

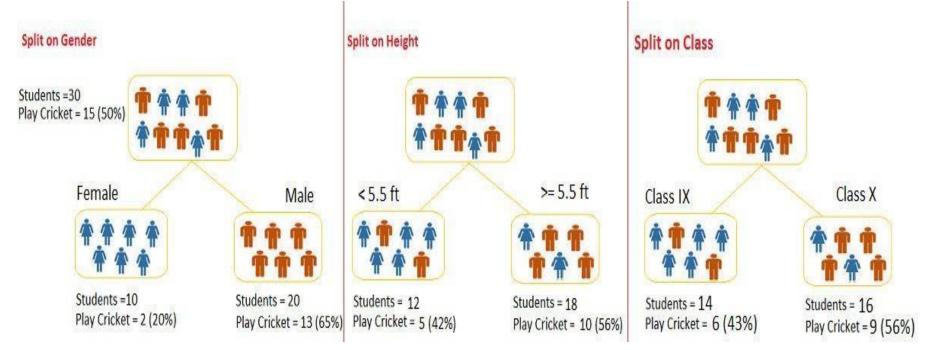
Makine Öğrenmesi ve İmge Tanıma

Karar Ağaçları- devam (2)

KARAR AĞAÇLARI

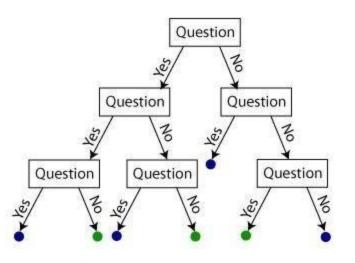


- Üç farklı özellikle tanımlanan 30 çocuğa bakalım.
 - Cinsiyet (E veya K),
 - Sınıf(IX. veya X.)
 - Boy (1.50 1.80 cm).
- Peki üç öznitelikten en belirgin/belirleyici olan girdi değişkeni hangisi?



BÖLÜMLEME KARARI





BAŞARILI - BAŞARISIZ

Araştırma: Gelir Dağılımındaki Eşitsizlik Değerleri (Gini Katsayıları)

1. Gini Endeksi

Gini: «Bir popülasyondan rastgele iki öğe seçildiğinde, aynı sınıfta olmaları gerekir» Bu durumda popülasyonlardan birinin saf olması olasılığı yüksektir.

"Başarı" (success) veya "Başarısızlık" (failure) gibi kategorik hedef değişkenle çalışır.

Yalnızca İkili bölümleme yapabilir.

Gini değeri arttıkça, homojenlik değeri de artar.

Gini yöntemi, CART (Sınıflandırma ve Regresyon Ağacı) tarafından ikili bölünmeler yapmak için kullanılır.



BÖLÜMLEME KARARI



1. Gini Endeksi

Başarı ve başarısızlık için olasılık karesinin toplamını hesapla, yani (p²+q²) Şimdi oluşturulan bölümün her bir düğümünün ağırlıklı Gini bölünümü için Gini'yi hesaplayın.



BÖLÜMLEME KARARI

Gini değeri arttıkça, homojenlik değeri de artar.



1. Gini Endeksi

Başarı ve başarısızlık için olasılık karesinin toplamını hesapla, yani (p²+q²) Şimdi oluşturulan bölümün her bir düğümünün ağırlıklı Gini bölünümü için Gini'yi hesaplayın.



