



# VERİ MADENCİLİĞİ

Fırat İsmailoğlu, PhD

Karar Ağaçları



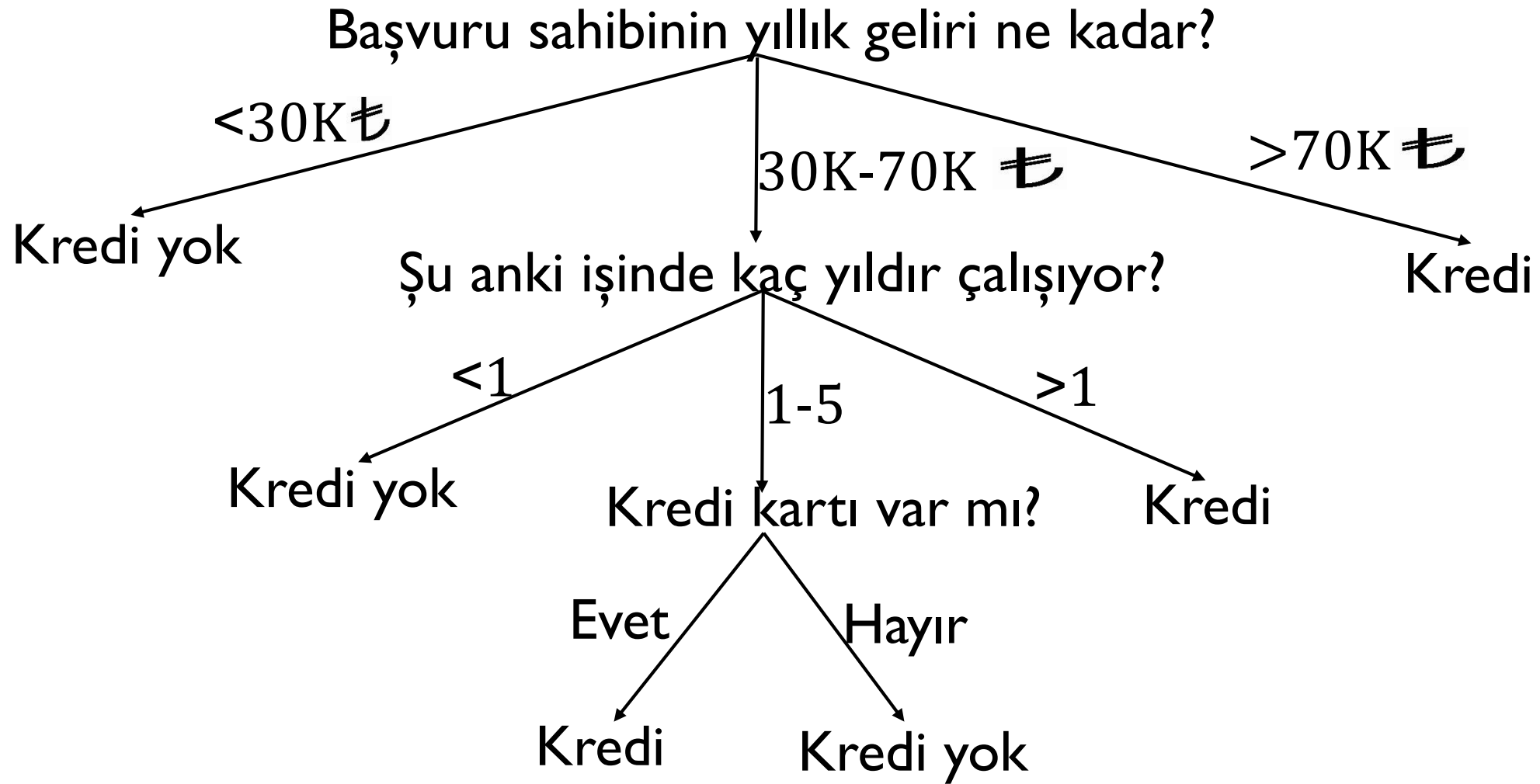


```
if vucutSicakligi==soguk{  
    return memeliDegil;  
else if doguruyorMu==evet{  
    return memeli;  
else  
    return memeliDegil;}}
```

En önemliden en önemsiz doğru ard arda sorular sorarak bir omurgalının memeli olup olmadığına karar veriyoruz.

Bir başka deyişle test örneği en yukarıdan aşağıya doğru yuvarlanır. Örnek bir sınıfa (memeli, memeli değil) geldiğinde durur. Bu sınıf onun tahmin edilen sınıfı olur.

