

**İSTANBUL TEKNİK ÜNİVERSİTESİ -FEN EDEBİYAT**  
**FAKÜLTESİ**  
**MATEMATİK MÜHENDİSLİĞİ PROGRAMI**



KARAR AĞACI VE LOJİSTİK REGRESYON ALGORİTMALARININ VERİ SETİ ÜZERİNE UYGULANMASI

**BİTİRME ÖDEVİ**

Elif UYGUN 090120406

Gamze KAPTAN 090120422

**Teslim Tarihi:** 29.05.2017

**Tez Danışmanı:** Doç.Dr.Atabey Kaygun

**MAYIS 2017**

## ÖNSÖZ

Bu çalışmada; Karar Ağaçları (Decision Tree) ve Lojistik Regresyon (Logistic Regression) Algoritmaları kullanılarak örnek bir veri seti üzerinde tahminleme çalışması yapılmış ve doğruluk oranları hesaplanmıştır.

Bitirme tezinde, çalışmaların planlanması, yürütülmesi konusunda bilgi ve tecrübelerini benimle paylaşarak destek olan değerli tez danışmanım Doç. Dr. Atabey Kaygun'a ve eğitim hayatım süresince maddi ve manevi desteklerini esirgemeyen anneme, babama, abime, ve kardeşime teşekkür eder sevgilerimi sunarım.

## İÇİNDEKİLER

ÖNSÖZ.....	HATA! YER İŞARETİ TANIMLANMAMIŞ.
İÇİNDEKİLER.....	2
SEMBOLLER / KISALTMALAR LİSTESİ.....	HATA! YER İŞARETİ TANIMLANMAMIŞ.
ÖZET .....	HATA! YER İŞARETİ TANIMLANMAMIŞ.
ABSTRACT .....	5
1. GİRİŞ.....	HATA! YER İŞARETİ TANIMLANMAMIŞ.
2. GENEL KISIMLAR.....	7
3. MATERYAL VE YÖNTEM (VARSA).....	HATA! YER İŞARETİ TANIMLANMAMIŞ.
4. BULGULAR (VARSA) .....	22
5. TARTIŞMA VE SONUÇ.....	23
KAYNAKLAR.....	24
EKLER .....	25

## ÖZET

**Uygun E. , Kaptan G.** Karar Ağaçları ve Lojistik Regresyon Algoritmalarının Veri Seti Üzerine Uygulanması. İstanbul Teknik Üniversitesi Fen Edebiyat Fakültesi, Matematik Mühendisliği, Bitirme Projesi. İstanbul. 2017.

Bilgisayar teknolojilerindeki gelişmeler, üretilen bilgi miktarlarında ve veri tabanı sistemlerinin hacminde artış meydana getirmiştir. Veri tabanlarında saklı tutulan, yararlı olma potansiyeline sahip verilerin keşfedilerek anlamlı örüntülerin ortaya çıkarılması, veri madenciliği kavramıyla ifade edilmektedir. Karar ağaçları ve Lojistik regresyon, sınıflandırma ve tahmin için sıkça kullanılan veri madenciliği yaklaşımlarından biridir. Bu çalışmada, bir bankadaki müşterilerin çeşitli bilgileri kullanılarak karar ağaçları ve lojistik regresyon algoritmaları ile bu bilgiler anlamlı hale getirilerek, hangi müşterilerin kampanyadaki ürünü alıp almayacağını tahminlemesi yapılmıştır.

Anahtar Kelimeler: Karar Ağaçları, Entropi, Lojistik Regresyon, Odds, R Studio

## **ABSTRACT**

**Uygun E. , Kaptan G.** Application of Decision Trees and Logistic Regression Algorithms on Dataset. İstanbul Technical University, Science&Letter Faculty, Mathematical Engineering, Graduation Project. İstanbul. 2017.

Developments in computer technologies caused increase in amount of information generated and volume of database systems. By discovering data kept stored in databases as ones having beneficial use potential and creation of meaningful patterns is expressed with the concept of data mining. Decision trees and Logistic Regression are one of the data mining approaches widely used for classification and forecasting. In this study, decision trees and logistic regression algorithms were used to make this information meaningful by using various information of customers in a bank to estimate which customers would buy the product in the campaign.

**Key Words:** Decision Trees, Entropy, Logistic Regression, Odds, R Stuido

## **Projenin Tanımı**

Bu projede Karar Ağaçları (Decision Tree) ve Lojistik Regresyon (Logistic Regression) algoritmaları ile veri analizinin nasıl yapıldığı gösterilecektir. Uygulama kısmında; veri madenciliği çözümleri ile istatistik kökenli algoritmaları bir programlama arayüzü altında sunan R programından faydalanılacaktır. R'ın kendisine ait özel kütüphaneleri sayesinde Decision Tree ve Logistic Regression algoritmaları veri kümesinde işlenerek veri kümesine ait analizler ortaya konulacaktır. Veri kümemiz bir bankanın veritabanından alınmış olup bankadaki müşterilere ait çeşitli bilgileri içermektedir. Veri kümesinde bulunan nitelikler (attribute) ve örnekler (instance) her iki algoritmada analize ulaşmak için oldukça önemlidir. Müşterilere ait bilgiler Decision Tree ve Logistic Regression algoritmaları ile işlenerek bu müşterilerden hangilerinin bankaya ait olan ürünü alırken hangilerinin bu ürünü almayacakları tahmin edilerek tahmin oranının doğruluğu hesaplanacaktır.