BÖLÜMLEME KARARI



1. Gini Endeksi

Başarı ve başarısızlık için olasılık karesinin toplamını hesapla, yani (p²+q²) Şimdi oluşturulan bölümün her bir düğümünün ağırlıklı Gini bölünümü için Gini'yi hesaplayın.

```
elma muz brokoli

Adet = 3 3 6

p = 3/12 3/12 6/12

= 1/4 1/4 1/2

GI = 1 - [(1/4)^2 + (1/4)^2 + (1/2)^2]

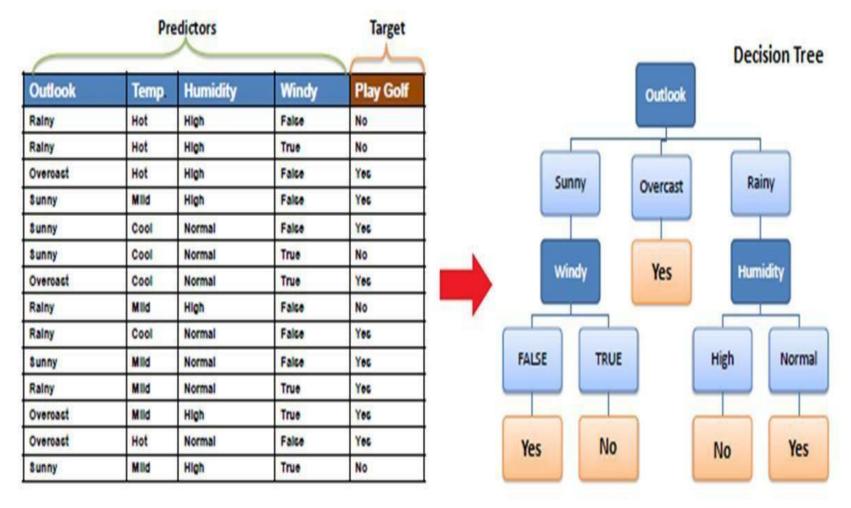
= 1 - [1/16 + 1/16 + 1/4]

= 1 - 6/16

= 10/16

= 0.625
```







	Predictors			
Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	Falce	No
Rainy	Hot	High	True	No
Overoact	Hot	High	Falce	Yes
Sunny	Mild	High	Falce	Yes
Sunny	Cool	Normal	Falce	Yes
Sunny	Cool	Normal	True	No
Overoact	Cool	Normal	True	Yes
Rainy	Mild	High	Falce	No
Rainy	Cool	Normal	Falce	Yes
Sunny	Mild	Normal	Falce	Yes
Rainy	Mild	Normal	True	Yes
Overoact	Mild	High	True	Yes
Overoact	Hot	Normal	Falce	Yes
Sunny	Mild	High	True	No

			play				
		yes		no		total	
	sunny		3		2		5
Outlook	overcast		4		0		4
	rainy		2		3		5
							14



		Pr	edictors		Target							
y		$\overline{}$)				play		
	Outlook	Temp	Humidity	Windy	Play Golf				yes	no	total	
	Rainy	Hot	High	Falce	No	⇔		sunny		3	2	5
 	Rainy	Hot	High	True	No	\alpha\rightarrow	Outlook	overcast		4	0	4
	Overoact	Hot	High	Falce	Yes		- Catioon	rainy		2	3	5
	Sunny	Mild	High	Falce	Yes			lality		2	J	
	Sunny	Cool	Normal	Falce	Yes]						14
	Sunny	Cool	Normal	True	No]	Gini(Outlook=Rainy)=					
	Overoact	Cool	Normal	True	Yes		Gini	Outloo	K=Kai	ny)=		
=	Rainy	Mild	High	Falce	No	⇔	1_	$(2/5)^2$ -	(3/5)	2 = 1 - 0	14 – n	.36 = 0.48
=	Rainy	Cool	Normal	Falce	Yes	\Diamond	- '	(2/0)	(0/0)	-1 ().IO U	.50 - 0.40
, i	Sunny	Mild	Normal	Falce	Yes							
⇒	Rainy	Mild	Normal	True	Yes	\Diamond						
	Overoact	Mild	High	True	Yes]			(1)[DİKKAT	(1)	
	Overoast	Hot	Normal	Falce	Yes				(:/-		(•)	
	Sunny	Mild	High	True	No	1			Cin	i Impu	rity - 1	l Cini
,		- 6	70		X	-			GII	i Impu	IILy = 1	L-GIIII



	Predictors				Target						
		_	_			1				play	
	Outlook	Temp.	Humidity	Windy	Play Golf				yes	no	total
	Rainy	Hot	High	Falce	No			sunny		3 2	. 5
	Rainy	Hot	High	True	No		Outlook	overcast		4 () 4
-	Overoact	Hot	High	Falce	Yes			rainy		2 3	
	Sunny	Mild	High	Falce	Yes			Talliy		2 3	
	Sunny	Cool	Normal	Falce	Yes	I					14
	Sunny	Cool	Normal	True	No		C:n:	/O+l.a.a.l	L. Dain	۸	
•	Overoact	Cool	Normal	True	Yes		Gini	(Outlool	k=Rainy	/)=	
	Rainy	Mild	High	Falce	No	'	1-	$(2/5)^2$ –	$(3/5)^2 =$: 1 - 0.16	- 0.36 = 0.48
	Rainy	Cool	Normal	Falce	Yes	ľ.		(2/0)	(0/0)	1 0.10	0.00 - 0.40
	Sunny	Mild	Normal	Falce	Yes	I	GinifO	hutlook-O	voreact)	-1-(1/1	$(0/4)^2 = 0$
	Rainy	Mild	Normal	True	Yes		Ollillo	JULIOUK-U	vercastj	- 1 - (4/4) - (0/4) - 0
-	Overoact	Mild	High	True	Yes	I	Gini	i(Outloo	k-Sunr	\\\\-	
+	Overoact	Hot	Normal	Falce	Yes			•			- 0.16 = 0.48
\	Sunny	Mild	High	True	No	I	1-	(3/3)	(2/5)-=	1 - 0.30	- 0.10 = 0.46
		76	X.		X						