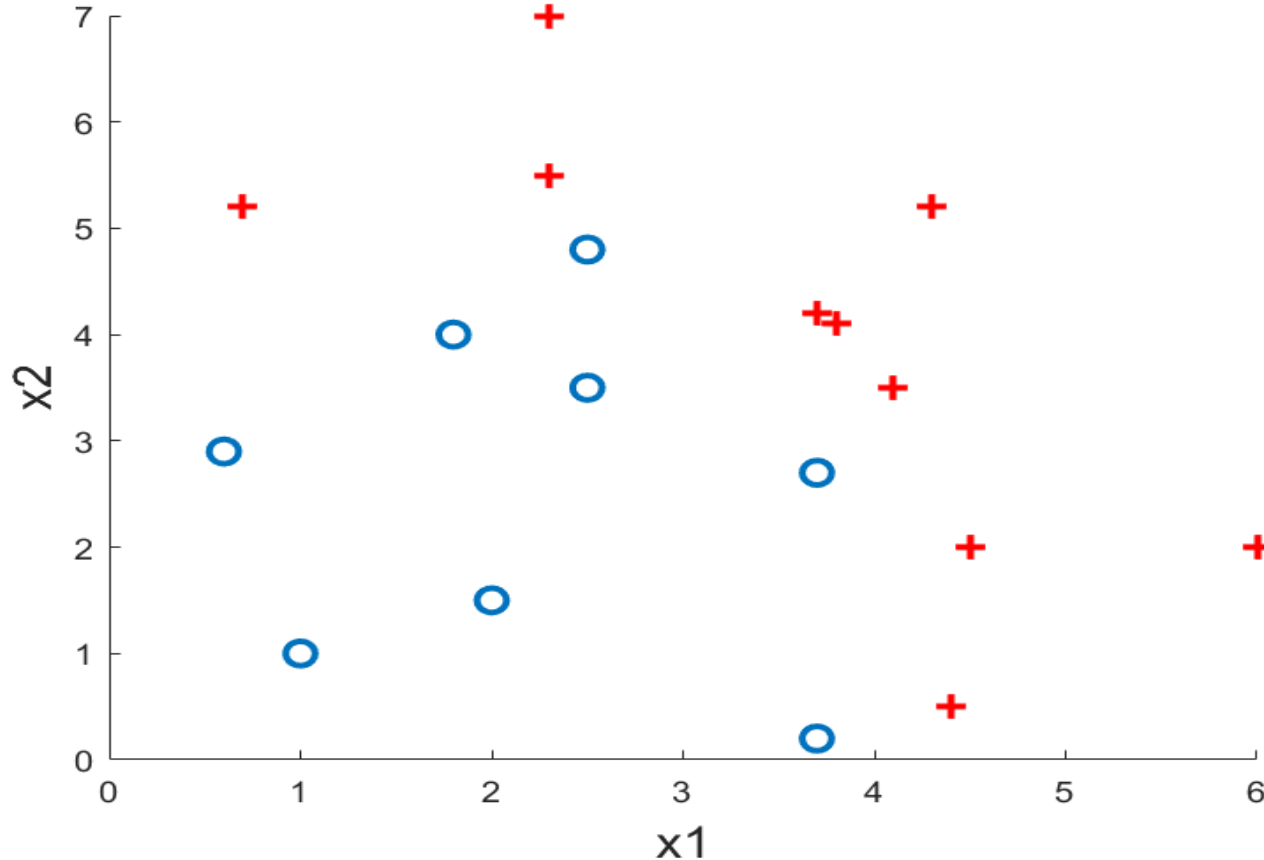


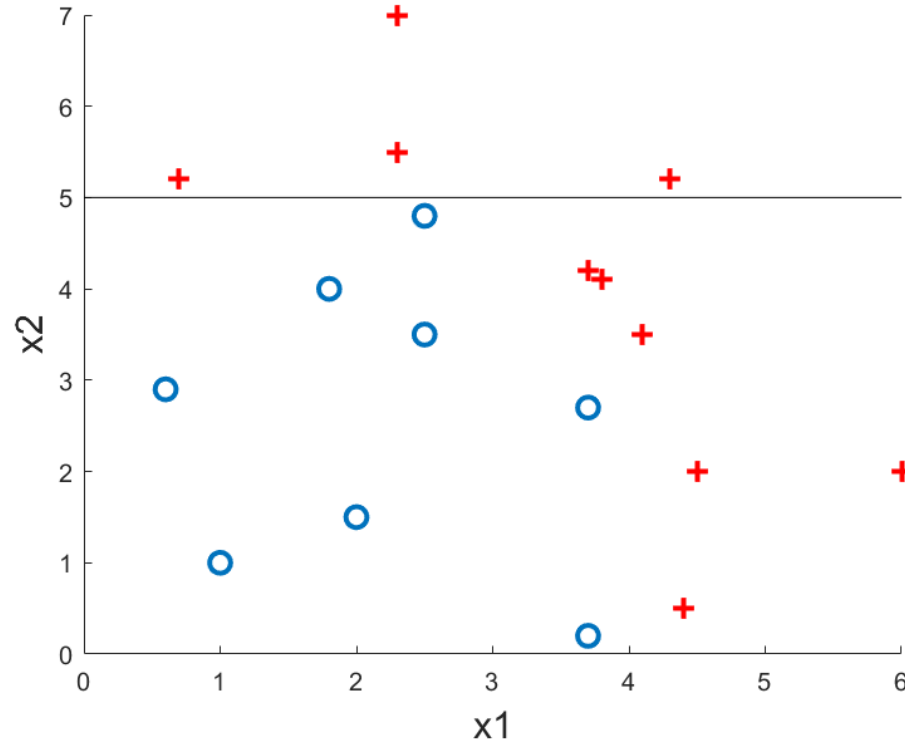


*Kredi verilip verilmeyeceğine karar veren mekanizma.*

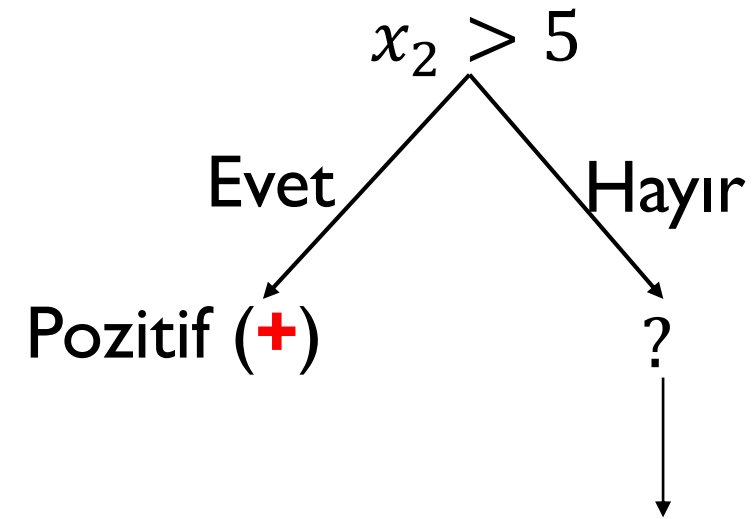
Karar ağaçları giriş uzayında doğrusal olarak ayrılmayan (yani bir hiperdüzlemle ayrılmayan) sınıfları birbirinden ayırmada kullanılan yöntemlerden biridir.

Örnek olarak bir giriş uzayında pozitif (+) ve negatif (○) örnekler aşağıdaki gibi dağılmış olsun. Burada pozitif ve negatif örnekleri bir hiperdüzlemle birbirinden ayırmak mümkün değildir.

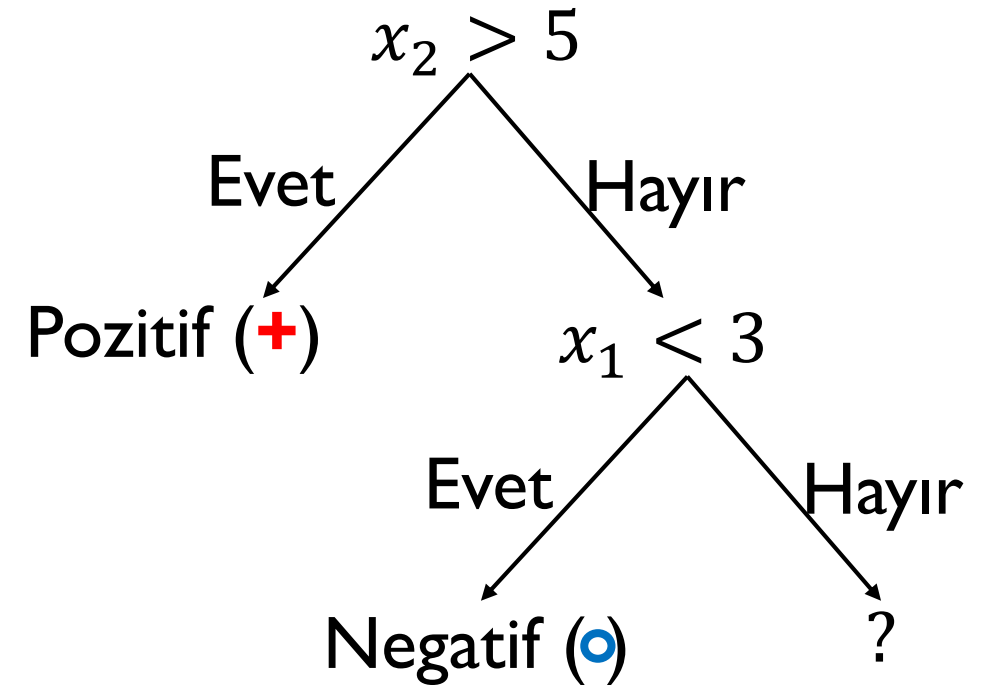
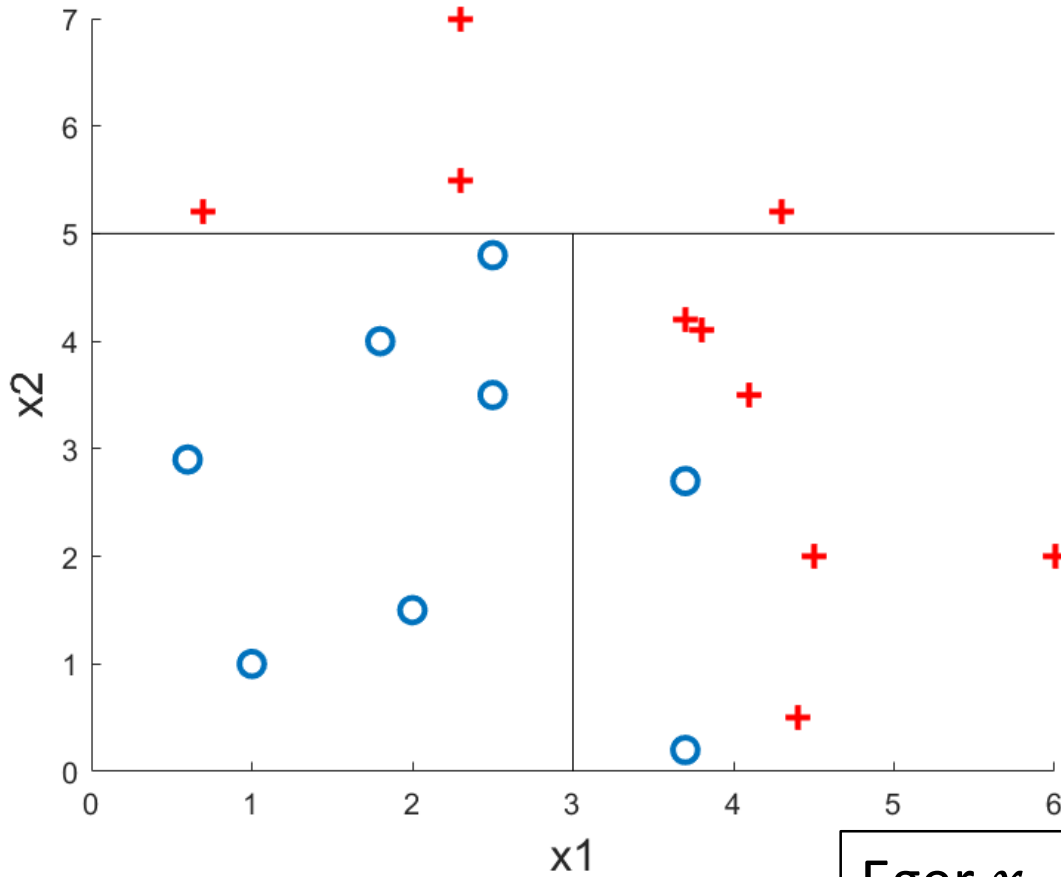




