

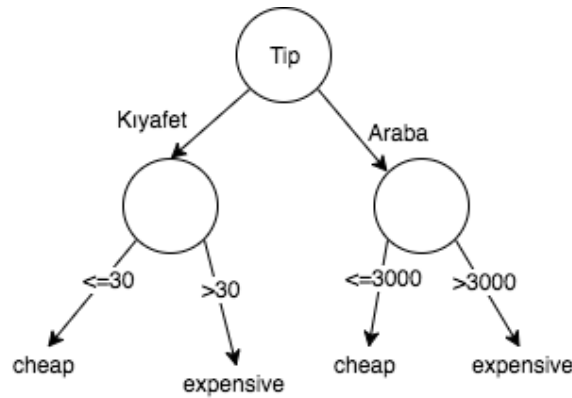
Machine Learning 101

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:

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

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

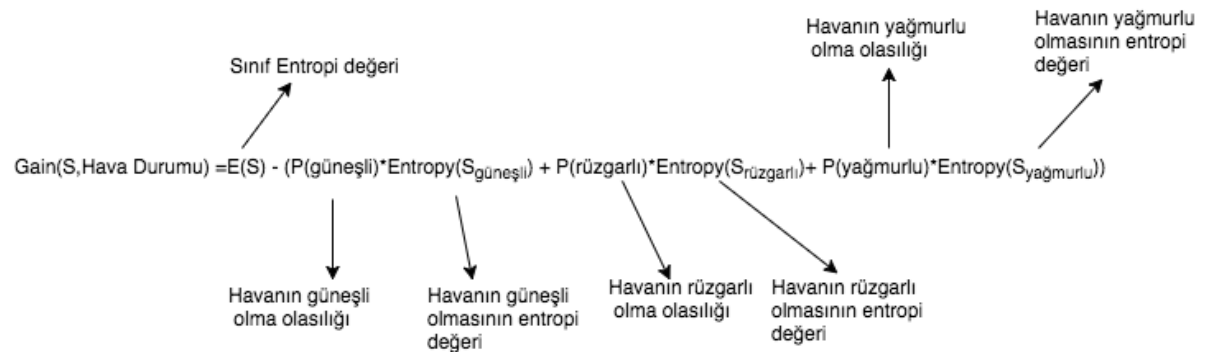
$$\text{Rüzgarlı} = 4 \text{ (3 Sinema + 1 Alışveriş)}$$

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

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

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

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



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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

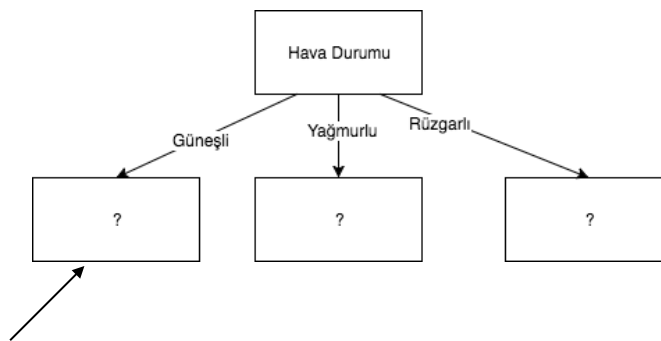
$$\text{Gain}(S, \text{Para}) = 0.2816$$

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

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

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

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

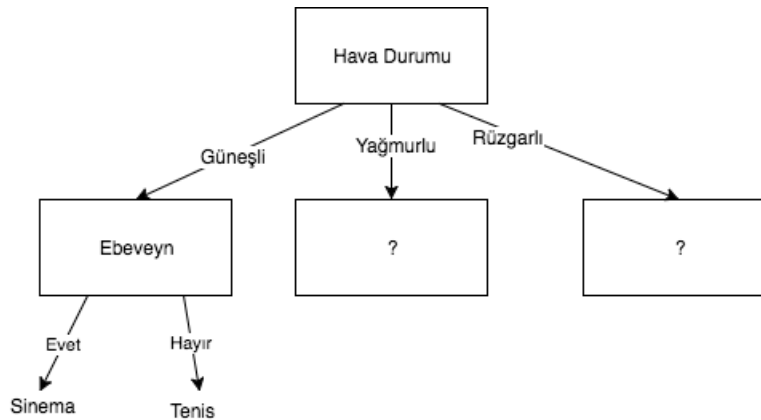
$$\text{Entropy}(S_{\text{evet}}) = -(1/3) \cdot \log_2(0) - 0 = 0$$

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

$$\text{Gain}(S_{\text{güneşli}}, \text{Ebeveyn}) = 0.928 - ((1/3) \cdot 0 + (2/3) \cdot 0) = \mathbf{0.928}$$

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

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.
- Eğitim verisi için çok başarılı sonuç verirken, test verisi için başarımın düşük olması overfitting örneğidir.