		Pre	edictors		Target							
V.		_							p	lay		
	Outlook	Temp.	Humidity	Windy	Play Golf				yes	no	total	
	Rainy	Hot	High	Falce	No			sunny	3	2	5	
	Rainy	Hot	High	True	No	Outl	ook	overcast	4	0	4	
-	Overoact	Hot	High	Falce	Yes			rainy	2	2	5	
	Sunny	Mild	High	Falce	Yes			laniy			14	
	Sunny	Cool	Normal	Falce	Yes						14	
	Sunny	Cool	Normal	True	No		ini/	Outlook	-Dainy	_		
•	Overoact	Cool	Normal	True	Yes	ا	11111	Outlook	-Kalliy,)=		
	Rainy	Mild	High	Faice	No		1 – 1	$(2/5)^2 - 1$	$3/5)^2 =$	1 - 0.16	- 0.36 = 0.48	
	Rainy	Cool	Normal	Falce	Yes		'	(2/0)	0/0) -	. 0.10	0.00 - 0.40	
Î	Sunny	Mild	Normal	Falce	Yes	Gi	niſ∩	utlook-Ov	oroact) -	1 _ (1/1)	$(0/4)^2 = 0$	
	Rainy	Mild	Normal	True	Yes	UI UI	UJIII	utiouk-uv	ercastj –	1- (4/4)	(0/4) - 0	
	Overoact	Mild	High	True	Yes	ے ا	ini	(Outloo	k=Sunn	v)=		
+	Overoact	Hot	Normal	Falce	Yes			•			- 0.16 = 0.48	
\	Sunny	Mild	High	True	No		1-	(3/3) - (2/0) -	1 - 0.30	- 0.10 - 0.40	

Ağırlıklı hesap

Gini(Outlook) = $(5/14) \times 0.48 + (4/14) \times 0 + (5/14) \times 0.48 = 0.171 + 0 + 0.171 = 0.342$



	Pro	edictors	Target	
Outlook	Temp.	Numidity	Windy	Play Golf
Rainy	Hot	High	Falce	No
Rainy	Hot	High	True	No
Overcast	Hot	High	Falce	Yes
Sunny	Mild	High	Falce	Yes
Sunny	Cool	Norma	Falce	Yes
Sunny	Cool	Normal	True	No
Overoast	Cool	Normal	True	Yes
Rainy	Mild	High	Falce	No
Rainy	Cool	Norma	Falce	Yes
Sunny	Mild	Normal	Falce	Yes
Rainy	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overoast	Hot	Nomal	Falce	Yes
Sunny	Mild	High	True	No

Temperature	Yes	No	Number of instances
Hot	2	2	4
Cool	3	1	4
Mild	4	2	6

Gini(Temp=Hot) = $1 - (2/4)^2 - (2/4)^2 = 0.5$

Gini(Temp=Cool) = $1 - (3/4)^2 - (1/4)^2 = 1 - 0.5625 - 0.0625 = 0.375$

Gini(Temp=Mild) = $1 - (4/6)^2 - (2/6)^2 = 1 - 0.444 - 0.111 = 0.445$

Gini(Temp) = (4/14) x 0.5 + (4/14) x 0.375 + (6/14) x 0.445 = 0.142 + 0.107 + 0.190 = 0.439

Gini(Outlook) = $(5/14) \times 0.48 + (4/14) \times 0 + (5/14) \times 0.48 = 0.171 + 0 + 0.171 = 0.342$