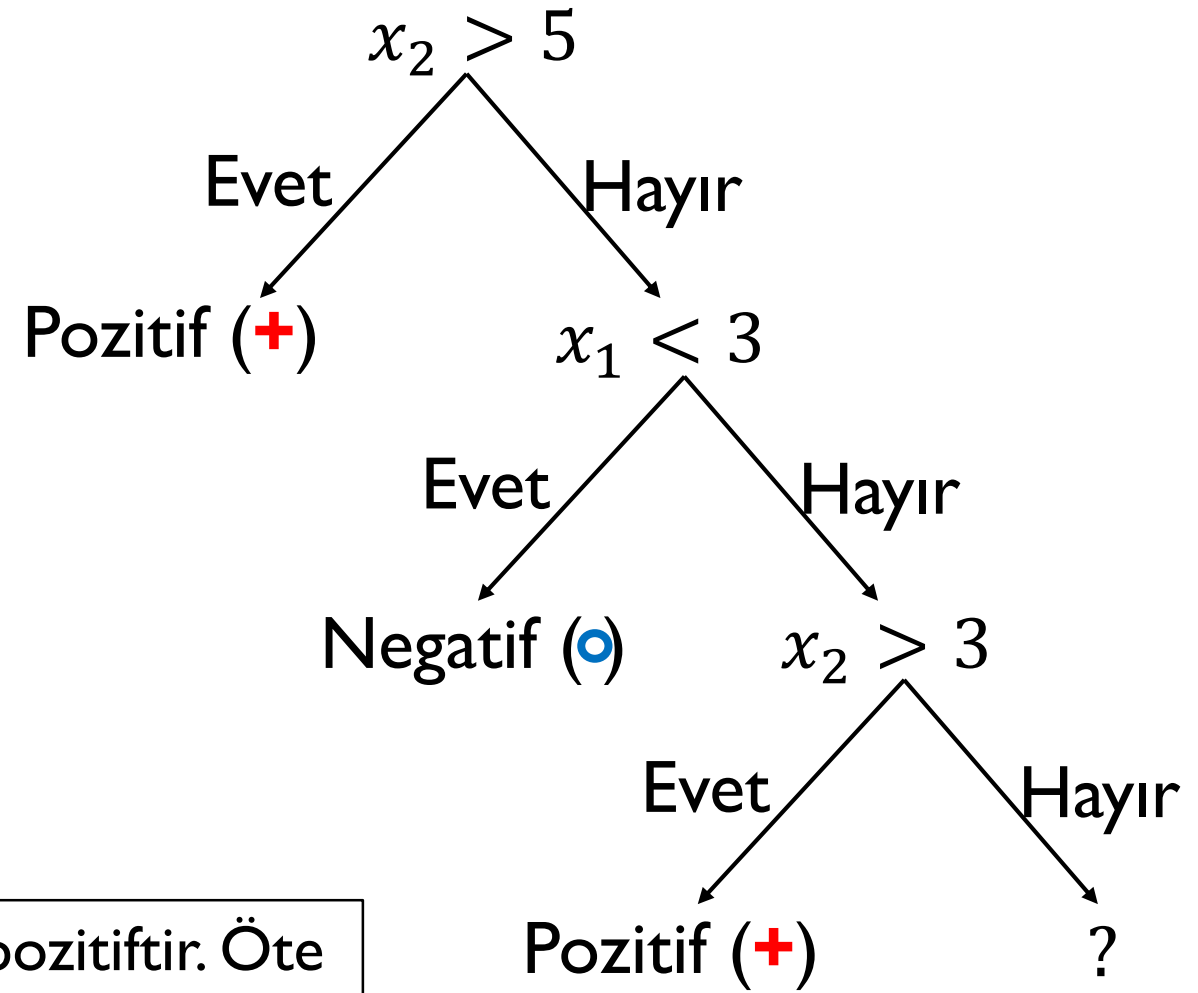
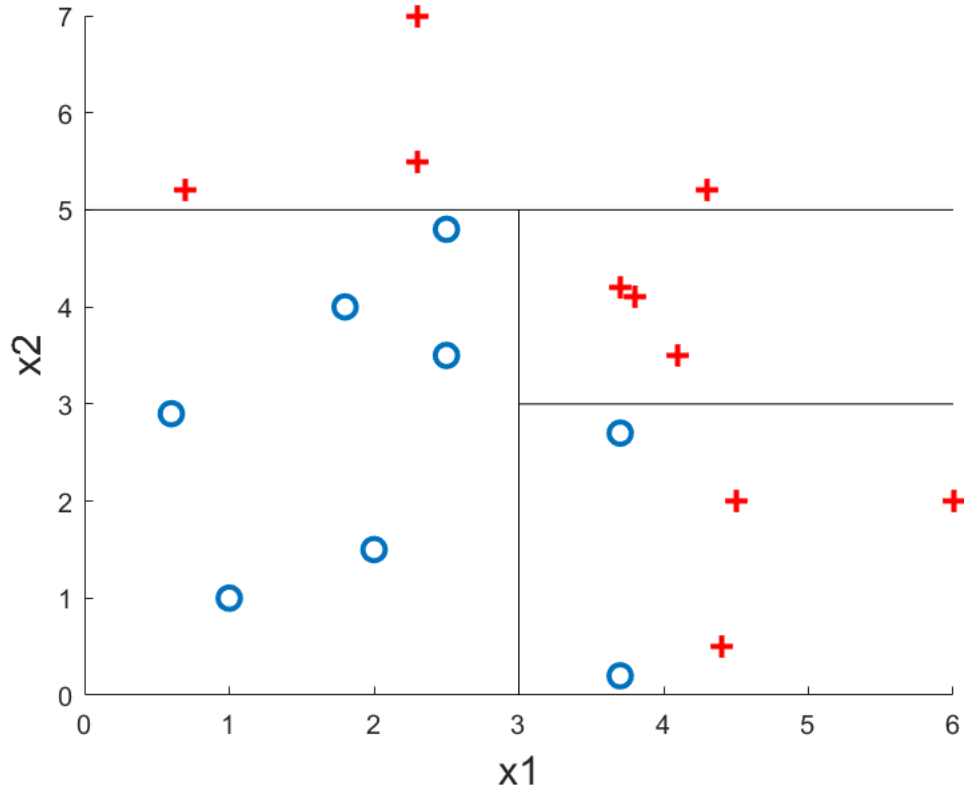
$x_2 = 5$  doğrusu ile uzayı iki parçaya boluyoruz. Öyleki yukarda kalan kısmın tamamı pozitif örneklerden oluşuyor. Bu durumda bir örnek için eğer  $x_2 > 5$  ise direkt bu örnek pozitifdir diyebiliriz.



$x_2 < 5$  olan örneklerde direkt negatif yada pozitif diyemeyiz. Daha ileri araştırma gerekir!

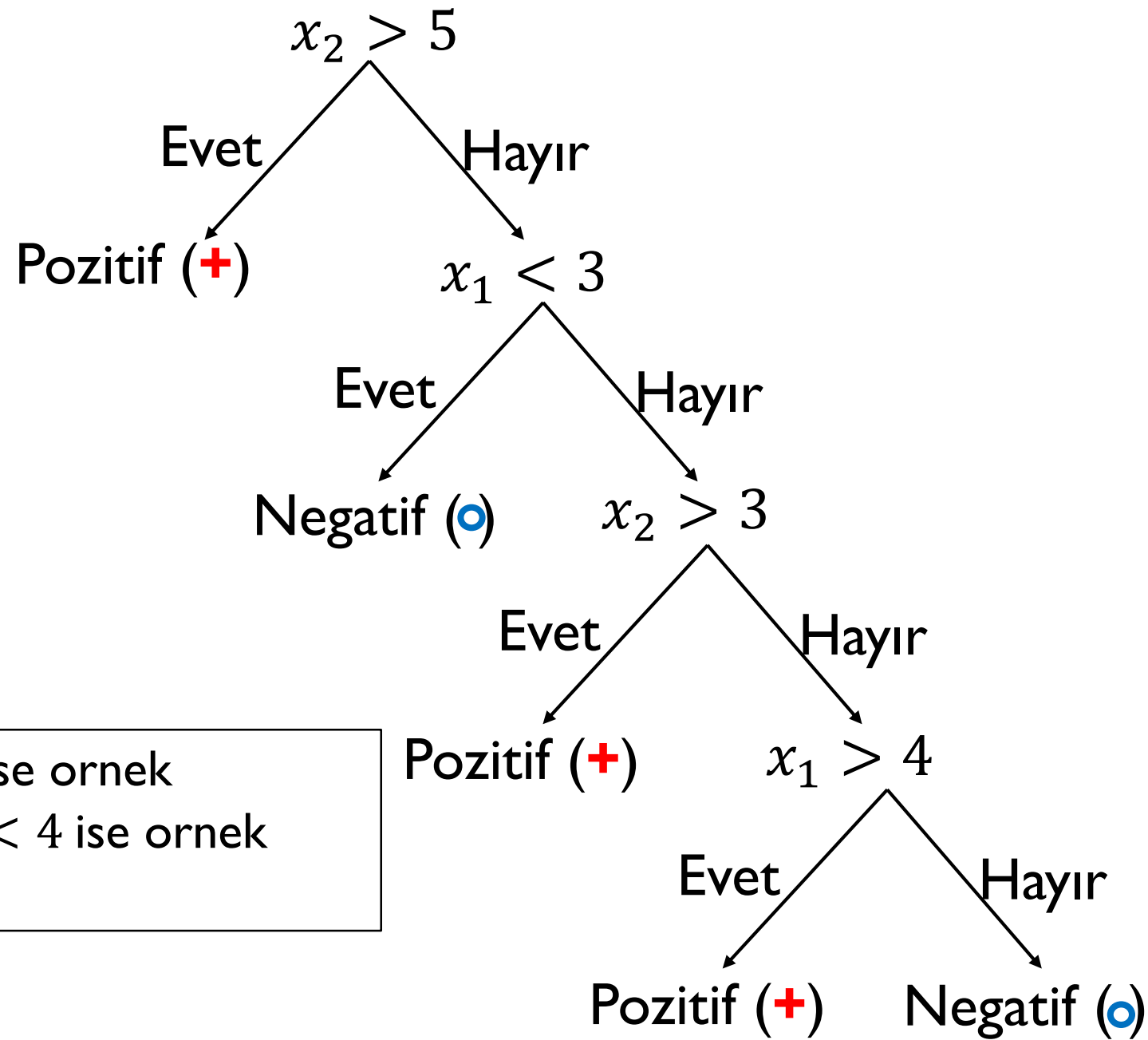
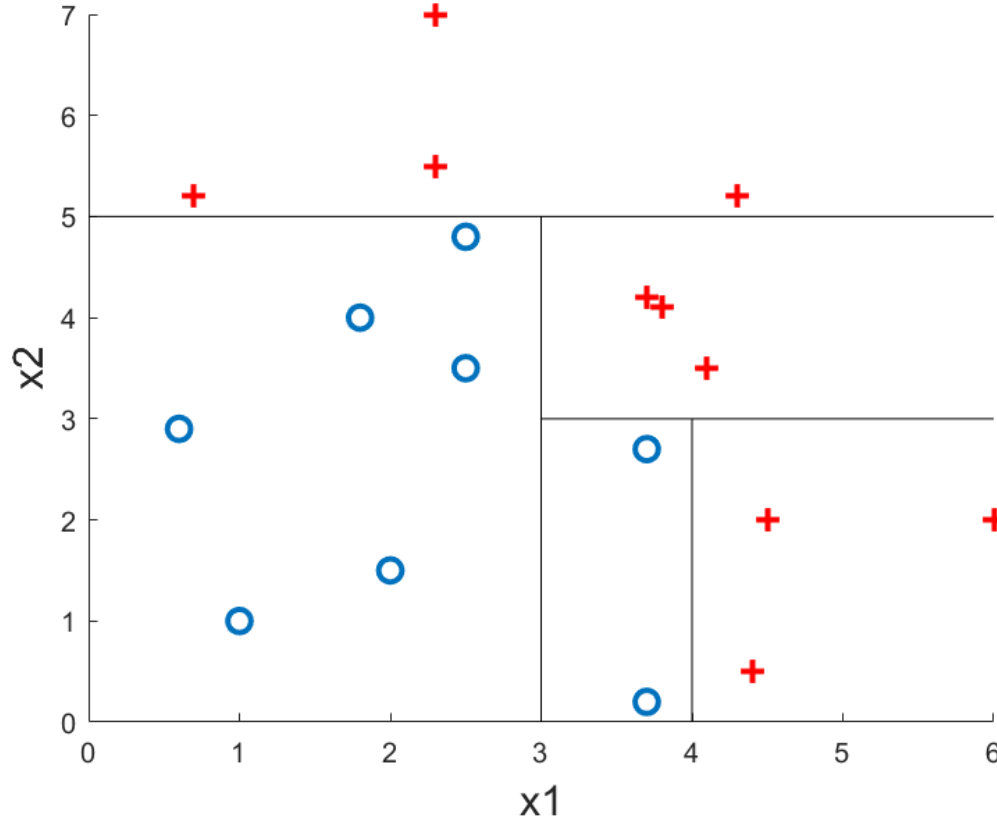


Eğer  $x_2 < 5$  ve  $x_1 < 3$  ise bu örnek negatiftir. olan örneklerde direkt negatif yada pozitif diyemeyiz.  $x_2 < 5$  fakat  $x_1 > 3$  ise örnek negatif yada pozitif olabilir. Bu noktada daha ileri araştırma gerekir!