## 1. GENEL KISIMLAR

### Karar Ağaçları ve Karar Ağaçları Algoritması

Karar Ağaçları sınıflama, özellik ve hedefe göre karar düğümleri (decision nodes) ve yaprak düğümlerinden (leaf nodes) oluşan ağaç yapısı formunda bir model oluşturan bir sınıflandırma yöntemidir. Sınıflandırma ve tahmin için sıkça kullanılan bir veri madenciliği (data mining) yaklaşımıdır. Bir karar alıcı tarafından tercihlerin, risklerin, kazançların ve hedeflerin tanımlanmasında yardımcı olmak için çeşitli karar düğümleri incelenebilir. Bu model, bir karar alıcıya karar alırken hangi faktörlerin göz önüne alınmasının belirlenmesinde yardımcı olur. Yapay Sinir Ağları(Artificial Neural Networks) gibi diğer metodolojilerin de sınıflandırma için kullanılmasına rağmen, karar ağaçları, kolay yorumu ve anlaşılabilirliği açısından karar vericiler açısından avantaj sağlamaktadır. Karar ağaçları;

- ✓ Düşük maliyetli olması,
- ✓ Anlaşılmasının, yorumlanmasının ve veritabanları ile entegrasyonunun kolay olması,
- ✓ Çoğu makine öğrenmesi algoritması ya sayısal ya da sınıflandırma verileri için kullanışlı olmasına rağmen karar ağaçlarının her iki veri tipinin işlenmesi için kullanılması,
- ✓ Güvenirliliklerinin iyi olması,
- ✓ Yüksek miktardaki veriyi kısa sürede işleyebilmesi gibi nedenlerden dolayı en yaygın kullanılan sınıflandırma tekniklerinden biridir.

Karar ağacı algoritması, veri setini küçük ve daha küçük parçalara bölerek ilerler. Karar ağacı tekniğini kullanarak verinin sınıflanması, öğrenme ve sınıflama olmak üzere iki aşamalıdır. Öğrenme aşamasında bir tane öğrenme kümesi oluşturulur. Bu öğrenme kümesindeki örnekleri (instance) en iyi ayıran nitelik (attribute) belirlenir. En iyi nitelik ağacın kök düğümü (root node) olmuş olur. Daha sonra bu kök düğümün alt düğümleri, alt düğümlerin de yine alt düğümleri veya yaprak düğümleri oluşturulur. Bu şekilde veriler küçük ve daha küçük verilere bölünmüş olur. Daha küçük verilere ayırma işlemi örneklerin hepsi aynı sınıfa ait, örnekleri bölecek nitelik kalmamış, kalan niteliklerin özelliklerini taşıyacak örnek kalmayana kadar devam eder ve yaprak düğümleriyle sonlanır. Sınıflama aşamasında ise test verisi oluşturulur ve test verisi karar ağacının doğruluğunu belirlemek

amacıyla kullanılır. Veri setleri birçok nitelikten oluşmaktadır. Bu nitelikler veriyi daha iyi tanımamıza yardımcı olmaktadır. Öğrenme aşamasında örnekleri en iyi ayıran niteliğin belirlenip o niteliğin kök düğüm olduğunu söylemiştik. Veri setindeki en iyi niteliğin belirlenmesi için kullanılan en yaygın yöntem Entropi ölçümüdür.

Entropi kimyada bilinen anlamıyla rastgeleliğin, belirsizliğin ve beklenmeyen durumun ortaya çıkma olasılığını gösterir. Veri madenciliği yaklaşımında da aynı anlama geldiğinden dolayı en iyi niteliğin belirlenmesinde entropi ölçüsü en az olan nitelikler kullanılır. Çünkü entropisi yüksek olan nitelik o oranda belirsiz ve kararsızdır. Verilen bir  $A_K$  alanının Entropi ölçüsünü bulan formüller şu şekildedir:

$$E(C|A_K) = \sum_{j=1}^{M_K} p(a_k, j) * [- \sum_{i=1}^N p(c_i|a_k, j) \log_2 p(c_i|a_k, j)]$$

Bu formülde;

$E(C|A_K)$  =  $A_K$  alanının sınıflama özelliğinin Entropi ölçüsü,  
 $p(a_k, j)$  =  $a_k$  alanının  $j$  değerinde olma olasılığı,  
 $p(c_i|a_k, j)$  =  $a_k$  alanı  $j$  değerindeyken sınıf değerinin  $c_i$  olma olasılığı,  
 $M_K$  =  $a_k$  alanının içerdiği değerlerin sayısı;  $j = 1, 2, \dots, M_K$ ,  
 $N$  = farklı sınıfların sayısı;  $i = 1, 2, \dots, N$ ,  
 $K$  = niteliklerin sayısı;  $k = 1, 2, \dots, K$

Eğer bir  $S$  kümesindeki elemanlar, kategorik olarak  $C_1, C_2, C_3, \dots, C_i$  sınıflarına ayrıştırılırsa,  $S$  kümesindeki bir elemanın sınıfını belirlemek için gereken bilgi şu formülle hesaplanmaktadır.

$$I(s_1, s_2, \dots, s_i) = - \sum_{i=1}^m p_i \log_2(p_i)$$

Bu formülde  $P_i$ ,  $C_i$  sınıfına ayrılma olasılığıdır. Entropi denklemi şu şekilde de ifade edilebilir.

$$E(A) = \sum_{i=1}^n \left( \frac{S_i}{S} \right) I(s_i)$$

Bu durumda  $A$  niteliği kullanılarak yapılacak dallanma işleminde, Bilgi Kazancı (information gain) şu formülle hesaplanmaktadır.



$$Kazanc(A) = E(S) - E(A)$$

Burada E(S); sistemin entropi değeri olurken E(A) ise A niteliğinin entropi değeridir. Her nitelik için ayrı ayrı bilgi kazancı hesaplanır ve bilgi kazancı en yüksek olan nitelik ağacın en üstünde konumlanarak kök düğüm olmuş olur. Bu işlemler her düğüm için örneklerin hepsi aynı sınıfa ait, örnekleri bölecek özellik kalmamış, kalan özelliklerin değerini taşıyan örnek olmayana kadar devam eder. Ve bu adımlar tamamlanarak karar ağacı tamamlanmış olur. Şimdi az örnekli bir veri kümesi üzerinde kök düğüm'ün nasıl bulunduğu bakalım:

### Örnek Uygulama

Bu örnek veri kümemizde hava, ebeveyn ve para durumuna göre o hafta ne yapılacağına karar verilecektir. Örnek veri kümesine bakıldığında;

<b>Hafta</b>	<b>Hava Durumu</b>	<b>Ebeveyn Durumu</b>	<b>Para Durumu</b>	<b>Karar</b>
1.Hafta	Güneşli	Var	Var	Sinema
2.Hafta	Güneşli	Yok	Var	Tenis
3.Hafta	Rüzgarlı	Var	Var	Sinema
4.Hafta	Yağmurlu	Var	Yok	Sinema
5.Hafta	Yağmurlu	Yok	Var	Ev
6.Hafta	Yağmurlu	Var	Yok	Sinema
7.Hafta	Rüzgarlı	Yok	Yok	Sinema
8.Hafta	Rüzgarlı	Yok	Var	Alışveriş
9.Hafta	Rüzgarlı	Var	Var	Sinema
10.Hafta	Güneşli	Yok	Var	Tenis

Veri setimizde toplam 10 örnek vardır. Bu 10 örnekten;

6 örnek için karar sinema (6/10),

2 örnek için karar tenis oynamak (2/10),

1 örnek için karar evde kalmak (1/10),

1 örnek için karar alışverişe gitmektir (1/10).

S kümesindeki elemanlar 4 sınıfa ayrılmıştır. Her biri için sınıflara ayrılma olasılığı parantez içinde belirtilmiştir. İlk önce sistemin entropisi bulunacaktır.

$$Entropi = -(6/10)\log_2(6/10) - 2/10\log_2(2/10) - 1/10\log_2(1/10) - 1/10\log_2(1/10); E(S) = 1,571\text{'dir.}$$

Kök düğümünün hangisi olduğunu bulmak için niteliklerimiz olan; hava, ebeveyn ve para özelliklerimizin bilgi kazançlarını hesaplanıp en yüksek kazanca sahip olan nitelik kök düğüm olarak seçilecektir. Hava niteliğimiz için bilgi kazancı değeri;

Kazanç (Hava Durumu)=  $E(S) - E(\text{Hava Durumu})$ ' dur. Bunun için hava durumunun entropi değerine ihtiyaç vardır.

Hava Durumu Güneşli, Rüzgarlı ve Yağmurlu şeklindedir. S kümesi Hava durumu için 3 sınıfa ayrılmıştır.  $S_1$  : Güneşli,  $S_2$ : Rüzgarlı,  $S_3$  : Yağmurlu

- ✓  $S_1$  toplam 3 tane olup bunlardan 1'i sinemaya gitmek 2'si de tenis oynamaktır.  $S_1$  2 sınıfa ayrılmış olup bu sınıfların  $S_1$  de bulunma olasılıklarına göre bilgi değerleri

$$I(S_1) = -1/3\log_2(1/3) - 2/3\log_2(2/3) = 0,918$$

- ✓  $S_2$  toplam 4 tane olup bunlardan 3'ü sinemaya gitmek 1'i de alışveriş yapmaktır.  $S_2$  2 sınıfa ayrılmış olup bu sınıfların  $S_2$ 'de bulunma olasılıklarına göre bilgi değerleri

$$I(S_2) = -3/4\log_2(3/4) - 1/4\log_2(1/4) = 0,811$$

- ✓  $S_3$  toplam 3 tane olup bunlardan 2'si sinemaya gitmek 1'i de evde kalmaktır.  $S_3$  2 sınıfa ayrılmış olup bu sınıfların  $S_3$ 'te bulunma olasılıklarına göre bilgi değerleri

$$I(S_3) = -2/3\log_2(2/3) - 1/3\log_2(1/3) = 0,918$$

Herhangi bir A niteliği için Entropi değeri

$$E(A) = \sum_{i=1}^n \left( \frac{S_i}{S} \right) I(s_i)$$

formülünden;

$$E(\text{HavaDurumu}) = 3/10*0,918 + 4/10*0,811 + 3/10*0,918 = 0,8752$$

Bu durumda hava durumunun Bilgi Kazancı =  $1,571 - 0,8752 = 0,70$  olarak bulunur.

Ebeveyn niteliği için bilgi kazancına bakılırsa; Ebeveyn durumu var veya yok şeklindedir. S kümesi için Ebeveyn niteliği 2 sınıfa ayrılmıştır.  $S_1$  : Var,  $S_2$  : Yok

- ✓  $S_1$  toplamda 5 tane olup 5'i de sinemaya gitmektir.  $S_1$  1 sınıfa ayrılmış olup bu sınıfın  $S_1$ 'te bulunma olasılıklarına göre bilgi değerleri

$$I(S_1) = 5/5 \log_2(5/5) = 0$$

- ✓  $S_2$  toplamda 5 tane olup 2'si tenis, 1'i ev, 1'i sinema, 1'i de alışveriştir.  $S_2$  4 sınıfa ayrılmış olup bu sınıfların  $S_2$ 'te bulunma olasılıklarına göre bilgi değerleri

$$I(S_2) = -2/5 \log_2(2/5) - 1/5 \log_2(1/5) - 1/5 \log_2(1/5) - 1/5 \log_2(1/5) = 1,922$$

Herhangi bir A niteliği için Entropi değeri

$$E(A) = \sum_{i=1}^n \left( \frac{S_i}{S} \right) I(s_i)$$

formülünden;

$$E(\text{ebevyn}) = (5/10) * 0 + (5/10) * 1,922 = 0,961$$

Bu durumda ebevyn durumunun Bilgi Kazancı =  $1,571 - 0,961 = 0,61$  olarak bulunur.

