayan dallarda özylenelemeli olarak çalıştırılır.

Karar Ağacından Karar Kurallarına

Bir karar ağacı, kök düğümden yaprak düğümlere tek tek eşlenerek kolayca bir dizi kurala dönüştürülebilir.

R_1 : IF (Outlook=Sunny) AND (Windy=FALSE) THEN Play=Yes

R_2 : IF (Outlook=Sunny) AND (Windy=TRUE) THEN Play=No

R_3 : IF (Outlook=Overcast) THEN Play=Yes

R_4 : IF (Outlook=Rainy) AND (Humidity=High) THEN Play=No

R_5 : IF (Outlook=Rain) AND (Humidity=Normal) THEN Play=Yes