- Elimizdeki örnekler içerisinde parazitli bir örnek varsa, overfitting'e sebep olabiliyor.
- Ağaç çok büyüdüğünde neredeyse her eğitim verisi için bir sonuç üretilmiş oluyor ve bu da overfitting'e neden oluyor.

Overfitting'den kaçınmak için ağacı ya büyürken ya da ağaç oluştuktan sonra budama yoluna gitmek gerekiyor.

Koşullu Olasılık

Koşullu olasılık A durumu olduğunda B durumunun oluşma olasılığını hesaplamamı sağlar ve bu sayede konuyu spesifik hale getirmemi sağlamış oluyor. Gösterimi $P(B|A)$ şeklindedir.

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Bayes Kuralı

Bayes kuralı bir rassal değişken için olasılık dağılımı içinde koşullu olasılıklar ile marjinal olasılıklar arasındaki ilişkiyi gösterir. Formülü aşağıdaki gibidir:

$$P(B|A) = \frac{P(A|B) * P(B)}{P(A)}$$

Diagram annotations:

- $P(B|A)$ is labeled "posterior probability".
- $P(A|B)$ is labeled "likelihood".
- $P(B)$ is labeled "prior probability".
- $P(A)$ is labeled "evidence".
- Arrows point from "class" to $P(B|A)$ and from "attribute" to $P(A|B)$.

Not: Terimlerin her yerde ingilizceleri görüleceği için açıklamalarında ingilizceleri kullanılmıştır.

Maximum A Posteriori (MAP)

Bilinmeyen bir kararın tahmini için kullanılır. Tüm verilen veri için tüm hipotezler arasında hangi sınıfın olasılığı daha yüksekse o sınıf seçilir.

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D) = \frac{P(D|h) * P(h)}{P(D)} = P(D|h) * P(h)$$

↓

hepsinde sabit olduğu için ihmal ediliyor/hesaplanmıyor

Örnek 3:

$P(\text{sınıf}=+ | \text{kanser})=0.98$ (gerçekte kanser olup, test sonucu da kanser(+) çıkan hasta oranı)

$P(\text{sınıf}=- | \text{kanser})=0.97$ (gerçekte kanser olmayan ve test sonucu kanser değil(-) çıkan hasta oranı)

$P(\text{kanser}) = 0.008$ (kanser hasta oranı)

Yukarıda verilen değerlere göre test sonucu pozitif yani kanser çıkan bir hasta gerçekte kanser sınıftan mıdır değil midir?

$P(? | \text{sınıf}=+)=?$

Çözüm:

Bir olasılığın toplamı her zaman 1 olmalıdır bu yüzden;

$P(\text{kanser}) = 1 - P(\text{kanser}) \Rightarrow P(\text{kanser}) = 1 - 0.008 = 0.992$

$$P(\text{kanser} | \text{sınıf} = +) = \frac{P(\text{sınıf} = + | \text{kanser}) * P(\text{kanser})}{P(\text{sınıf} = +)}$$

$P(\text{kanser} | \text{sınıf}=+) = 0.98 * 0.008 / P(\text{sınıf}=+) = \mathbf{0.008}$

$$P(\text{kanser} | \text{sınıf} = +) = \frac{P(\text{sınıf} = + | \text{kanser}) * P(\text{kanser})}{P(\text{sınıf} = +)}$$

$P(\text{kanser} | \text{sınıf}=+) = 0.03*0.992/P(\text{sınıf}=+) = \mathbf{0.02976}$

sonuç: kanser

Naive Bayes Teoremi

Naive Bayes Teoremi tüm özellikleri birbirlerinden bağımsız kabul eder. Metin sınıflandırma başarısı oldukça yüksektir.

$$V_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) * \prod_{i=1}^n P(x_i | v_j)$$

Örnek 4:

Aşağıda verilen tablodaki veri setine göre Dergi Promosyonuna Evet, Saat Promosyonuna Evet, Hayat Sigortasına Hayır, Kredi Kartı Sigortasına Hayır cevabını vermiş bir kişinin cinsiyeti nedir?