Para durumu için de aynı işlemler yapıldığında Bilgi Kazancı **0,2816** şeklinde bulunur. Hava durumu, ebevyn ve para niteliklerinden bilgi kazancı en yüksek olan 0,70 ile hava durumudur. Bu nedenle hava durumu kök düğüm olmuş olur. Böylece az örnekli bir veri setinde kök düğümün bu şekilde hesaplama ile bulunduğunu göstermiş olduk. Bizim üzerimizde çalışacağımız veri kümesi 4521 örnekli bir veri olduğundan dolayı bu şekilde hesaplama yaparak kök düğüme ulaşmamız çok uzun zaman alacaktır. Bu nedenle R studio üzerinde R'ın kendine ait kütüphaneleri kullanılarak sonuç alınacaktır.

## Veri Seti:

Veri setimiz Portekiz'e ait bir bankanın müşteri bilgilerini içermektedir. Bankada olan kampanyalar için bu müşteriler telefon ile aranmaktadır. Eğer müşteri kampanyaya ait ürünü almayı kabul ederse birkaç defa daha iletişime geçilmektedir. Bank.csv tablomuzda 4521 tane örnek (instance) olup, 16 + karar (y) niteliği (attribute) bulunmaktadır. Bu nitelikler örneklere ait özelliklerdir ve veriyi tanımamıza yardımcı olur. Bu niteliklere bakılacak olursa;

1. age: Müşterilere ait yaş bilgisini içerir. Sayısal bir değerdir.
2. job: Müşterilerin meslek bilgisini içerir. Kategorik olup admin, unknown, unemployed, management, housemaid, entrepreneur, student, blue-collar, self-employed, retired, technician, services şeklinde 12 çeşittir.
3. marital: Müşterilerin medeni durumlarını içermektedir. Kategorik olup married, divorced, single şeklinde 3 çeşittir.
4. education: Müşterilere ait eğitim bilgilerini içermektedir. Kategorik olup unknown, secondary, primary, tertiary şeklinde 4 çeşittir.
5. default: Müşterinin bankada kredisi olup olmadığının bilgisini içerir. Kategorik olup Yes, No şeklinde 2 çeşittir.
6. balance: Müşterinin ortalama yıllık bakiyesini euro olarak ifade etmektedir. Sayısal bir değerdir.
7. housing: Müşterinin konut kredisi olup olmadığının bilgisini içerir. Kategorik olup Yes, No şeklinde 2 çeşittir.
8. loan: Müşterinin bireysel kredisinin olup olmadığının bilgisini içerir. Kategorik olup Yes, No şeklinde 2 çeşittir.
9. contact: Müşteriye ait iletişim tipini belirtir. Kategorik olup unknown, telephone, cellular şeklinde 3 çeşittir.
10. day: Bu kampanya için müşteriyle en son iletişime geçirilen ayın gün bilgisini içerir. Sayısal bir değerdir.
11. month: Bu kampanya için müşteriyle en son iletişime geçirilen ayın bilgisini içerir. Kategorik bir değer olup jan,feb,mar,apr,may,june,july,aug,sep,oct,nov,dec şeklinde 12 çeşittir.

12. duration: Bu kampanya için müşteriyle en son iletişime geçirilen görüşmenin saniye olarak süresini içerir. Sayısal bir değerdir.
13. campaign: Bu kampanyada bir müşteri için kaç kez iletişime geçtiğinin bilgisini içerir. Sayısal bir değerdir.
14. pdays: Bir önceki kampanyada müşteri ile en son iletişime geçilmesinden sonra kaç gün geçtiğini gösterir. Sayısal bir değerdir. Burada bulunan -1 değeri daha önce iletişime geçilmediğini göstermektedir.
15. previous: Bir önceki kampanyada müşteri ile iletişime kaç defa geçirildiğinin bilgisini içerir. Sayısal bir değerdir.
16. poutcome: Müşterinin bir önceki kampanyadaki ürünü alıp almadığının sonucunu içerir. Kategorik olup unknown, other, failure, success şeklinde 4 çeşittir.
17. y: Müşteri ile ilgili kredi verilip verilmemesinin kararının verildiği alandır. Kategorik olup Yes, No şeklinde 2 çeşittir.

Uygulama kısmında veri madenciliği çözümleri ile hem istatistik kökenli uygulamaları hem de yapay zeka kökenli algoritmaları bir programlama arayüzü altında sunan R Studio'dan yararlanılmıştır. Banka verisine ait veri kümesi txt olarak düzenlenmiştir.

```
data<-read.table("Banka.txt",header=T,sep="\t")
head(data)
```

age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
30	unemployed	married	primary	no	1787	no	no	cellular	19	oct	79	1	-1	0	unknown	no
33	services	married	secondary	no	4789	yes	yes	cellular	11	may	220	1	339	4	failure	no
35	management	single	tertiary	no	1350	yes	no	cellular	16	apr	185	1	330	1	failure	no
30	management	married	tertiary	no	1476	yes	yes	unknown	3	jun	199	4	-1	0	unknown	no
59	blue-collar	married	secondary	no	0	yes	no	unknown	5	may	226	1	-1	0	unknown	no
35	management	single	tertiary	no	747	no	no	cellular	23	feb	141	2	176	3	failure	no

## read.table Kullanımı

read.table fonksiyonu dosyayı tablo formatında okur ve dosyadaki alanlara ve satırlara karşılık gelen durumlarla ilgili olarak bir veri çerçevesi oluşturur.

read.table(**file**, header=FALSE, **sep**="", ) şeklindedir. Burada **file parametresi**; okunacak verinin dosya ismidir. Buraya dosyanın tam yolu yazılabileceği gibi sadece adı da yazılabilir. Eğer ki dosya R'nin çalıştığı dizinde bulunuyorsa dosyanın sadece adını yazmak yeterlidir. Getwd() komutu ile R'nin çalıştığı dizinin yolu öğrenilir. Banka dosyamızı bu dizinin altına koyduğumuz için sadece dosya adını ve uzantısını yazmamız yeterli oldu. **header parametresi**; dosyadaki değişkenlerin adlarını ilk satırda içerip içermediğinin bilgisini içeren mantıksal bir değerdir. Default olarak FALSE değerindedir. Banka dosyamızda ilk satırda değişkenlerin ismi yazılmıştır. Bu nedenle TRUE değerini alır. **sep parametresi**; ayırıcı karakterdir. Satırlardaki her bir değer birbirinden nasıl ayrıldığını belirtmek için kullanılır. Default değeri "" şeklindedir. Bizim verimizde değerler birbirinden tab ile ayrıldıklarından dolayı sep="\t" 'dir. Veri setimizi data isimli bir değişkene atadık. Veri setimiz 4521 örnekli olduğundan dolayı bütün satırları yazmak yerine head fonksiyonu ile ilk 6 satır yazılmış oldu.