	Predictors					
Outlook	Temp	Humidity	Windy	Play Golf		
Rainy	Hot	High	Falce	No		
Rainy	Hot	High	True	No		
Overcast	Hot	High	Falce	Yes		
Sunny	Mild	High	Falce	Yes		
Sunny	Coo	Normal	Falce	Yes		
Sunny	Cool	Normal	True	No		
Overoact	Cocl	Normal	True	Yes		
Rainy	Mild	High	Faice	No		
Rainy	Coo	Normal	Falce	Yes		
Sunny	Mild	Normal	Falce	Yes		
Rainy	Mild	Normal	True	Yes		
Overcast	Mild	High	True	Yes		
Overoact	Hot	Normal	Falce	Yes		
Sunny	Mild	High	True	No		

Humidity	Yes	No	Number of instances
High	3	4	7
Normal	6	1	7

Gini(Humidity=High) = $1 - (3/7)^2 - (4/7)^2 = 1 - 0.183 - 0.326 = 0.489$

Gini(Humidity=Normal) = $1 - (6/7)^2 - (1/7)^2 = 1 - 0.734 - 0.02 = 0.244$

Gini(Humidity) = $(7/14) \times 0.489 + (7/14) \times 0.244 = 0.367$

Wind	Yes	No	Number of instances
Weak	6	2	8
Strong	3	3	6

Gini(Wind=Weak) = $1 - (6/8)^2 - (2/8)^2 = 1 - 0.5625 - 0.062 = 0.375$

Gini(Wind=Strong) = $1 - (3/6)^2 - (3/6)^2 = 1 - 0.25 - 0.25 = 0.5$

Gini(Wind) = (8/14) x 0.375 + (6/14) x 0.5 = 0.428

Gini(Temp) = $(4/14) \times 0.5 + (4/14) \times 0.375 + (6/14) \times 0.445 = 0.142 + 0.107 + 0.190 = 0.439$

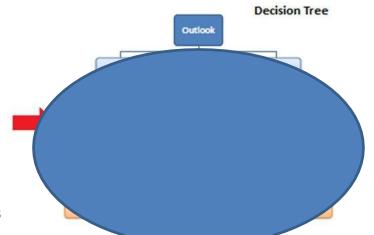
Gini(Outlook) = (5/14) x 0.48 + (4/14) x 0 + (5/14) x 0.48 = 0.171 + 0 + 0.171 = 0.342



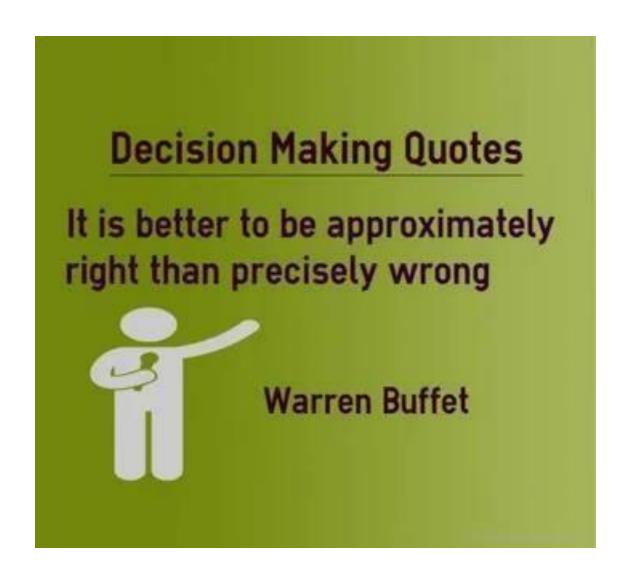
	larget			
Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	Falce	No
Rainy	Hot	High	True	No
Overoact	Hot	High	Falce	Yes
Sunny	Mild	High	Falce	Yes
Sunny	Cool	Normal	Falce	Yes
Sunny	Cool	Normal	True	No
Overoast	Cool	Normal	True	Yes
Rainy	Mild	High	Falce	No
Rainy	Cool	Normal	Falce	Yes
Sunny	Mild	Normal	Falce	Yes
Rainy	Mild	Normal	True	Yes
Overoast	Mild	High	True	Yes
Overoast	Hot	Normal	Falce	Yes
Sunny	Mild	High	True	No

