

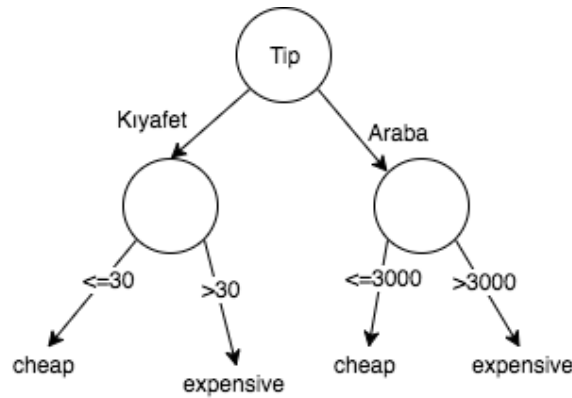
# Machine Learning 101

## Öğrenme Yöntemleri

1. Eğitimli Öğrenme
2. Denetimsiz Öğrenme
3. Destekli Öğrenme

## Karar Ağacı (ID3)

Parametrik olmayan, eğitimli bir yöntemdir. Sınıflandırma ve regresyon için kullanılırlar. Bir karar ağacı ara karar düğümü ve terminal yapraklardan oluşur. Yaprak düğümlerin bir output değeri vardır. Bu output sınıflandırmada sınıf kodu, regresyonda da nümerik bir değer olarak görünür. Karar ağaçlarında alt kümeler bölmenin amacı her alt kümeyi olabildiğince homojen hale getirmektir. Dezavantajı karar ağacı algoritmalarının greedy yöntemler olmasıdır.



## Entropi

$$2\text{-sınıflı Entropi}(S) = -(p_+ \cdot \log_2 p_+ + p_- \cdot \log_2 p_-)$$

$$n\text{-sınıflı Entropi}(S) \Rightarrow E(S) = \sum_{i=1}^n -p_i \cdot \log_2 p_i$$

### Örnek 1:

Örnek Sayısı	Örnek Sınıfı
9	+
5	-

$p_+ = 9/(9+5) = 9/14$  (aynı eğitim sınıfında bulunan bir örneğin + sınıfta bulunma olasılığı)

$p_- = 5/14$  (aynı eğitim sınıfında bulunan bir örneğin - sınıfta bulunma olasılığı)

$E = -9/14 * \log_2(9/14) - 5/14 * \log_2(5/14) = 0.94$

## Kazanım

Karar ağaçlarında kökler, düğümler ve yapraklar kazanım değerine göre oluşturulur.

$$\text{Gain}(S,A) = E(\text{before}) - G(\text{after\_splitting})$$

↑      ↑      ↑  
eğitim örneği   özellikler   S örneği için hesaplanmış olan orjinal entropi değeri

### Örnek 2:

Aşağıdaki tabloyu eğitim seti olarak kabul edip, tabloyu ID3 karar ağacı algoritmasına göre sınıflandırınız:

Haftasonu	Hava Durumu	Ebeveyn Durumu	Para Durumu	Karar(Sınıf)
H1	Güneşli	Evet	Zengin	Sinema
H2	Güneşli	Hayır	Zengin	Tenis
H3	Rüzgarlı	Evet	Zengin	Sinema
H4	Yağmurlu	Evet	Fakir	Sinema
H5	Yağmurlu	Hayır	Zengin	Ev
H6	Yağmurlu	Evet	Fakir	Sinema
H7	Rüzgarlı	Hayır	Fakir	Sinema
H8	Rüzgarlı	Hayır	Zengin	Alışveriş
H9	Rüzgarlı	Evet	Zengin	Sinema
H10	Güneşli	Hayır	Zengin	Tenis

**Not:** ID3 algoritması WEKA aracında C4.5 olarak gösterilmektedir.

Öncelikle tüm tabloyu en doğru şekilde ikiye bölecek olan özellik seçilmelidir. Bunun için de en yüksek kazanım veren özellik belirlenmelidir.

10 adet eğitim örneği için değerler şu şekilde bölünmektedir.

\* 6 adet Sinema

\* 2 adet Tenis

\* 1 adet Ev

\* 1 adet Alışveriş

Başlangıç için bu değerler üzerinden Entropi değeri hesaplanmalıdır.

$$E(S) = -((6/10)*\log_2(6/10) + (2/10)*\log_2(2/10) + (1/10)*\log_2(1/10) + (1/10)*\log_2(1/10))$$

$E(S) = 1.571$  (bu değer başlangıç entropi değeri olarak Information Gain'i hesaplamak için kenarda tutulacak.)