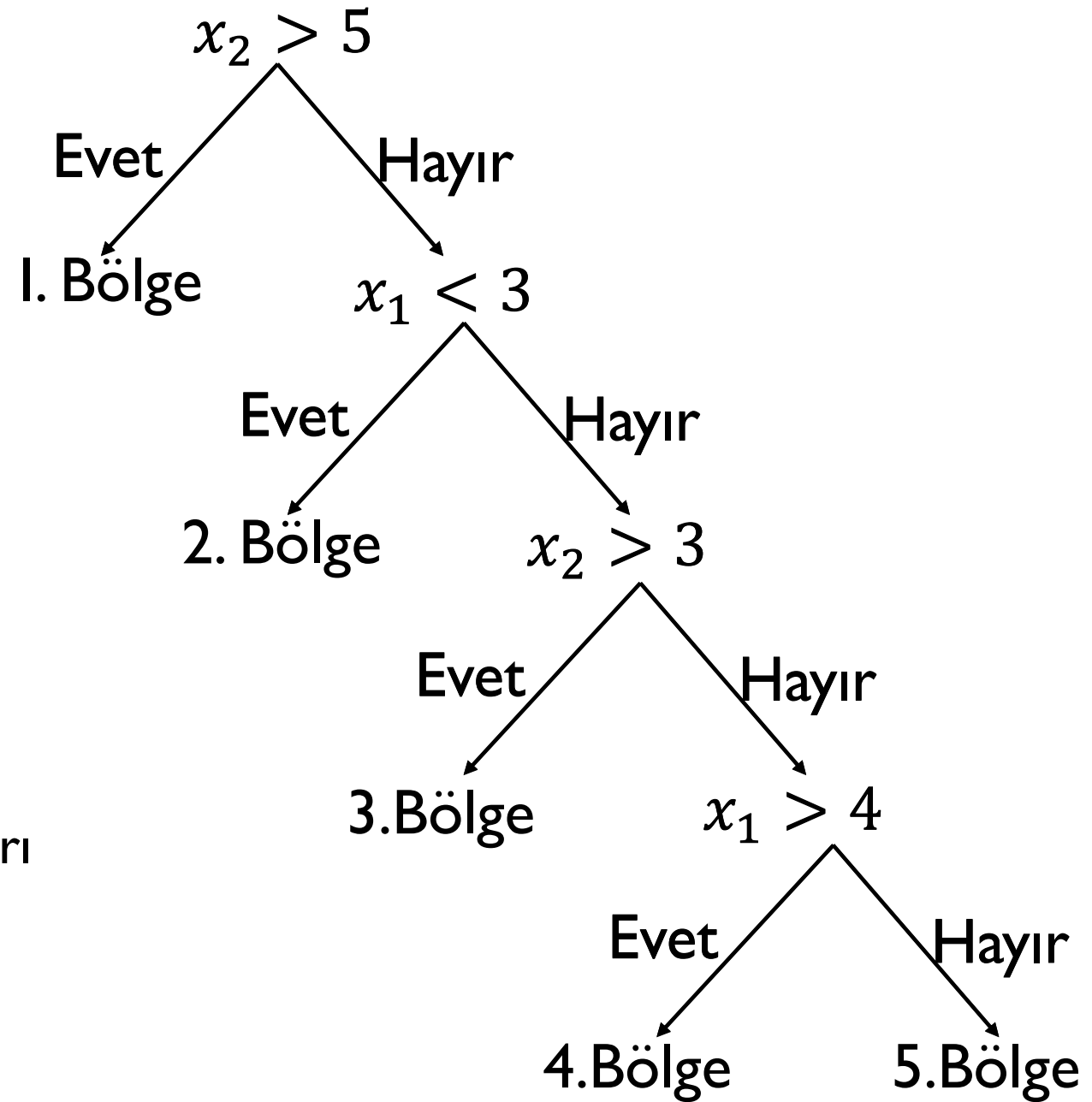
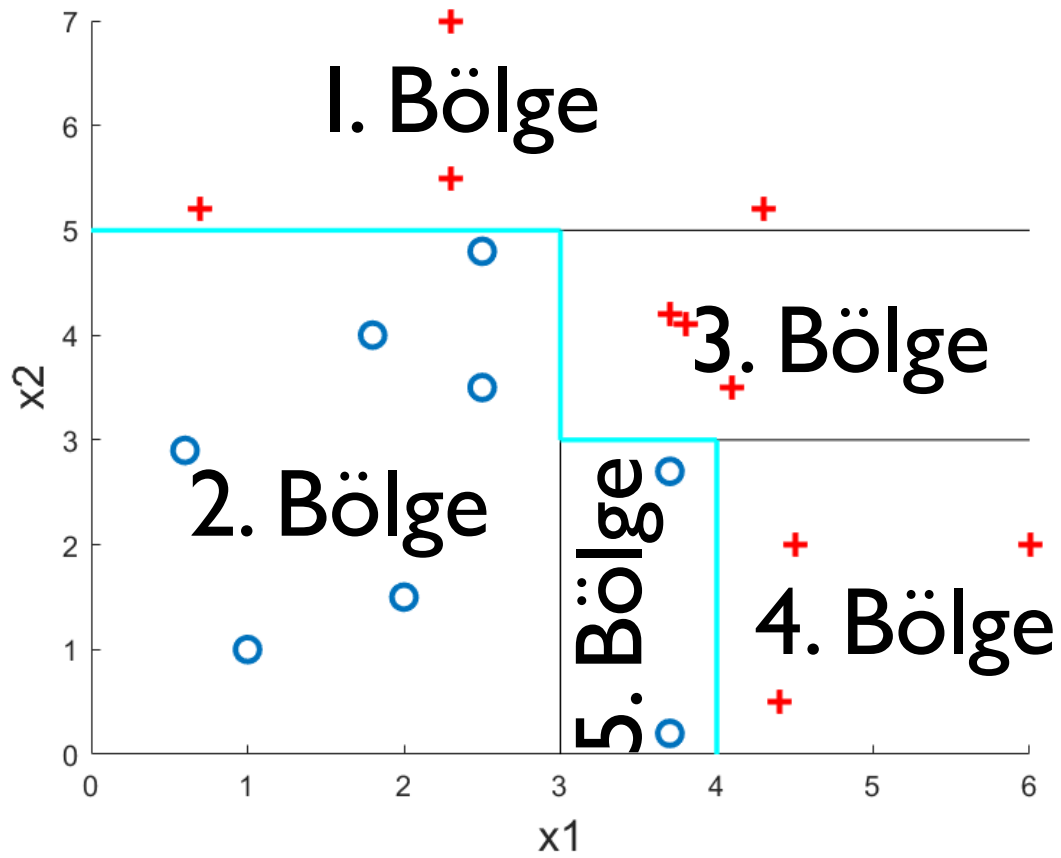


Eğer  $x_2 < 5, x_1 > 3$  ve  $x_2 > 3$  ise bu örnek pozitiftir. Öte yandan örnek  $x_2 < 5, x_1 > 3$  ve  $x_2 < 3$  ise bu örnek pozitif yada negatiftir diyemeyiz. Bu noktada daha ileri araştırma gerekir!



