

T.C.  
BEYKENT ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ

VERİ MADENCİLİKİNDE  
KATEGORİK VERİLERİN ÇIKARILMASI İÇİN  
AÇIK YÖNTEMLE YENİ BİR ALGORİTMA UYGULAMASI

YÜKSEK LİSANS TEZİ  
BURAK ÇAKIR

**Enstitü Ana Bilim Dalı: Matematik - Bilgisayar**

**Tez Danışmanı: Yrd. Doç. Dr. Gökhan S LAHTARO LU**

EYLÜL, 2008  
STANBUL

T.C.  
BEYKENT ÜNİVERSİTESİ  
FEN BİLİMLERİ ENSTİTÜSÜ MÜDÜRLÜĞÜ  
TEZLİ YÜKSEK LİSANS TEZ SINAV TUTANAĞI

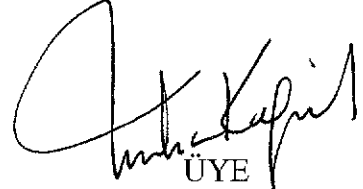
21.11.2008

Enstitümüz Matematik Bilgisayar Anabilim dalı Bilgi Teknolojileri Programı yüksek lisans öğrencilerinden 060862003 numaralı **Burak ÇAKIR'a** "Beykent Üniversitesi Lisansüstü Eğitim - Öğretim Yönetmeliği'nin ilgili maddesine göre hazırlayarak, Enstitümüze teslim ettiği "**VERİ MADENCİLİĞİNDE KATEGORİSEL DEĞİŞKENLER İÇİN VERİ KÜMELEMEDE AĞAÇ YÖNTEMİ İLE YENİ BİR ALGORİTMA UYGULAMASI**" tezini, Yönetim Kurulumuzun 13.10.2008 tarih ve 2008/17 sayılı toplantısında seçilen ve Fakülte binasında toplanan jüri üyeleri huzurunda, ilgili yönetmeliğin (c) bendi gereğince aday tarafından savunulmuş ve sonuçta adayın tezi hakkında *oybirliği* ile **Kabul** kararı verilmiştir.

İşbu tutanak, 4 nüsha olarak hazırlanmış ve Enstitü Müdürlüğü'ne sunulmak üzere tarafımızdan düzenlenmiştir.

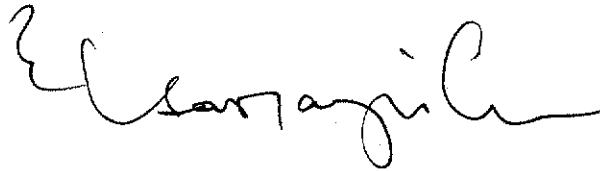
  
DANIŞMAN

YRD. DOÇ. DR. GÖKHAN SİLAHTAROĞLU

  
ÜYE

YRD. DOÇ. DR. TURHAN KARAGÜLER

ÜYE  
PROF. DR. ESAT HAMZAOĞLU



## ÖNSÖZ

Veri madencili i, ülkemizde henüz yeterli kaynak bulunmayan, eksikli i hissedilen bir alandır. Kurumların elinde çok ciddi veritabanları olu maya ba lamı tır. Hemen her disiplin tarafından çe itli amaçlar için kullanılan veri madencili i alanındaki geli meler, akademisyen ve ara tırmacı bilim çevreleri tarafından, uluslararası bilimsel dergiler ve kongreler aracılı ıyla takip edilmektedir.

Kümeleme ve a aç yöntemleri, veri madencili i alanında sık kullanılan yöntemlerdendir. Bu alanda sürekli yeni geli meler olmaktadır. Bu çalı mada, veri madencili i alanında kullanılmak üzere kategorisel de i kenlere sahip veri tabanlarında a aç yöntemiyle veri kümeleme için yeni bir algoritmanın kullanılabilirli inin gösterilmesi için hazırlanmı tır.

Bu tezin, veri madencili i konusu ile ilgilenenlere yardımcı olmasını ve farklı bir bakı açısı göstermede yararlı olmasını dilerim.