Dergi Promosyonu	Saat Promosyonu	Hayat Sigortası	Kredi Kartı Sigortası	Cinsiyet
Evet	Hayır	Hayır	Hayır	Erkek
Evet	Evet	Evet	Evet	Kadın
Hayır	Hayır	Hayır	Hayır	Erkek
Evet	Evet	Evet	Evet	Erkek
Evet	Hayır	Evet	Hayır	Kadın
Hayır	Hayır	Hayır	Hayır	Kadın
Evet	Evet	Evet	Evet	Erkek
Hayır	Hayır	Hayır	Hayır	Erkek
Evet	Hayır	Hayır	Hayır	Erkek
Evet	Evet	Evet	Hayır	Kadın

Çözüm:

Soruyu kısaltırsak eğer; (Evet, Evet, Hayır, Hayır) -> Cinsiyet=?

Kontrol edilmesi gereken 2 hipotez var elimizde kadın mı yoksa erkek mi? Bunun için her ikisi için de sonucu hesaplayıp hangisinin değeri daha büyükse o cinsiyettir cevabını vereceğiz.

Hipotez 1:

$$P(cins = erkek|X) = \frac{P(X|cins = erkek) * P(cins = erkek)}{P(X)}$$

P(X | cins=erkek)=? cevabı için tüm olasılıklar ayrı ayrı bulunarak çarpılacak.

P(dergi_prom=evet | cins=erkek) (cinsiyeti erkek olup dergi promosyonuna evet deme olasılığı)
= 4/6 (cinsiyeti erkek olup evet diyenler/cinsiyeti erkek olanların toplam sayısı)

P(saat_prom=evet | cins=erkek) = 2/6

P(hayat_sig=hayır | cins=erkek) = 4/6

P(kredi_kart=hayır | cins=erkek)=4/6

P(X | cins=erkek) = P(dergi_prom=evet | cins=erkek)*P(saat_prom=evet | cins=erkek)*P(kredi_kart=hayır | cins=erkek)*P(hayat_sig=hayır | cins=erkek)

$$P(X | \text{cins}=\text{erkek}) = (4/6) * (2/6) * (4/6) * (4/6) = 0.0988$$

$$P(\text{cins}=\text{erkek}) = 0.6$$

$$P(\text{cins}=\text{erkek} | X) = 0.0988 * 0.6 = 0.0593$$

Hipotez 2:

$$P(\text{cins} = \text{kadın} | X) = \frac{P(X | \text{cins} = \text{kadın}) * P(\text{cins} = \text{kadın})}{P(X)}$$

$$P(X | \text{cins}=\text{kadın}) = (3/4) * (2/4) * (1/4) * (3/4) = 0.07$$

$$P(\text{cins}=\text{kadın}) = 0.4$$

$$P(\text{cins}=\text{kadın} | X) = 0.0281$$

0.0593 > 0.0281 olduğu için Bayes sınıflandırıcı cinsiyetin erkek olma olasılığının daha yüksek olduğunu söylüyor.

Not: Örneğin elektronik posta sınıflandırmada her elektronik posta içerisindeki kelimeler üzerinden naive bayes ile bir sınıflandırma yapıp postanın spam mi, gerçek posta mı olduğuna karar veriyorsak, eğitim seti içerisinde bulunmayan bir kelime ile karşılaştığımızda karar aşamasında o değerın eğitim setindeki sayısı 0 olduğu için tüm sonucu çarpımdan ötürü 0'a çekecektir. Bu durumda **Laplacian(1-up) Smoothing** kullanılabilir:

$$P(x_i | v_j) = \frac{n_i + 1}{n + |\text{vocab}|}$$

Vj sınıfına ait olan xi kelimelerinin sayısı

↑

n_i + 1

↓

n + |vocab|

↓ ↓

postadaki sınıfı Sözlükteki
olan kelimelerin kelimelerin
sayısı sayısı

Örnek 5:

Class A: “The cat crabs the crolls off the stairs”

Class B: “It is raining cats and dogs”

Yukarıdaki sınıflara göre X=“cats eat mice and dogs bunny bones” metnini sınıflandırın.

Sözlük={cats, crab, croll, stair, rain, dog}

$P(V_j) = 1/2$ (2 adet sınıf olduğu için ya A ya B) | Sözlük|=6 $n_A=4$ $n_B=3$

X metninde cat ve dog kelimeleri sadece sözlükte bulunduğu için onlar üzerinden hareket edeceğiz:

$$P(x_i | v_j) = (n_{ij} + 1) / (n_j + 6)$$

1. Hipotez:

$$P(A | X) = P(A) * P(\text{cat} | A) * P(\text{dog} | A) = 1/2 * (1+1)/(4+6) * (0+1)/(4+6) = \mathbf{0.01}$$

2. Hipotez:

$$P(B | X) = P(B) * P(\text{cat} | B) * P(\text{dog} | B) = 1/2 * (1+1)/(3+6) * (1+1)/(3+6) = \mathbf{0.024}$$

$0.024 > 0.01$ olduğu için X dökümanı B sınıfına **aittir**.