Örneğin; ilk örneğimiz 30 yaşında, işsiz, evli, ilkokul mezunu, aylık geliri 1787 euro olan, en son iletişime 19 ekimde 79 sn'lik bir telefon görüşmesiyle geçilmiş, önceki kampanyadaki ürünü alıp almadığı bilinmeyen bir müşteridir.

```
: summary(data)
```

```

      age                job                marital                education                default
Min.   :19.00    management :969    divorced: 528    primary   : 678    no :4445
1st Qu.:33.00    blue-collar:946    married  :2797    secondary:2306    yes: 76
Median :39.00    technician :768    single   :1196    tertiary  :1350
Mean   :41.17    admin.     :478
3rd Qu.:49.00    services   :417
Max.   :87.00    retired    :230
              (Other) :713

      balance    housing    loan                contact                day
Min.   : -3313    no :1962    no :3830    cellular :2896    Min.   : 1.00
1st Qu.: 69      yes:2559    yes: 691    telephone:301    1st Qu.: 9.00
Median : 444
Mean   : 1423
3rd Qu.: 1480
Max.   : 71188
              unknown :1324    Median :16.00
              Mean   :15.92
              3rd Qu.:21.00
              Max.   :31.00

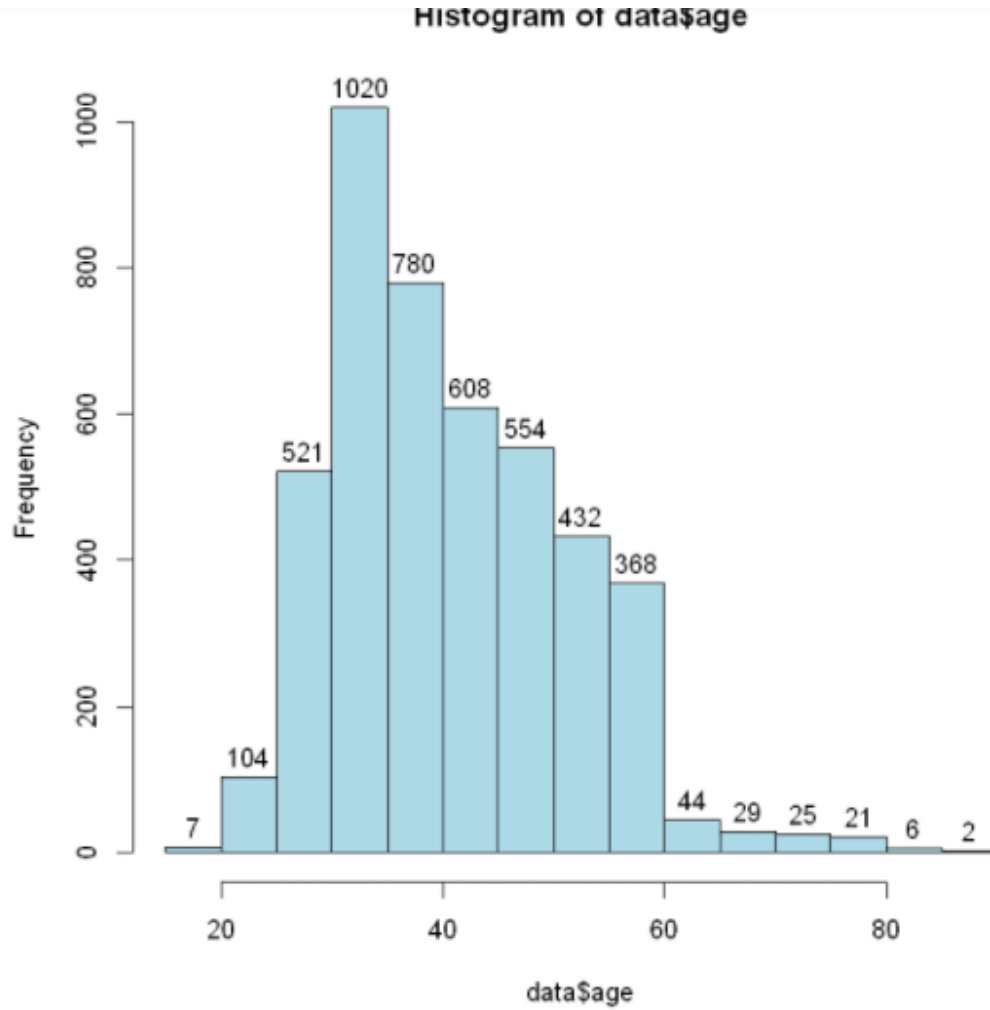
      month                duration                campaign                pdays
may    :1398    Min.   : 4    Min.   : 1.000    Min.   : -1.00
jul    : 706    1st Qu.: 104    1st Qu.: 1.000    1st Qu.: -1.00
aug    : 633    Median : 185    Median : 2.000    Median : -1.00
jun    : 531    Mean   : 264    Mean   : 2.794    Mean   : 39.77
nov    : 389    3rd Qu.: 329    3rd Qu.: 3.000    3rd Qu.: -1.00
apr    : 293    Max.   :3025    Max.   :50.000    Max.   :871.00
(Other): 571

      previous                poutcome                y
Min.   : 0.0000    failure: 490    no :4000
1st Qu.: 0.0000    other  : 197    yes: 521
Median : 0.0000    success: 129
Mean   : 0.5426    unknown:3705
3rd Qu.: 0.0000
Max.   :25.0000