Tezin hazırlanmasında ve algoritmanın olu turulması a amasında benden yardımlarını esirgemeyen; veri madencili i konusunda ülkemizde büyük bir bo lu u dolduran de erli danı man hocam Yrd. Doç. Dr. Gökhan S LAHTARO LU'na te ekkür ederim. Tezin yazım a amasında bana sa ladı ı katkılardan ve gösterdi i ho görüden dolayı e im Aynur ÇAKIR'a te ekkür ederim.

## ÖZET

Veri madenciliğinde veri farklı sınıflar ve gürültü içerdiği zaman büyük popülasyonları ayıklama önemli bir problemdir. İyi bir algoritma mekanizması veya metodu kümeleri bulmak açısından etkili olmalıdır. Ayrıca boyut büyüdükçe uzayın karmaşıklığı ve zaman karmaşıklığı önemli hale gelmeye başlar.

Bu tez çalışmasında, veri madenciliği alanında kullanılmak üzere kategorisel veri kümelerine sahip veri tabanlarında Aaç yöntemiyle veri kümeleme için yeni bir algoritmanın kullanılabilirliğini gösterilmesi gerçekleştirilmiştir. Aaç yöntemi kullanarak sınıflara sahip oldukları veri kümelerine ayrıştırılmıştır. Bu veri kümeleri ve bölümler sayesinde bir veri tabanını mümkün olan en az parçadan en çok parçaya doğru sıralayarak bir Aaçolu yapı oluşturulması gerçekleştirilmiştir.

Çalışmanın birinci bölümünde veri madenciliği, kullanılan alanlar ve veri madenciliğinin gelişimi konularına değinilmiştir. İkinci bölümünde veri madenciliğinde kullanılan yöntemler, algoritmalar ve bu tezin konusu olan algoritmaya yakın olanları konu edilmiştir.

Tez çalışmasının üçüncü bölümünde ise söz konusu olan algoritma ve uygulanması gösterilmiştir. Kategorisel verilerden oluşan bir veri tabanı ile gerçekleştirilen sonuçlar elde edilmiş, aynı kategoride olan diğer algoritmalarla karşılaştırılması gerçekleştirilmiştir. Sentetik bir veri tabanı ile elde edilen sonuçlar gösterilmiştir.

Dördüncü bölümde ise, elde edilen bilgiler doğrultusunda, sonuçlar incelenmiş ve algoritmanın uygulanabilirliği hakkında yorumlar yapılmıştır.

Anahtar Kelimeler: Veri madenciliği, kümeleme, Aaç, ayırım, budama, entropi, gini

## ABSTRACT

The major reason that data mining became one of the hottest current technologies of the information age is the wide availability of huge amounts of data and the need for turning such data into useful information and knowledge. As computer systems getting cheaper and computer power increases, the amount of data available to be collected and processed increases. Therefore using techniques that operates very well with large amounts of data becomes an obvious choice. The information and knowledge gained can be used for applications ranging from business management, production control, and market analysis, to engineering design and science exploration.

In this study, a new data mining algorithm used and tested for categorical variable. This algorithm improved by Yrd. Doç. Dr. Gökhan S LAHTARO LU. This algorithm to call "A Tree Approach to Clustering Data with Categorical Variables". In the literature there are different approaches to form tree. To determine the best attribute, used an equal-split parameter. After forming the clusters, used another clustering algorithm such as PAM, CLARA or K-Means to reduce the number of leaves to the number required by the user.

Keywords: Data mining, clustering, tree, split, pruning, entropy, gini.