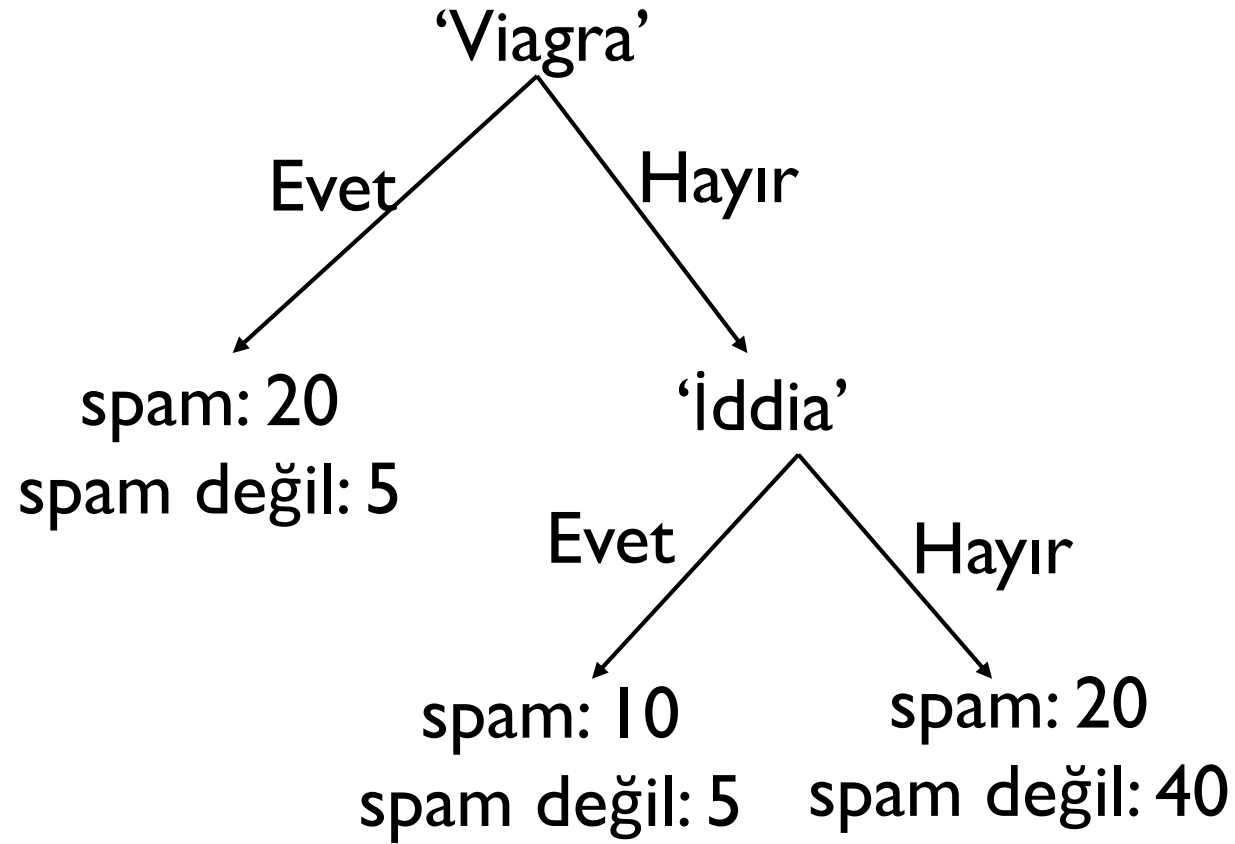


Eğer  $x_2 < 5$  ,  $x_1 > 3$  ,  $x_2 < 3$  ve  $x_1 > 4$  ise örnek pozitifdir;  $x_2 < 5$  ,  $x_1 > 3$  ,  $x_2 < 3$  ve  $x_1 < 4$  ise örnek negatiftir.

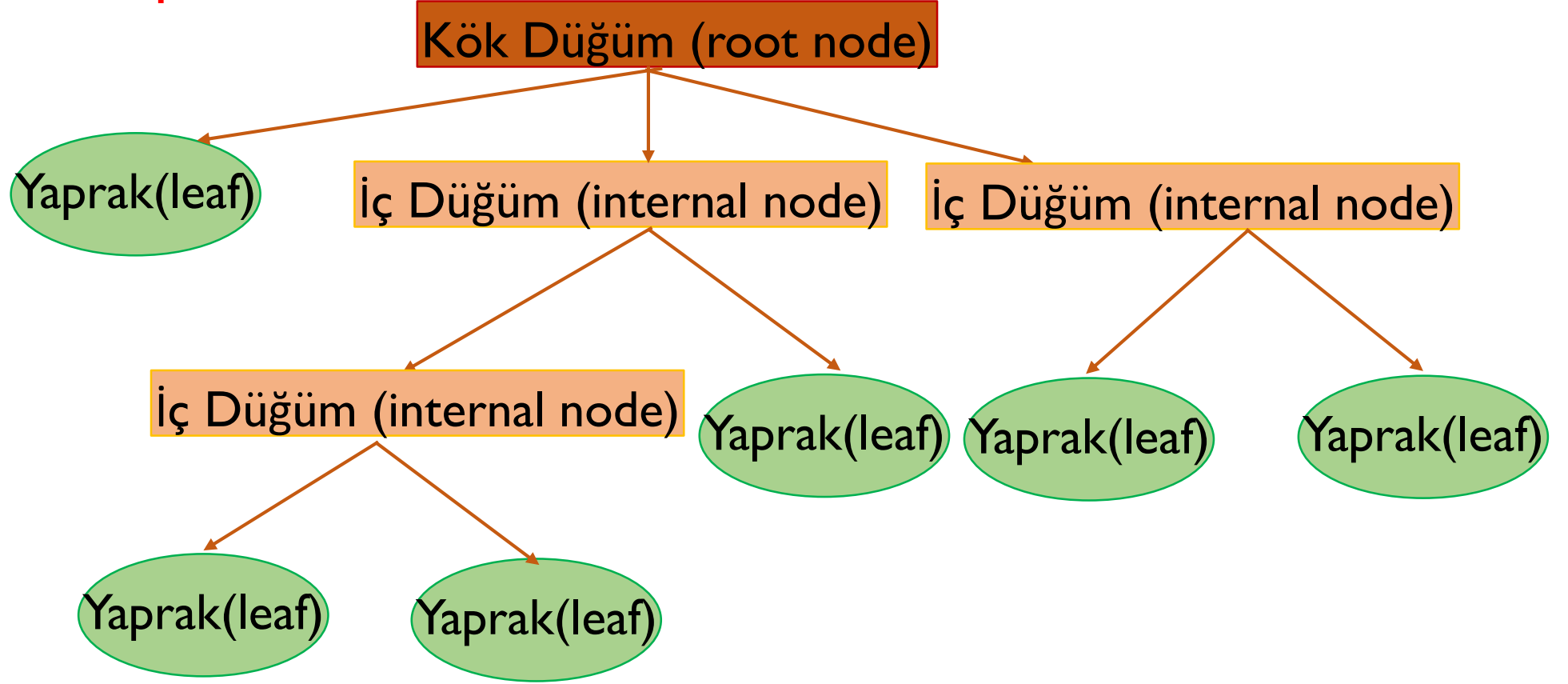


— Karar ağacının oluşturduğu karar sınırı (decision boundary)

**Not:** Dikkat edilirse oluşturulan her bölgede yalnızca bir sınıftan örnekler var (yalnızca pozitif yada yalnızca negatif örnekler). Bu haliye bu bölgelere *saf* (pure) bölgeler diyeceğiz. Fakat karar ağaçlarıyla cogunlukla tamamen saf bölgeler oluşturmak mümkün değildir.



# Karar Ağacının Yapısı



- Bir örneği sınıflandırmaya kök düğümden başlarız.
- İç düğüm, henüz sınıflandırma kararını vermediğimiz noktadır. Bir yaprağa varıncaya kadar ilerleriz. Her iç düğümden bir özellik test edilir. (iç düğüm de aslında bir köktür).
- Bir yaprağa vardığımızda sınıflandırma işlemi durur. Örnek, bu yaprağın ilişkilendirildiği sınıfa eşlenir. Yaprakla ilişkilendirilen sınıf, yaprakta en çok görülen sınıftır.



## Karar Ağacı Nasıl İnşaa Edilir?

Bir karar agaci rekürsif (yinelemeli- kendini çağiran) olarak insa edilir.

Eger herhangi bir anda elimizdeki veri seti yeterince homojen ise, yani aşağı yukarı aynı sınıfa ait örneklerden oluşuyorsa, ağacın bu kısımdan buyumesi durur. Burada yaprak oluşturulur.

Yaprak veri setindeki en yaygın sınıf ile etiketlenir.

Bir noktada veri seti yeterince homojen değilse büyüme devam eder. İlk olarak o noktada bir kök oluşturulur. Ve bu noktadaki veri seti için *en iyi özellik* belirlenir.

Bu özelliğin aldığı farklı değerlere göre veri seti parçalara ayrılır. Her bir parça için tekrar ağaç büyültme algoritması çağrılarak bir ağaç oluşturulur. Bu ağaç köke bağlanır..



```
agacBuyult( $V, O$ ) //  $V$ : veri seti,  $O$ : özellik seti
if homojenMi( $V$ ) :
     $Y$  yaprağını oluştur;
     $Y'$ 'yi  $V'$ 'deki en yaygın sınıfla etiketle;
    return  $Y$ ;
else: //  $V$  yeterince homojen değilse
     $K$  kökünü oluştur;
    enIyiOzellik( $V, O$ ) ile en iyi özelliği bul:  $O_{enIyi}$ 
     $K$  köküne bulunan en iyi özelliğin ( $O_{enIyi}$ ) adını ver;
     $V'$ 'yi  $O_{enIyi}$ 'nin her bir değeri için  $V_i$  parçalarına böl;
    for each  $V_i$ :
        cocuk=agacBuyult( $V_i, O$ ) ;
        cocuk'u  $K$  köküne ekle, bağlantıya  $O_{enIyi}$ 'nin  $i$ . değerini yaz;
    end for;
return  $K$ ;
```

$\text{homojenMi}(V)$  bir boolean fonksiyondur: doğru yada yanlış döner.  $V$  veri seti yeterince homojen ise doğru'ya değilse yanlış'a döner.

$\text{enIyiOzellik}(V, O)$ ,  $V$  veri setini en iyi bölen özelliğe döner. Bu noktayı daha detayli inceleyeceğiz.

```
agacBuyult_v2( $V, O$ ) //  $V$ : veri seti,  $O$ : özellik seti
if  $\text{homojenMi}(V)$  ;;
    return  $V$ 'deki en yaygın sınıf;
else: //  $V$  yeterince homojen değilse
     $\text{enIyiOzellik}(V, O)$  ile en iyi özelliği bul:  $O_{\text{enIyi}}$ 
     $V$ 'yi,  $O_{\text{enIyi}}$ 'nin her bir değeri için  $V_i$  parçalarına böl;
    for each  $V_i$ :
         $\text{cocuk\_i} = \text{agacBuyult}(V_i, O)$ ;
    end for;
return çocukları çocuk_i olan ağac;
```



Kiři	Ev sahibi	Evlilik durumu	Yıllık gelir	Krediyi ödemiş mi?
	Evet	Bekar	125K	Hayır
	Hayır	Evli	100K	Hayır
	Hayır	Bekar	70K	Hayır
	Evet	Evli	120K	Hayır
	Hayır	Bekar	95K	Evet
	Hayır	Evli	60K	Hayır
	Evet	Bekar	220K	Hayır
	Hayır	Bekar	85K	Evet
	Hayır	Evli	75K	Hayır
	Hayır	Bekar	90K	Evet



$V = \{ \text{[Image 1]}, \text{[Image 2]}, \text{[Image 3]}, \text{[Image 4]}, \text{[Image 5]}, \text{[Image 6]}, \text{[Image 7]}, \text{[Image 8]}, \text{[Image 9]}, \text{[Image 10]} \}$

$O = \{\text{ev sahibi, evlilik durumu, yıllık gelir}\}$

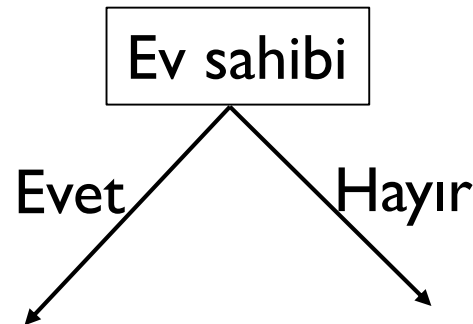
**agacBuyult** ( $V, O$ ) algoritmasına başlayalım...