Dradictors

Gini index	
0.342	
0.439	
0.367	
0.428	



Target

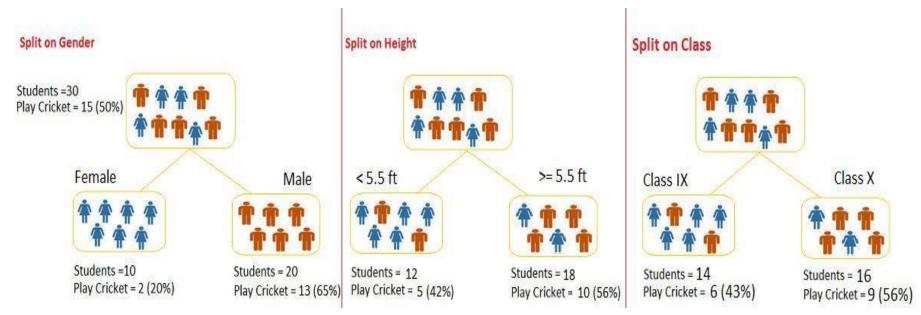




KARAR AĞAÇLARI



- Üç farklı özellikle tanımlanan 30 çocuğa bakalım.
 - Cinsiyet (E veya K),
 - Sınıf(IX. veya X.)
- − Boy (1.50 − 1.80 cm).
- Peki üç öznitelikten en belirgin/belirleyici olan girdi değişkeni hangisi?



KARAR AĞAÇLARI



- Üç farklı özellikle tanımlanan 30 çocuğa bakalım.
 - Cinsiyet (E veya K),
 - Sınıf(IX. veya X.)
- − Boy (1.50 − 1.80 cm).
- Peki üç öznitelikten en belirgin/belirleyici olan girdi değişkeni hangisi?

Cinsiyet:

Gini(Kız) = (0.2)*(0.2)+(0.8)*(0.8)=0.68

Gini(Erkek) = (0.65)*(0.65)+(0.35)*(0.35)=0.55

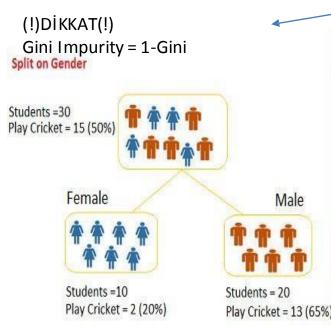
Gini (Cinsiyet)= (10/30)*0.68+(20/30)*0.55 = **0.59**

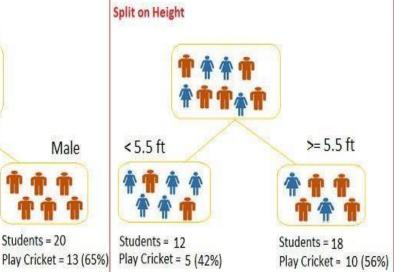
Sinif:

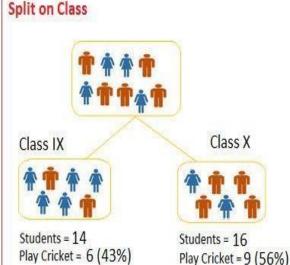
Gini(IX) = (0.43)*(0.43)+(0.57)*(0.57)=0.51

Gini(X) = (0.56)*(0.56)+(0.44)*(0.44)=0.51

Gini(Sınıf) = (14/30)*0.51+(16/30)*0.51 = 0.51







BÖLÜMLEME KARARI

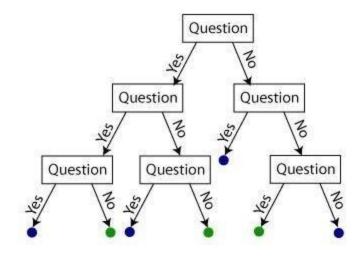


- 1. Gini Endeksi
- 2.Ki-Kare(Chi-Square):
- üst düğüm ve alt düğüm arasındaki istatistiksel farkı bulmak için kullanılır. Hedef değişkenin beklenen frekansları ile hedef değişkenin gözlemlenen/elde edilen frekansları arasındaki tüm farkların karelerinin toplamı ile hesaplanır.
- Ki-Karenin değeri arttıkça, ana düğüm ile alt düğüm arasındaki istatistiksel fark değeri de artar.
- Ki kare = ((Gerçek Beklenen)² / Beklenen)^{1/2}