Tek tek tüm özelliklerin kazanım değerleri hesaplanarak en yüksek kazanım değerine sahip olan özellik kök düğümü olarak seçilir:

$$\text{Gain}(S, \text{Hava Durumu}) = ?$$

$$\text{Güneşli} = 3 \text{ (1 Sinema + 2 Tenis)}$$

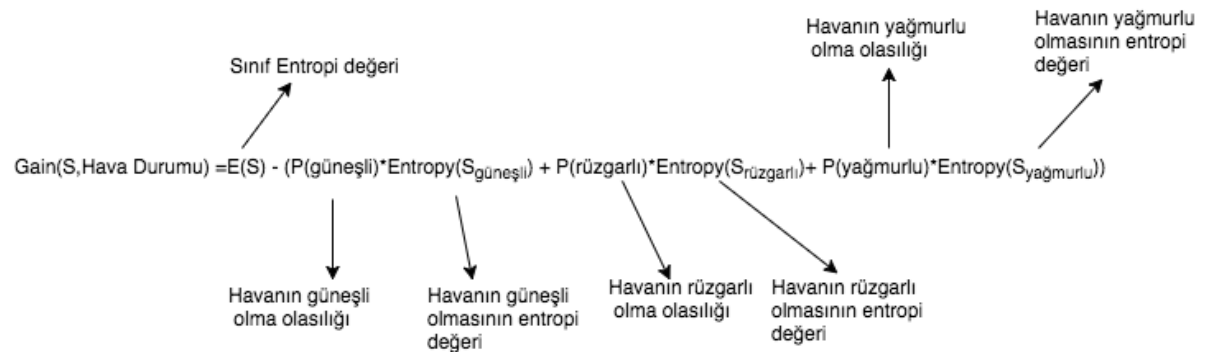
$$\text{Rüzgarlı} = 4 \text{ (3 Sinema + 1 Alışveriş)}$$

$$\text{Yağmurlu} = 3 \text{ (2 Sinema + 1 Ev)}$$

$$\text{Entropy}(S_{\text{güneşli}}) = - (1/3)*\log_2(1/3) - (2/3)*\log_2(2/3) = 0.918$$

$$\text{Entropy}(S_{\text{rüzgarlı}}) = - (3/4)*\log_2(3/4) - (1/4)*\log_2(1/4) = 0.811$$

$$\text{Entropy}(S_{\text{yağmurlu}}) = - (2/3)*\log_2(2/3) - (1/3)*\log_2(1/3) = 0.918$$



$$\text{Gain}(S, \text{Hava Durumu}) = 1.571 - (((1+2)/10)*0.918 + ((3+1)/10)*0.811 + ((2+1)/10)*0.918)$$

$$\text{Gain}(S, \text{Hava Durumu}) = 0.70$$

$$\text{Gain}(S, \text{Ebeveyn}) = ?$$

$$\text{Evet} = 5 \text{ (5 adet Sinema)}$$

$$\text{Hayır} = 5 \text{ (2 adet Tenis + 1 adet Sinema + 1 adet Alışveriş + 1 adet Ev)}$$

$$\text{Entropy}(S_{\text{evet}}) = - (5/5) * \log_2(5/5) = 0$$

$$\text{Entropy}(S_{\text{hayır}}) = -(2/5) * \log_2(2/5) - 3 * (1/5) * \log_2(1/5) = 1.922$$

$$\text{Gain}(S, \text{Ebeveyn}) = \text{Entropy}(S) - (P(\text{evet}) * \text{Entropy}(S_{\text{evet}}) + P(\text{hayır}) * \text{Entropy}(S_{\text{hayır}}))$$

$$\text{Gain}(S, \text{Ebeveyn}) = 1.571 - ((5/10) * \text{Entropy}(S_{\text{evet}}) + (5/10) * \text{Entropy}(S_{\text{hayır}}))$$

$$\text{Gain}(S, \text{Ebeveyn}) = 0.61$$

$$\text{Gain}(S, \text{Para}) = ?$$

$$\text{Zengin} = 7 \text{ (3 Sinema + 2 Tenis + 1 Alışveriş + 1 Ev)}$$

$$\text{Fakir} = 3 \text{ (3 Sinema)}$$

$$\text{Entropy}(S_{\text{zengin}}) = 1.842$$

$$\text{Entropy}(S_{\text{fakir}}) = 0$$

$$\text{Gain}(S, \text{Para}) = \text{Entropy}(S) - (P(\text{zengin}) * \text{Entropy}(S_{\text{zengin}}) + P(\text{fakir}) * \text{Entropy}(S_{\text{fakir}}))$$