```

**Summary** fonksiyonu çeşitli model oluşturma fonksiyonlarının sonuç özetlerini üretmek için kullanılan genel işlevli bir fonksiyondur. Summary fonksiyonu ile sayısal değerli niteliklerin önemli istatistiksel değerleri elde edilmiş olur. Sayısal değerlerin en küçük değeri (min), en yüksek değeri (max), ortalaması (mean), ortanca değeri (median), ortanca değerinden küçük olan değerlerin ortancası (1st Quartile), ortanca değeriden büyük olan değerlerin ortancası (3rd Quartile) gibi önemli istatistiksel değerleri elde edilir. Kategorik değerli niteliklerin ise her bir kategoriden kaç tane olduğunun özeti verilmektedir. Örneğin yaş niteliğine (sayısal) bakıldığında müşterilerden en genç kişi 19 yaşında iken, en yaşlı kişi 87 yaşındadır. 4521 müşterinin yaş ortalaması 41.17'dir. Median değeri 39 iken, 1st Quartile 33 ve 3rd Quardile değeri 49'dur. Evlilik niteliğine (kategorik) bakıldığında ise 528 tane boşanmış (divorced), 2797 tane evli (married), 1196 tane de bekar (single) müşteri olduğu görülmektedir.

```
hist(data$age, col = "light blue", freq = TRUE, labels=TRUE)
```

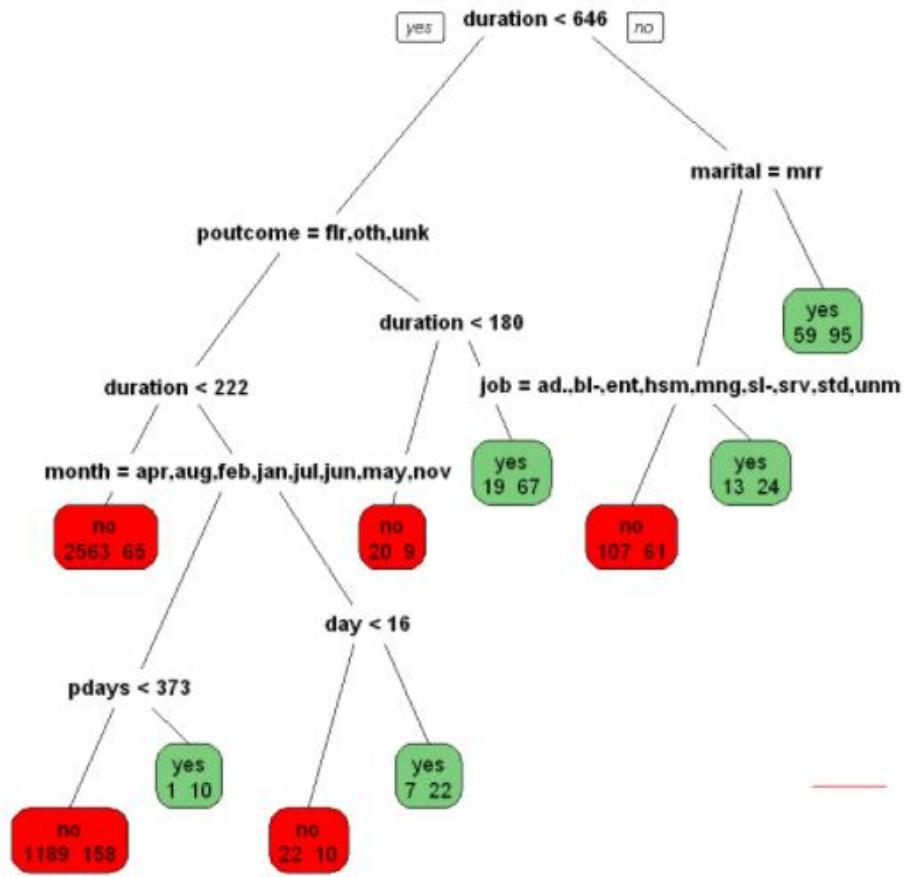


**hist** fonksiyonu ile yaş-frekans arasındaki ilişki gösterilerek hangi yaş grubunda kaç tane müşteri olduğu görülmektedir.

Müşteri kararının tahmin edilmesi için karar ağacı algoritmalarından faydalanılmıştır. R'ın kendisine ait olan rpart ve party kütüphaneleri kullanılmıştır. Brieman, Friedman, Olshen ve Stone tarafından 1984 tarihinde geliştirilen sınıflandırma ve regresyon ağaçları(classification and regression trees(CART)) rpart paketi aracılığıyla üretilir.

```
library(rpart)
library(rpart.plot)
data<-read.table("Banka.txt",header=T,sep="\t")
elif<-rpart(y~.,data=data,method="class")
prp(elif,cex=0.9,box.palette=c("red","palegreen3"),extra=1)
legend("bottomright",legend=c("yes","no"),box.col=c("red","green"))
```





rpart ve rpart.plot paketleri R'nin çalıştığı dizine `install.packages("rpart")` ve `install.packages("rpart.plot")` kodları ile yüklenmiştir.

### rpart Genel Kullanımı

rpart ile sınıflandırma ve regresyon analizleri yapılmaktadır. `rpart(formula,data,method)`. Parametrelerin kullanımı;

**formula** =sonuç~nitelik1+nitelik2+..... Banka verisinde sonucumuz y alanında saklıdır. 16 tane nitelik olduğundan tek tek yazmak yerine . yazılmıştır. **data**=Veri kümesi, **method** ="class" (sınıflandırma ağacı için), "anova" (regresyon ağacı için) Sınıflandırma ağacı oluşturulacağı için class yazılmıştır. Yaygın kullanımı 3 parametrelidir. `control`, `weights`, `subset` şeklinde opsiyonel parametreler de alabilir.