3. Varyans

Variance =
$$\frac{\sum (X - \overline{X})^2}{n}$$

X⁻: X in ortalaması

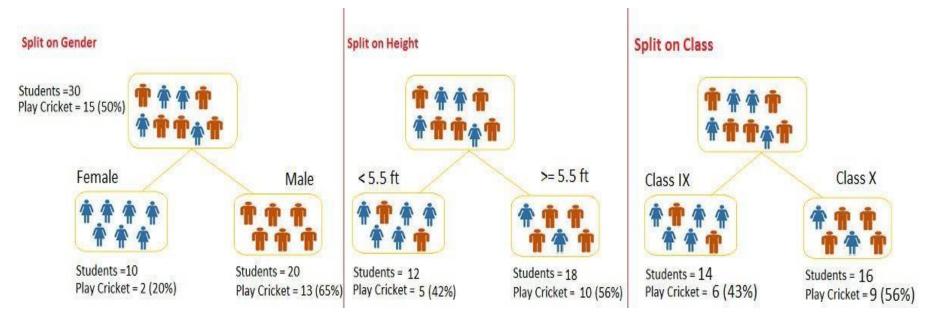


BAŞARILI - BAŞARISIZ

Kİ-KARE



- Üç farklı özellikle tanımlanan 30 çocuğa bakalım.
 - Cinsiyet (E veya K),
 - Sınıf(IX. veya X.)
 - *− Boy (1.50 − 1.80 cm).*
- Peki üç öznitelikten en belirgin/belirleyici olan girdi değişkeni hangisi?



Kİ-KARE

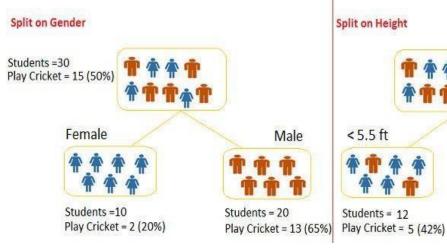


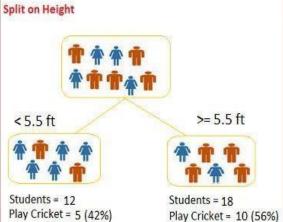
- Üç farklı özellikle tanımlanan
 30 çocuğa bakalım.
 - Cinsiyet (E veya K),
 - Sınıf(IX. veya X.)
 - Boy (1.50 1.80 cm).
- Peki üç öznitelikten en belirgin/belirleyici olan girdi değişkeni hangisi?

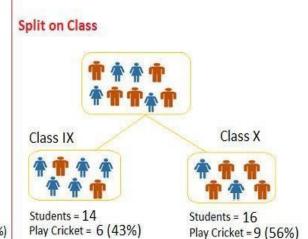
Node	Play Cricket	Not Play Cricket	Total	Expected Play Cricket	Expected Not Play Cricket		Deviation Not Play Cricket	Chi-Square			
								Play	Not Play		
				,	,	,	,	Cricket	Cricket		
Female	2	8	10	5	5	-3	3	1.34	1.34		
Male	13	7	20	10	10	3	-3	0.95	0.95		
							Total Chi-Square	e 4.58			

	Node	Play Cricket	Not Play Cricket	Total	Expected Play Cricket	Expected Not Play Cricket	I	Deviation Not Play Cricket	Chi-Square	
									Play	Not Play
L					,	,	,	,	Cricket	Cricket
	IX	6	8	14	7	7	-1	1	0.38	0.38
	X	9	7	16	8	8	1	-1	0.35	0.35
١	11 \ 1/2							Total Chi-Square	1.46	

Ki kare = ((Gerçek - Beklenen) 2 / Beklenen) 1/2







Kİ-KARE



Node	Play Cricket	Not Play Cricket	Total	Expected Play Cricket	Expected Not Play Cricket		Deviation Not Play Cricket	Chi-Square	
								Play	Not Play
	1			2	_	-	-	Cricket	Cricket
Female	2	8	10	5	5	-3	3	1.34	1.34
Male	13	7	20	10	10	3	-3	0.95	0.95
• Ki kare = ((Gerçek - Beklenen) ² / Beklenen) ^{1/2}							Total Chi-Square	4.58	

Node Play Cricket		Not Play		Expected	Expected Not	Deviation	Deviation Not	Chi-Square	
	Cricket	Total	Play Cricket				Play	Not Play	
								Cricket	Cricket
IX	6	8	14	7	7	-1	1	0.38	0.38
X	9	7	16	8	8	1	-1	0.35	0.35
-/2 E\2/ Poklonon\ %							Total Chi-Square	1.46	

=(2-5)²/ Beklenen) ^{1/2}

 $= (9/5)^{1/2} = 1.8 \% = 1.34$

Kara ağaçların AVANTAJLARI



Anlaması kolay:

Analitik altyapısı olmayan insanlar için bile, karar ağacı algoritmasının anlaşılması çok kolaydır. Bir kullanıcının ağaçları incelemesi, okuması ve yorumlaması için herhangi bir istatistiksel bilgi veya bilgiye sahip olması gerekmez. Kullanıcılar verileri kolayca okuyabilir.