$\text{homojenMi}(V) = \text{hayır}$ . (krediyiodemeyen 7 kişi, krediyiödeyen 3 kişi)

$\text{enIyiOzellik}(V, O) = O_{\text{enIyi}} = \text{Ev sahibi}$

‘Ev sahibi’ (ilk) kök olur.

Ev sahibi’nin değerleri: Evet, Hayır. O halde Ev sahibi kökünden ‘Evet’ ve ‘Hayır’ dalları çıkar.



Ev sahibi özelliği veri setini ikiye ayırır. Ev sahibi olanlar:  $V_{\text{Evet}}$ , ev sahibi olmayanlar:  $V_{\text{Hayır}}$

$$V_{\text{Evet}} = \{ \text{[Image of a man]}, \text{[Image of a woman]}, \text{[Image of a woman]} \}$$

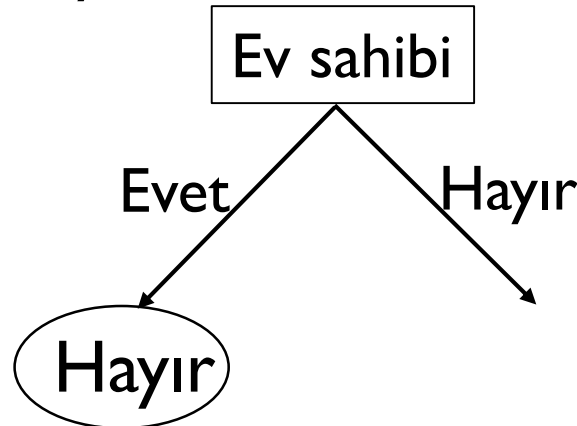
$$V_{\text{Hayır}} = \{ \text{[Image of a woman]}, \text{[Image of a woman]}, \text{[Image of a man]}, \text{[Image of a man]}, \text{[Image of a man]}, \text{[Image of a woman]}, \text{[Image of a man]} \}$$

Bu elde edilen iki veri seti için **agacBuyult** ( $V, O$ ) algoritmasını çağıralım.

**agacBuyult** ( $V_{\text{Evet}}, O$ ) :

**homojenMi** ( $V_{\text{Evet}}$ ) = evet. ( $V_{\text{Evet}}$ 'tekilerin tamamı aldığı krediyi ödemiş).

Şu halde bir yaprak oluşturalım. Oluşan bu yaprağın sınıfı 'Hayır' tir; çünkü bu yapraktakiler aldığı krediyi ödememistir. Ağacın buyumesi bu noktada durur.



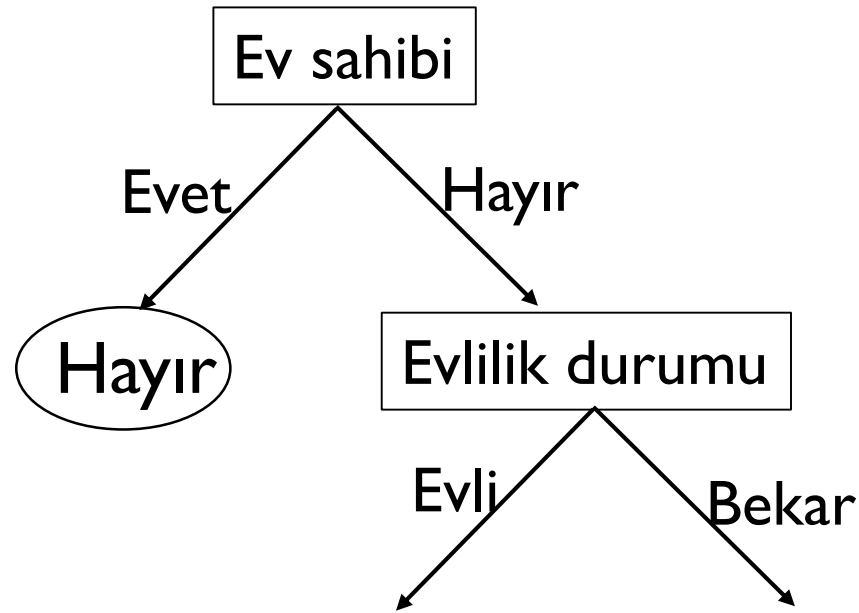
**agacBuyult** ( $V_{\text{Hayır}}, 0$ ) :

homojenMi ( $V_{\text{Hayır}}$ ) = Hayır . ( $V_{\text{Hayır}}$ 'dakilerin bazıları aldığı krediyi ödemiş, bazıları ödememiş).

enIyiOzellik ( $V_{\text{Hayır}}, 0$ ) =  $O_{\text{enIyi}}$  = Evlilik durumu

'Evlilik durumu' kök olur.

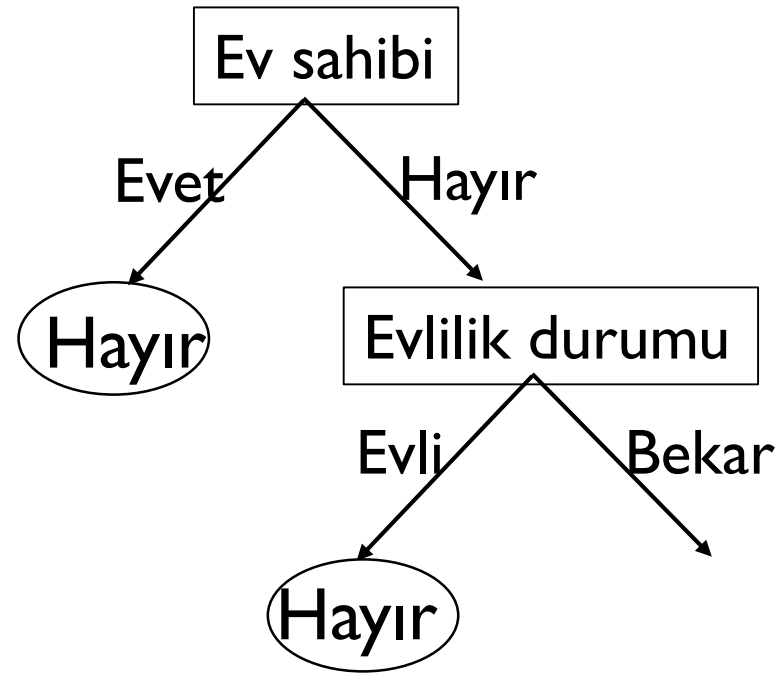
Evlilik durumu'nun değerleri: Evli, Bekar. O halde Evlilik durumu kökünden 'Evli' ve 'Bekar' dalları çıkar.



$V_{\text{Hayır,Evli}} = \{ \text{[Woman]}, \text{[Man]}, \text{[Woman]} \}$

**agacBuyul t** ( $V_{\text{Hayır,Evli}}, 0$ ) :

homojenMi ( $V_{\text{Hayır,Evli}}$ ) = evet . O halde bir yaprak oluşturalım. Oluşan bu yaprağın sınıfı 'Hayır'dır; çünkü ev sahibi olmayıp evli olanların aldığı krediyi ödememistir.



$V_{\text{Hayır,Bekar}} = \{ \text{Hayır}, \text{Bekar} \}$

**agacBuyul t** ( $V_{\text{Hayır,Bekar}}, 0$ ) :

$\text{homojenMi}(V_{\text{Hayır,Bekar}}) = \text{hayır}$ . (Ev sahibi olmayıp bekar olanların bir kısmı (3 kişi) aldığı krediyi odemiş, bazisi (1 kişi) ödememistir).

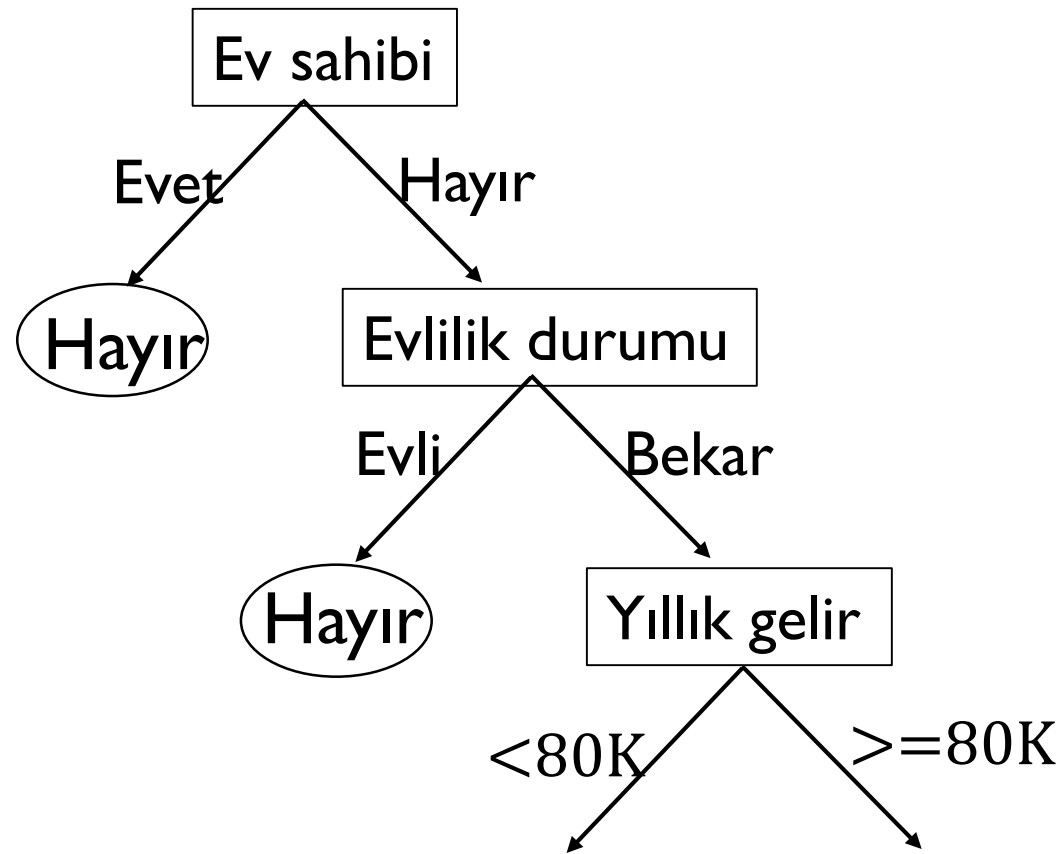
$\text{enIyiOzellik}(V_{\text{Hayır,Bekar}}, 0) = O_{\text{enIyi}} = \text{Yıllık gelir}$ .

