

103

- Her özdeşliğin bilgi kazancını hesapla
- Karar ağacı döğümüne seailen özdeşliğin değeriyle genişlet ↓

◊ Başta sınıflandırılmamış örnek vardı mı?

- Karar özdeşlikler ile devam et

Bilgi kazancı (Gain) :

- Bilgi kazancı ne kadar yüksekse, bir özdeşlik o kadar iyi bir bölünme kriteri olarak kabul edilir

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

+ ↓
ism veriseli Feature

Entropy = Değündeki veri portolodanın sınıflarının ne kadar karışık olduğunu ölçen bir metriktir.

E ↓ Homojenite ↑ it ↑ itetetejen ↑

* Karışık olursa bilgi kazancı yüksek olur.

$$E(S) = \sum_{i=1}^C -P_i \log_2 P_i$$

Play Golf	
Yes	No
9	5

$$\begin{aligned} \text{Entropy}(S, G) &= -Y \log_2 Y - N \log_2 N \\ &= -0.36 \cdot \log_2 0.36 - (0.64 \log_2 0.64) \\ &= 0.917 \end{aligned}$$

$Y = \frac{9}{14} = 0.64$ $N = \frac{5}{14} = 0.36$

[max entro, min entro]
[0, 1]

Entropisi yüksek bir değer
nerdeyse 1 olacaktır

CART

- İter öz niteliğin Gini impurity değerini, En düşük Gini impurity değerine sahip öz niteliği seç
- Karar ağacı düğümüne seçilen öz niteliğin değeriyle genişlet
- Başa sınıflandırılmamış örnek kaldı mı?
- Kaldı öz niteliklerle devam et

Gini Impurity (Gini belirsizliği) 1 Başta bir karar ağacı algoritması olan CART tarafından kullanılan bir ölçüdür. Verilen sınıfların dağılımının ne kadar karışık olduğunu [0, 0.5] aralığında karışıkta olabilir (2 etiket varsa) [0, 1] Modet etiket varsa gini değeri bu olur. Gini ne kadar olursa o kadar homojenizedir.

$$\text{Gini Impurity} = 1 - \sum_{i=0}^n p_i^2$$

* Gini ve Entropy karar ağaçlarında node 'lar' belirlemek için kullandığımız karmaşıklık ölçütleridir. Karmaşıklığın en az olmasını tercih ederiz ama karmaşıklık en fazla olanı node olarak seçeriz.

Örnek

Heart Disease	
Yes	No
105	39

144

$$\text{Gini} = 1 - \left[\left(\frac{105}{144} \right)^2 + \left(\frac{39}{144} \right)^2 \right]$$

$$= 0.4 \rightarrow \text{belirsizlik yarıya}$$