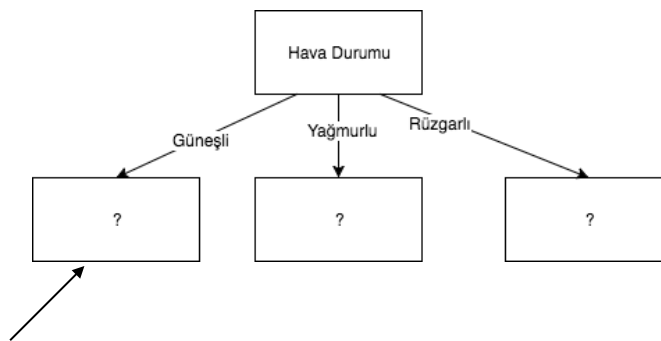
$$\text{Gain}(S, \text{Para}) = 0.2816$$

Son durumda tüm kazanım değerleri alt alta sıralanıp içlerinden en yüksek kazanım değerine sahip olan özellik kök düğüm olarak seçilir:

$$\text{Gain}(S, \text{Hava Durumu}) = 0.70$$

$$\text{Gain}(S, \text{Ebeveyn}) = 0.61$$

$$\text{Gain}(S, \text{Para}) = 0.2816$$



Hafta sonu	Hava	Ebeveyn	Para	Karar (Sınıf)
H1	Güneşli	Evet	Zengin	Sinema
H2	Güneşli	Hayır	Zengin	Tenis
H10	Güneşli	Hayır	Zengin	Tenis

Yukarıdaki örnekte görüldüğü gibi kök seçildikten sonra özelliğe ait tekil değerlerin her biri yaprak olarak belirlenmiş ve bu kez veri seti bu tekil değerler ile özelleştirilmiştir. Yani yukarıdaki tablo hava durumunun “Güneşli” olması durumunda diğer özelliklerin alabileceği

değerleri gösteren ayrı bir veri setine dönüştürülmüştür. Böylelikle hava durumu güneşli olduğunda bir alt yaprağın hangi özelliğe ait olacağı bu tabloya göre aşağıdaki gibi bulunacaktır:

\* 1 adet Sinema

\* 2 adet Tenis

$$\text{Entropy}(S_{\text{güneşli}}) = -(1/3) \cdot \log_2(1/3) - (2/3) \cdot \log_2(2/3) = 0.918$$

$$\text{Gain}(S_{\text{güneşli}}, \text{Ebeveyn}) = ?$$

$$\text{Gain}(S_{\text{güneşli}}, \text{Ebeveyn}) = \text{Entropy}(S_{\text{güneşli}}) - (P(\text{evet} | S_{\text{güneşli}}) \cdot \text{Entropy}(S_{\text{evet}}) + P(\text{hayır} | S_{\text{güneşli}}) \cdot \text{Entropy}(S_{\text{hayır}}))$$

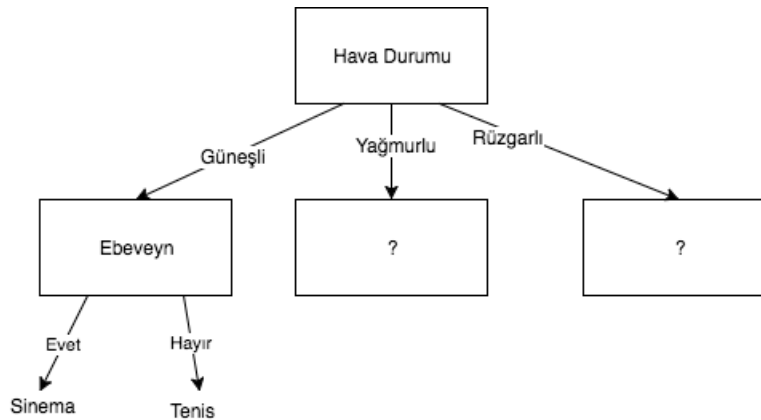
$$\text{Entropy}(S_{\text{evet}}) = -(1/3) \cdot \log_2(0) - 0 = 0$$

$$\text{Entropy}(S_{\text{hayır}}) = -(2/3) \cdot \log_2(0) - 0 = 0$$

$$\text{Gain}(S_{\text{güneşli}}, \text{Ebeveyn}) = 0.918 - ((1/3) \cdot 0 + (2/3) \cdot 0) = \mathbf{0.918}$$

$$\text{Gain}(S_{\text{güneşli}}, \text{Para}) = 0.918 - ((3/3) \cdot 0.918 + (0/3) \cdot 0) = \mathbf{0}$$

Ebeveyn özelliğinin kazanım değeri daha fazla olduğu için Güneşli durumunun yaprağı Ebeveyn olacaktır. Güneşli durumu için oluşturulan veri setine göre de eğer Evet ise Sinema kararı, hayır ise de Tenis kararı çıkacaktır:



Rüzgarlı ve Yağmurlu durumları için de aynı işlemler yapılarak, son durumda karar ağacı tamamlanacaktır.

**Overfitting:** Öğrenmenin ezberlemeye dönüşmesidir. Başarım grafiğinde belli bir andan sonra başarımın yükselmesinin durup, düşmeye başlamasıdır.