Grafik **gösterimi** son derece sezgisel (bütünsel) ve kullanıcı dostudur .

Veri araştırmalarında faydalı:

Karar ağacı, en hızlısı olmasa bile, kesinlikle en önemli değişkenin tanımlanmasının en hızlı yönteminden olduğuna inanılmaktadır. Karar ağacı, kullanıcıların, özelliklerin yanı sıra yeni değişkenler oluşturmalarına yardımcı olabilir. Bu yeni özellikler, hedef değişkeni tahmin etmek için daha fazla güce sahip olacaktır. Veri arama aşamasında da kullanılabilir.

AVANTAJLARI



Daha az veri temizliği gerekiyor:

makine öğreniminde, bir kullanıcı zamanının çoğunu veri temizlemeye ve iyi verileri kötü verilerden ayırmaya harcamak zorundadır. Bununla birlikte, karar ağacı söz konusu olduğunda, bu süreç oldukça kolaydır ve çok zaman almaz. Uç değerlerin yanı sıra aykırı değerlerden etkilenmez ve böylece temizleme işlemi kolaylaşır.

Veri türü bir kısıtlama değil:

Karar ağacı çok yönlü bir algoritmadır ve kategorik ve sayısal veri değişkenlerini kolaylıkla işleyebilir.

Parametrik Olmayan Yöntem:

Parametrik olmayan bir yöntem, sınıflandırıcı yapıları veya uzamsal dağılım ile ilgili varsayımları olmayan bir yöntem anlamına gelir. Karar ağacı parametrik olmayan bir yöntemdir.

NOT: parametrik yöntemler sınırlı sayıda parametre alır parametrik olmayan yöntemlerde veri arttıkça parametre artablik olmayan 29

DEZAVANTAJI



Aşırı uyum (overfitting) gösterme:

Karar ağacı modellerinde, aşırı uyumun en yaygın ve karşılaşılan problemdir. Modeli eldeki veriyle son derece yüksek başarılı sonuç üretecek şekilde oluşturur ama hiç görmediği yeni bir veriye çok yüksek hata verir. (Soruları ezberleme örneği) Bununla birlikte, bu sorun budamanın yanı sıra model parametreleri üzerindeki kısıtlamalar kullanılarak çözülebilir.

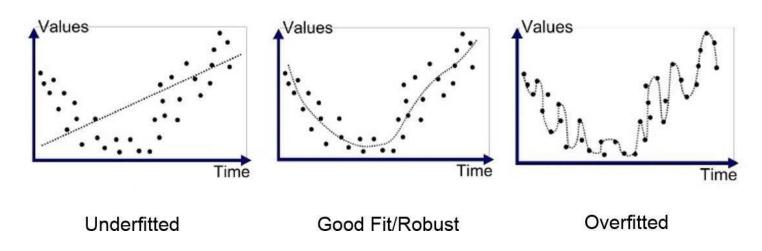
Sürekli değişkenler için uygun değil:

Sürekli değişkenlerle çalışabilse de, hiç uygun değildir. Karar ağacı, de değişkenleri gittikçe daha fazla sınıflandırmaya başladığında kategoride bilgileri kaybetmeye başlar.