‘Yıllık gelir’ kök olur.

Diğer iki özelliğin aksine Yıllık gelir kategorik değil, sayısal bir özelliktir. Sonsuz değer alır. Teorik olarak alacağı her değer için bir dal çıkarmak gerekir. Bu mümkün olmadığından bir eşik değeri buluruz. Bu eşik değerinin nasıl bulunacağını daha sonra anlatacağız.

Bu eşik değeri 80K olsun. Yıllık geliri 80K dan az olanlar için bir dal, 80K dan fazla olanlar için başka bir dal oluşturacağız.





$V_{\text{Hayır,Bekar},80K\_az} = \{ \text{[Image of a young woman]} \}$

$V_{\text{Hayır,Bekar},80K\_fazla} = \{ \text{[Image of Brad Pitt]}, \text{[Image of Bill Gates]}, \text{[Image of George Clooney]} \}$

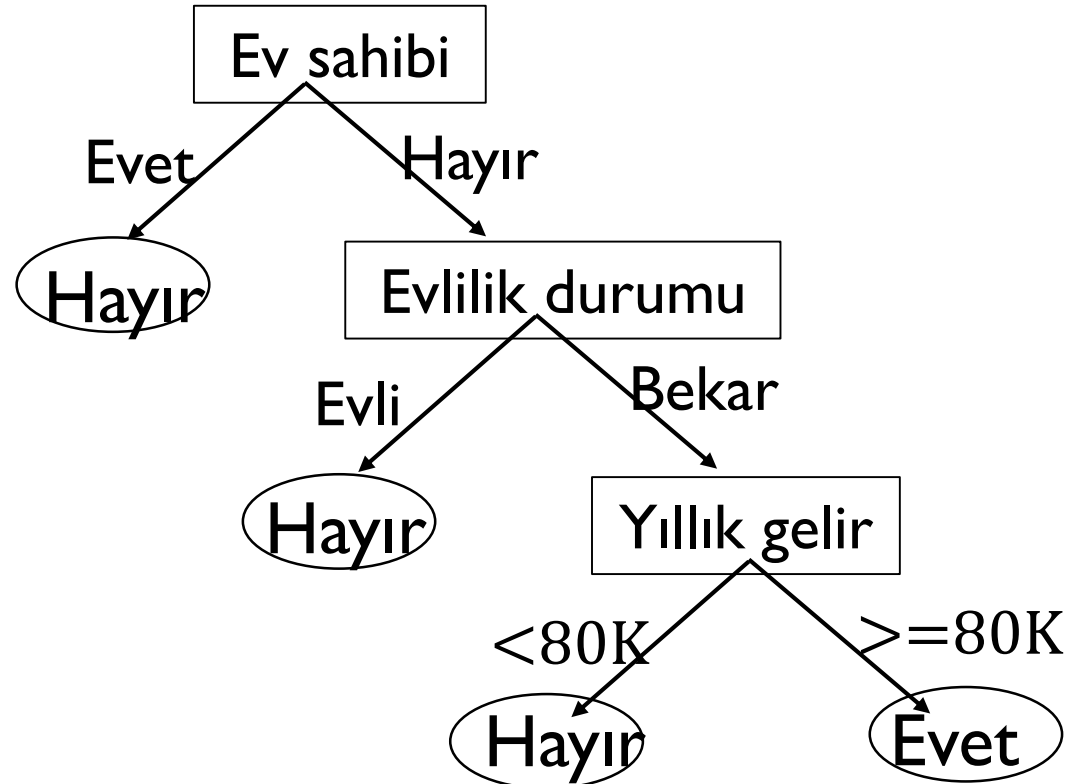


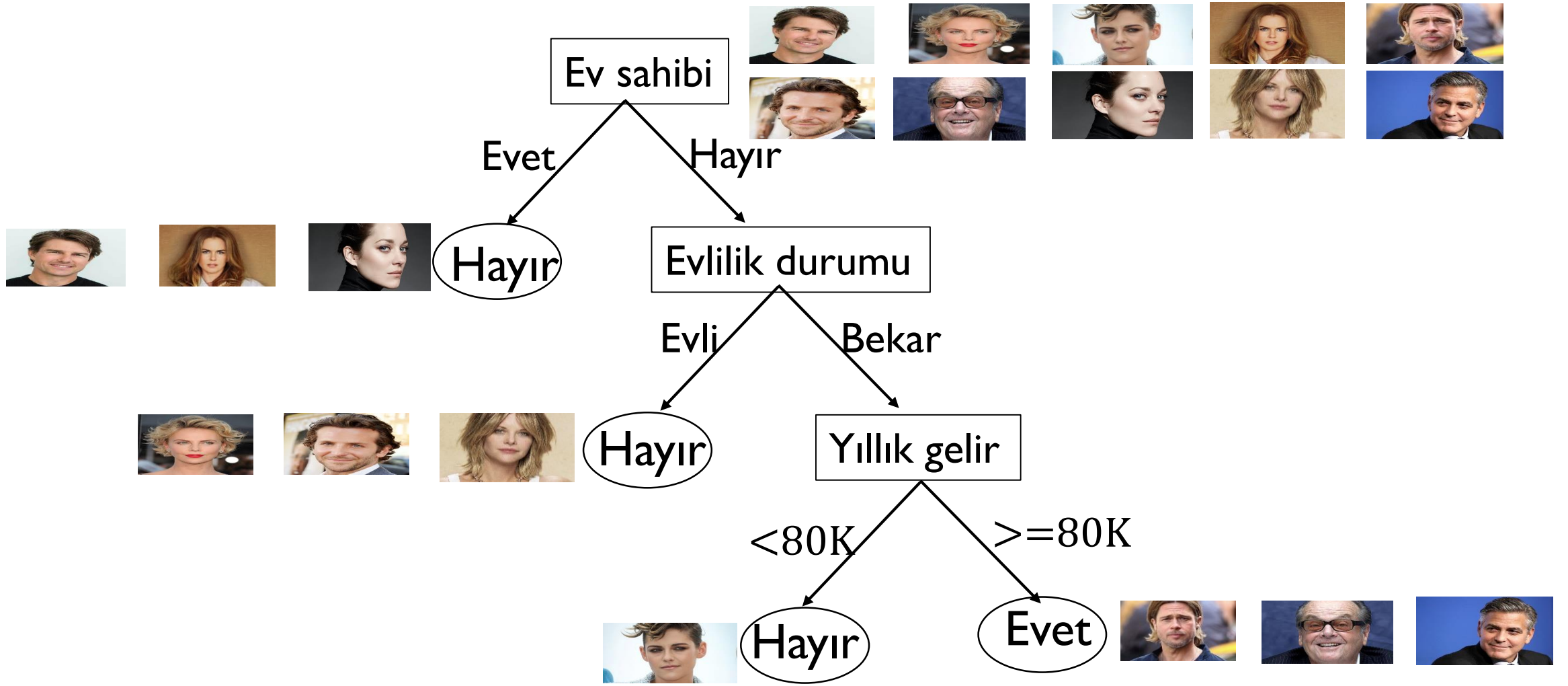
**agacBuyult** ( $V_{\text{Hayır,Bekar,80K\_az}},0$ ) :

$\text{homojenMi}(V_{\text{Hayır,Bekar,80K\_az}}) = \text{evet}$  . Yaprak olusturalim. Oluşan bu yaprağın sınıfı 'Hayır'dır; çünkü ev sahibi olmayıp, bekar ve geliri 80K dan az olanlar aldığı krediyi odememistir.

**agacBuyult** ( $V_{\text{Hayır,Bekar,80K\_fazla}},0$ ) :

$\text{homojenMi}(V_{\text{Hayır,Bekar,80K\_fazla}}) = \text{evet}$  . Yaprak olusturalim. Oluşan bu yaprağın sınıfı 'Evet'tir'dır; çünkü ev sahibi olmayıp, bekar ve geliri 80K dan fazla olanlar aldığı krediyi odemistir.