### prp Genel Kullanımı

`rpart.plot` paketinin içinde kullanılan bir plot fonksiyondur. `prp( x ,cex,box.pallet, extra )` Parametrelerin kullanımı; `x` = rpart nesnesidir. Tek zorunlu argumdur. `cex` ile yazı boyutu ayarlanır. `box.palette` ile düğüm sonucunda bulunan farklı değerler renklendirilir. `extra` ile düğüm sonucunda bulunan sayısal değerler yazılır.

### **Karar Ağacının Yorumlanması:**

- ✓ Ağacın kök düğümünü mevduatı satmaya çalışan banka çalışanlarının müşterilerle telefonda yaptığı görüşmenin süresi(duration) oluşturmaktadır. Müşteri sayısı 4521 olan bir ortamda müşterilerin 4162'si telefonla 646 sn'den daha az konuşurken 359'u 646 sn'den daha fazla konuşmuştur.
- ✓ 646 sn'den daha az konuşan 4162 müşterinin yaklaşık 4047'si bir önceki kampanyadaki ürünü almamış (failure), ürünü alıp almadığı bilinmeyen (unknown) ve diğer durumda (other) olan müşteriler iken 115'i satın alan (success) müşterilerdir. Bu 4047 müşterinin 2628'nin telefonla yaptığı görüşme süresi 222 sn'den daha azken 1419'nun telefonla yaptığı görüşme süresi 222 sn'den daha fazladır. Telefonla yaptığı görüşme süresi 222 sn'den daha az olan 2628 müşterinin 2563'ü mevduatı satın almazken 65'i mevduatı satın almıştır (1. yaprak düğüm).
- ✓ Görüşme süresi 222 sn'den daha fazla olan 1419 müşteriden 1358'i bu telefon görüşmesini ocak, şubat, nisan, mayıs, haziran, temmuz, ağustos ve kasım aylarında gerçekleştirmiştir. Bu 1358 müşteriden 1347'sinin bir önceki kampanya çalışması için aranmasının üzerinden 373 günden daha az bir zaman geçmiştir. Bu 1347 müşteriden 1189'u mevduatı satın almazken 158'i mevduatı satın almıştır (2. yaprak düğüm). Bu 1358 müşteriden 11 kişinin bir önceki kampanya çalışması için aranmasının üstünden 373 günden fazla bir zaman geçmiştir. Bu müşterilerden 10'u mevduatı satın alırken 1'i mevduatı satın almamıştır (3. yaprak düğüm).
- ✓ Görüşme süresi 222 sn'den daha fazla konuşan 1419 müşterinin 61'i bu telefon görüşmesini mart, eylül, ekim ve aralık aylarında gerçekleştirmiştir. Bu 61 müşteriden 32'si mart, eylül, ekim ve aralık aylarının 16'sından önce aranırken 29'u bu ayların 16'sından sonra aranmıştır. 16'sından önce aranan 32 müşteriden 22'si mevduatı satın almazken 10'u mevduatı satın almıştır (4. yaprak düğüm).
- ✓ 16'sından sonra aranan 29 müşteriden 22'si mevduatı satın alırken 9'u mevduatı satın almamıştır (5. yaprak düğüm).
- ✓ Görüşme süresi 646 sn'den daha az olan ve bir önceki kampanyadaki ürünü satın aldığı bilinen 115 müşteriden 29'u telefonla 180 sn'den daha az konuşurken 86'sı 180 sn'den fazla konuşmuştur.(180 ile 646 sn arasında). 180 sn'den az daha az konuşan 29 müşteriden 20'si mevduatı satın almazken 9'u satın almıştır (6. yaprak düğüm).

- ✓ 180 sn'den fazla konuşan 86 kişiden 67'si mevduatı satın alırken 19'u satın almamıştır (7. yaprak düğüm).
- ✓ Görüşme süresi 646 sn'den daha fazla olan 359 müşteriden 205'nin medeni durumu evli değı iken 154'nün medeni durumu evlidir . Bu 205 müşterinin 168'nin meslekleri yönetici(admin), mavi yakalı(blue-collar), girişimci(entrepreneur), hizmetçi(housemaid), idareci(management), kendi işi olan(self employed), hizmet eden(services), öğrenci(student) veya işsiz(unemployed)'dir. Bu 168 müşteriden 107'si mevduatı satın almazken 61'i satın almıştır. (8. düğüm)
- ✓ Bu 359 müşteriden 37'sinin ise meslekleri bilinmeyen(unknown), emekli(retired) ve teknisyen(technician)'dir. Bu müşterilerden 24'ü mevduatı satın alırken 13'ü satın almamıştır. (9. düğüm)
- ✓ Medeni durumu evli olan 154 müşteriden 95'i mevduatı satın alırken 59'u mevduatı satın almadır. (10.düğüm)

## Veri Kümesini Bölme ve Tahminleme

```
size<- nrow(data)*0.8
index<-sample(1:nrow(data),size=size)
training<-data[index,]
test<-data[-index,]
head(training)
head(test)
nrow(training)
```

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcc
1066	30	unemployed	married	tertiary	no	0	yes	no	cellular	18	nov	756	1	-1	0	unknc
1700	49	technician	married	unknown	no	323	yes	no	cellular	18	nov	41	2	-1	0	unknc
3027	39	management	divorced	tertiary	no	26	no	no	unknown	18	jun	311	5	-1	0	unknc
4194	48	services	divorced	secondary	no	3186	no	yes	cellular	9	apr	104	1	-1	0	unknc
3755	36	management	single	tertiary	no	2944	no	no	cellular	18	aug	882	8	-1	0	unknc
2796	36	technician	divorced	tertiary	no	1174	yes	no	cellular	18	nov	192	2	-1	0	unknc