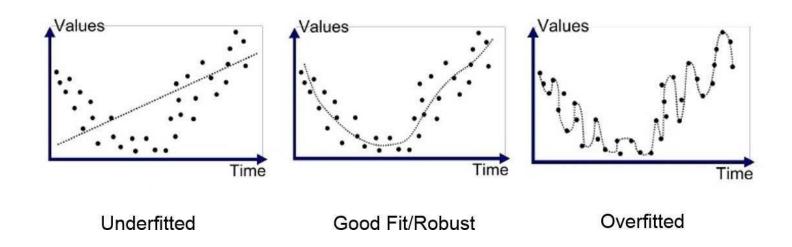
NOT: Bağımlı değişken sürekli olduğunda, regresyon ağaçları kullanılırken bağımlı değişken kategorik olduğunda; sınıflandırma ağaçları kullanılır

Aşırı uyum (overfitting) gösterme:

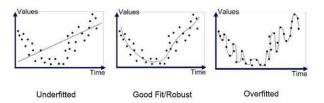
Karar ağacı modellerinde, aşırı uyumun en yaygın ve karşılaşılan problemdir. Modeli eldeki veriyle son derece yüksek başarılı sonuç üretecek şekilde oluşturur ama hiç görmediği yeni bir veriye çok yüksek hata verir. (Soruları ezberleme örneği) Bununla birlikte, bu sorun budamanın yanı sıra model parametreleri üzerindeki kısıtlamalar kullanılarak çözülebilir.



- Ağaç boyutuna ilişkin kısıt koyma
- Ağaç budaması



Ağaç boyutuna ilişkin kısıt koyma



- Ağaç budaması
- Düğüm bölümlemesi için Minimum Örnekler (minimum samples).
- Bölümlemek için gereken gözlem adedi tanımlanabilir.
- Bir yaprak için minimum örnek tanımlanabilir.
- Azınlık sınıfının çoğunluk sınıfı olduğu bölgeler çok sınırlı olduğundan, dengesiz sınıf problemleri için daha düşük değerler seçilir.

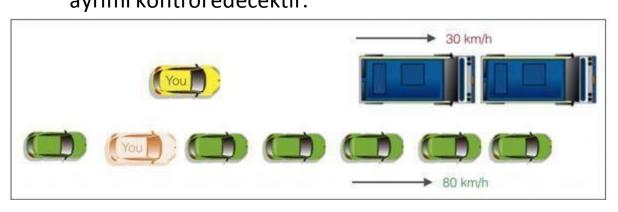
Ağacın Maksimum Derinliği (maximum depth):

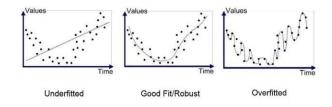
- Bu, bir ağacın derinliğini kontrol etmek için kullanılabilir.
- Çapraz kontrol (cross validation) kullanılarak ayarlanması gerekir.
- Çok sayıda terminal düğüm oluşur.

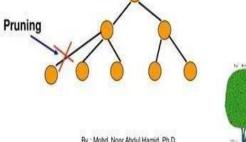
Ağaç budaması

Açgözlü yaklaşım (greedy)

Yöntem, verilen en iyi durma koşullarından biri elde edilene kadar sadece en iyi ayrımı kontrol edecektir.









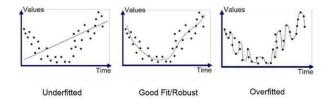
Gerçek hedef ne?

- -geçmek mi?
- -Mesafe almak mı?
- Maksimum yol kat etmek mi? PAÜ CENG420 Makine Öğrenmesi ve İmge Tanıma

budama olsaydı, yöntem birkaç adım geri alacak ve üzerinde işlem yapmadan önce durumu düşünme şansına sahip olacaktı.

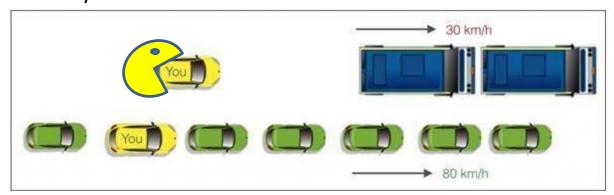
Ders Notları 2024 40

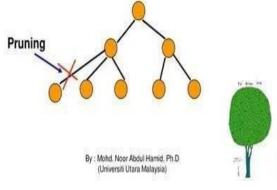
Ağaç budaması



Açgözlü yaklaşım (greedy)

Yöntem, verilen en iyi durma koşullarından biri elde edilene kadar sadece en iyi ayrımı kontrol edecektir.





Budamayı kullanmak için, büyük derinlikli karar ağacı üretilir.

Sonra alttan başlayıp yavaş ve kademeli olarak yukarıda, karşılaştırıldığında negatif sonuçlar veren yaprakları kesmeye başlarız.



PAÜ CENG420 Makine Öğrenmesi ve İmge Tanıma Ders Notları 2024