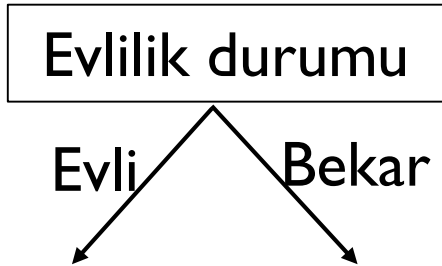
# Özelliklerin Bölünmesi

## Özellikler

### Kategorik

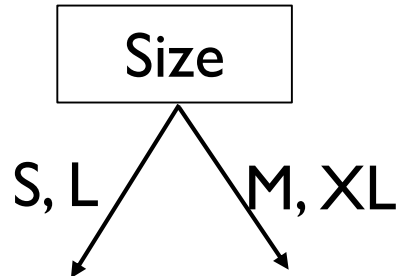
#### Sembolik

Her bir değer için ayrı ayrı bölünür



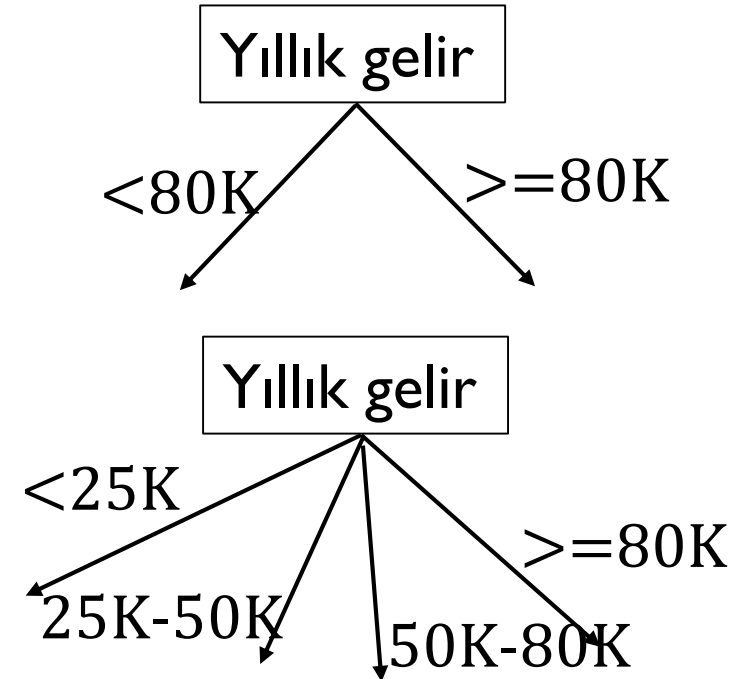
#### Sıralı

Ya her bir değer için ayrı ayrı bölünür; yada bazı değerler gruplandıktan sonra gruplar arasında bölünme olur



### Sayısal

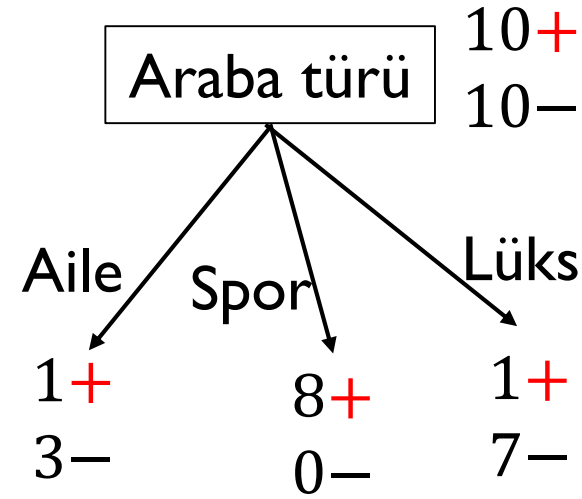
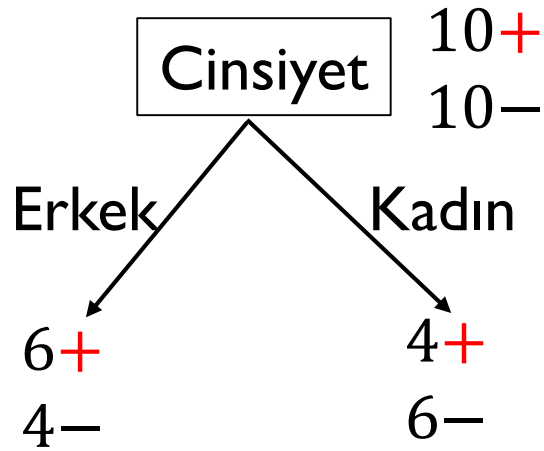
Bir yada birkaç eşik değeri belirlenir. Bu eşiklere göre bölünme yapılır



# Karar Ağacı Oluştururken En İyi Özellik Nasıl Bulunur?

Karar agacinda bir kök için, bu kökü parçalara ayıracak özellik seçerken, her bir parçada sınıf dağılımı en saf/en homojen olacak şekilde özellik secilimi yapılır. Her bir parçada sürpriz/belirsizlik miktarı en düşük olmalıdır.

**ör.** Bir karar agacinda belirli bir kockte elimizde 20 tane örnek olsun. Bu örneklerin 10 tanesi pozitif sınıfa, 10 tanesi negatif sınıfa ait olsun. Bu kök homojen olmadigindan bunu homojen olacak parçalara ayırırız. Diyelimki elimizde iki özellik olsun: cinsiyet ve araba türü.



Araba türü'nün oluşturduğu parçalar daha homojen; hangi sınıfın olacağı belirsizliği daha düşük!

## Süprizin (Belirsizliğin) Ölçülmesi

Information theory'de  $p$  olasiligina sahip bir olay gerçekteştigindeki şaşkınlığımız  $-\log p$  ile ölçülür. Bir başka ifadeyle  $-\log p$ ,  $p$  olasiligina sahip bir olayin gerceklesmesinin bizim icin ne kadar *supriz* oldugunu gosterir.

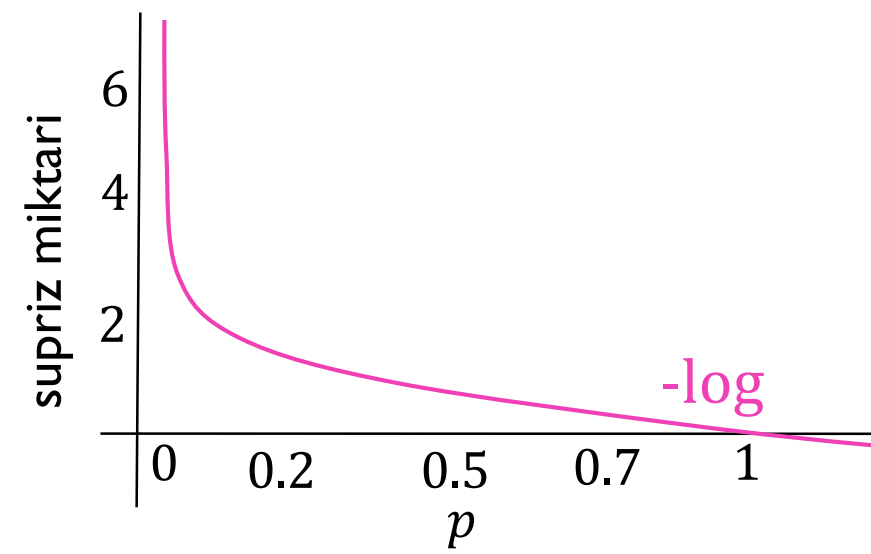
**ör.** Diyelimki elimizde hileli bir bozuk para var. Bu parada tura gelmesi olasılığı  $P(tura) = 0.2$ , yazı gelmesi olasılığı  $P(yazı) = 0.8$  olsun.

Tura gelmesindeki supriz miktarı:  $-\log 0.2 = 1.6$

Yazı gelmesindeki supriz miktarı:  $-\log 0.8 = 0.2$

Atılan bu parada tura gelmesi daha büyük suprizdir (şaşkınlık miktarımız daha fazladır); öte yandan yazı gelmesi pek de süpriz değildir.





$p = 0$  iken (yani olayın gerçekleşme olasılığı yok iken), eğer olay gerçekleşirse supriz sonsuz (infinity) olur.

$p = 1$  iken (yani olayın gerçekleşmesi kesin iken), olay gerçekleşirse supriz 0 olur.

## Entropy (Ortalama Supriz)

$p_i$ 'ler  $K$  tane olayın gerçekleşme olasılığı olsun ( $\sum_{i=1}^K p_i = 1$ ). Bu durumda entropy:

$$\sum_{i=1}^K p_i \cdot -\log p_i$$

her bir suprizi, bu suprizin görülme olasılığı ile carpiyoruz



## Karar ağaçlarına dönersek;

Entropy'yi homojen olmamanın ölçüsü olarak kullanacağız.

(Entropy safsızlığın/belirsizliğin ölçüsüdür).

Yani bir düğümde entropy yüksekse, burada baskın bir sınıf yoktur. Sınıflardan benzer sayılarda örnekler vardır (örneğin 10 pozitif, 10 negatif örnek).

$p_i$ , bir  $V$  kumesinde  $i$ . sınıfın görülme olasılığı olsun ve bu düğümde toplam  $K$  tane sınıf olsun. Bu durumda  $V$  kumesindeki entropy  $H(V) = \sum_{i=1}^K p_i \cdot -\log p_i$

Diyelim ki bir  $O$  özelliğinin  $n$  farklı değeri olsun:  $O_1, O_2, \dots, O_n$

$n$  değer  $V$  kumesini  $n$  parçaya ayırır:  $V_1, V_2, \dots, V_n$

$V$  düğümü bir  $O$  özelliğine göre  $n$  parçaya ayrılsın.

Bu parçaların herbirindeki entropiyi hesaplarız:  $H(V_1), H(V_2), \dots, H(V_n)$ .

Daha sonra bulunan entropileri  $P(O_1), P(O_2), \dots, P(O_n)$  olasılık sırası ile çarpılarak toplanır:



$V$  kumesinin bir  $O$  özelliği ile  $n$  parçaya bolunmesi ile ortaya cikan entropy:

$$\sum_{i=1}^n P(O_i)H(V_i)$$

$V$  kumesinin  $O$  ile bolunmeden onceki entropysi (belirisizligi)  $H(V)$  idi. Yukarida hesaplanan yeni entropi ile arasindaki fark *bölünme ile oluřan kazanç miktarini* verir. Bu da řu řekilde hesaplanir:

$$Kazanç miktarı (O) = H(V) - \sum_{i=1}^n P(O_i)H(V_i)$$

En iyi özellik maksimum kazanç miktarı sađlayan özelliştir; yani belirsizligi en çok dusuren ozelliştir.