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcc
6	35	management	single	tertiary	no	747	no	no	cellular	23	feb	141	2	176	3	failure
13	36	technician	married	tertiary	no	1109	no	no	cellular	13	aug	328	2	-1	0	unknow
26	41	management	married	tertiary	no	5883	no	no	cellular	20	nov	182	2	-1	0	unknow
31	68	retired	divorced	secondary	no	4189	no	no	telephone	14	jul	897	2	-1	0	unknow
34	32	management	single	tertiary	no	2536	yes	no	cellular	26	aug	958	6	-1	0	unknow
38	32	blue-collar	married	secondary	no	2089	yes	no	cellular	14	nov	132	1	-1	0	unknow

Daha önce Karar Ağaçları'nda bahsedildiği üzere; Karar ağacını tekniğini kullanarak verinin sınıflanması, öğrenme ve sınıflama olmak üzere iki basamaklı bir işlemdir. Öğrenme basamağında önceden bilinen bir eğitim verisi, model oluşturmak amacıyla sınıflama algoritması tarafından test edilir. Sınıflama basamağında ise test veri seti ile karar ağacının doğruluğu test edilir. 4521 örneklik veri kümesi yapılması amacıyla %80 oranında (3616) eğitim (training) ve %20 oranında (905) da test veri kümesine bölünmüştür. Training veri seti ile model oluşturulurken, test veri seti ile de bu modelin doğruluğu test edilecektir. Sample fonksiyonu ile toplam veri üzerinden örneklem oluşturulmuştur. Training ve test veri kümelerini data veri kümesi üzerinden rastgele (random) seçilen örnekler oluşturmaktadır. Yani kodları her çalıştırdığımızda farklı örnekler gelmektedir. Örnek olması amacıyla daha önce de kullanılan head fonksiyonu ile ilk 6 örnek gösterilmiştir.

```
elif<-rpart(y~.,data=training,method="class")
prediction<- predict(elif,test,type="class")
head(prediction)
```

```
6    no
13   no
26   no
31  yes
34  yes
38   no
```

### predict Genel Kullanımı

predict fonksiyonu rpart nesnesinin tahmin edilen sonucunu bir vektör olarak döndürür. predict (object ,newdata, type=c("vector","prob","class","matrix") Parametrelerin kullanımı;

**object:** rpart sınıfından rpart fonksiyonu tarafından oluşturulan nesnedir. Banka veri kümesi için rpart fonksiyonundan üretilen nesnemiz elif idi. **newdata:** Tahmin ediciler için gerekli olan değerleri içeren veri çerçevesidir. Kullanılacak newdata %80 oranında bölme işlemi yapılan test veri kümesidir. **type:** Tahmin edilen değerlerin döndürüldüğü karakter dizisidir. Vector, prob, class, matrix olarak 4 çeşidi vardır. Tercih olarak "class" kullanılmıştır. 6 örnek için bakıldığında; tahmin ediciler kullanılarak ; 31 ve 34 numaralı müşterilerin ürünü satın alacağı tahmin edilirken, 6, 13, 26, 38 numaralı müşterilerin ürünü satın almayacağı tahmin edilmiştir. Peki gerçekte de durum bu şekilde mi?

```
df<- data.frame(data)
df[31,]
df[13,]
df[34,]
```

	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
31	68	retired	divorced	secondary	no	4189	no	no	telephone	14	jul	897	2	-1	0	unknown	yes
	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
13	36	technician	married	tertiary	no	1109	no	no	cellular	13	aug	328	2	-1	0	unknown	no
	age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	y
34	32	management	single	tertiary	no	2536	yes	no	cellular	26	aug	958	6	-1	0	unknown	yes

31, 13 ve 24 numaralı müşterilere bakıldığında tahminlerin doğru olduğu gözükmektedir.

```
confusion.matrix <- table(prediction, test$y)
confusion.matrix
```

```
prediction  no yes
no      768  72
yes     26  39
```

*Confusion Matrix* ilişkili istatistiklerle gözlemlenen ve tahmin edilen sınıfların çapraz tablolarını hesaplar.

- ✓ Test veri kümesinde toplamda 840 müşterinin mevduatı satın almadığı (no) bilinirken; bunların 768 tanesi satın almaz (no), 72 tanesi de satın alır (yes) şeklinde tahmin edilmiştir. (1.satır)
- ✓ Training veri kümesinde toplamda 65 müşterinin mevduatı satın aldığı bilinirken; bunların 26 tanesi satın almaz (no), 39 tanesi de satın alır (yes) şeklinde tahmin edilmiştir. (2.satır)
- ✓ 905 örneklemlili veri kümesinden  $768+39=807$  örnek doğru tahmin edilerek başarı oranı %89.17 olarak bulunmuştur.

## **2. BULGULAR (VARSA)**

### 3. TARTIŞMA VE SONUÇ

## KAYNAKLAR

Tezde kaynak gösterimi yazar-tarih sistemine göre olmalıdır.

Kaynaklar metin içinde (Yazar Soyadı, yıl) şeklinde gösterilmelidir.

Kaynak listesinde ise alfabetik sırada olmalıdır. Basılı dergide makale gösterimi:

Yazar Soyadı İsim baş harf/harfleri, (Yıl). Makale adı. *Dergi adı kısaltması (İtalik)*; **Cilt no (Koyu)** (Sayı no zorunlu değil): sayfa no-sayfa no.

Örnek olarak;

Hoogstraal H (1985). Argasid and Nuttalliellid as Parasites and Vectors. *Adv Parasit*; **24**: 135-238.

Melikoğlu G, Bitiş L, Meriçli AH (2004). Flavonoids of *Crataegus microphylla*. *Nat Prod Res* 2004; **18**: 211-213.

Diğer kaynakların kaynak listesinde gösterimi için “Bitirme Projesi Tez Yazım Klavuzu” sayfa 8’e bakılabilir.



## **EKLER**