ÖNSÖZ.....	I
ÖZET.....	II
ABSTRACT.....	III
Ç NDEK LER.....	IV
EK LLER L STES .....	VII
TABLÖLAR L STES .....	VIII
BÖLÜM 1 .....	1
1.1 Veri Madencili i.....	1
1.2 Veri Madencili inin Uygulama Alanları.....	2
1.3 Di er Uygulamalar.....	3
1.3.1 aret leme.....	3
1.3.2 Biyoloji.....	3
1.3.3 Tıp.....	3
1.4 Veri Madencili i için Verilerin Hazırlanması.....	5
1.5 Veri Madencili i Modelleri.....	6
1.5.1 De er Tahmini Modeli.....	8
1.5.2 Ba lantı Analizi.....	10
1.5.3 Birliktelik Kuralları.....	10
1.5.4 Örüntü Tanıma.....	12
1.5.5 Ardı ık Zaman Örüntüleri.....	13
1.5.6 Kümeleme Analizi.....	14
1.6 Sınıflandırma Teknikleri ve Algoritmaları.....	16
1.6.1 Karar A açları.....	16
1.6.2 istatisti e Dayalı Algoritmalar.....	19

1.6.2.1 Bayesyen Sınıflandırma.....	19
1.6.2.2 Regresyon.....	19
1.6.3 Mesafeye Dayalı Sınıflandırma Algoritmaları.....	21
1.6.3.1 K- En Yakın Kom u.....	21
1.7 Yapay Sinir A ları.....	22
1.8 Birliktelik Kuralları ve li ki Analizi.....	22
1.8.1 AIS Algoritması.....	23
1.8.2 SETM Algoritması.....	24
1.8.3 Apriori Algoritması.....	25
1.8.4 AprioriTid Algoritması.....	26
BÖLÜM 2. ....	28
2.1 Kümeleme Analizi.....	28
2.1.1 Benzerlik ve Uzaklık.....	29
2.1.2 Kümeleme Analizinin Sınıflandırılması.....	31
2.1.3 Hiyerar ik Yöntemler.....	32
2.1.3.1 SLINK Algoritması ve Tek Ba lantı Tekni i.....	34
2.1.3.2 CURE Algoritması.....	34
2.1.3.3 CHAMELEON Algoritması.....	35
2.1.3.4 BIRCH.....	36
2.1.4 Bölümlemeli Yöntemler.....	38
2.1.4.1 K- Ortalama (K-Means) Algoritması.....	38
2.1.4.2 PAM Algoritması.....	40
2.1.4.3 CLARA Algoritması.....	40
2.1.4.4 CLARANS Algoritması.....	41
BÖLÜM 3.....	42

3.1 Çalışmanın Konusu.....	42
3.2 Çalışmanın Amacı.....	42
3.3 Çalışmaya Konu Olan Algoritma.....	43
3.4 Uygulama.....	43
3.4.1 Kullanılan Veri Tabanı.....	44
3.4.2 Delphi Programlama Dili ile Örnek Uygulama - 1.....	45
3.4.3 Excel'de Elde Edilen Sonuçların İncelenmesi.....	48
3.4.4 Algoritmanın ve Uygulamanın Sonuçları.....	49
3.5 Çalışmanın Önemi.....	53
3.6 Çalışmanın Kısıtları.....	53
3.7 Araştırmanın Modeli ve Hipotez.....	54
3.8 Algoritmaya Ait Özellikler.....	54
3.9 Hipotezin Testi.....	55
3.10 SPSS Programı ile Elde Edilen Sonuçlar.....	55
BÖLÜM 4 .....	60
4.1 Sonuç ve Öneriler.....	60
4.1.1 Genel Sonuçlar.....	60
4.1.2 Gelecek Araştırmalar için Öneriler.....	61
KAYNAKLAR.....	62



## EK LLER L STES

ekil - 1.1 Veri madencili i ve bilgi ke fi süreci.....	6
ekil - 1.2 Euclid uzayında veriler.....	15
ekil- 3.1 Delphi Programı Uygulaması - 1.....	45
ekil- 3.2 Delphi Programı Uygulaması - 2.....	50
ekil- 3.3 Budama E i inin 8 'den Küçük Oldu u Durum.....	52
ekil- 3.4 Budama E i inin 8 - 12 Arasında Oldu u Durum.....	52
ekil- 3.5 Budama E i inin 12 'den Büyük Oldu u Durum.....	53
ekil- 3.6 SPSS A acı.....	56
ekil- 3.7 Algoritmadan Elde Edilen A aç - 1.....	57

## TABLOLAR İÇİNDİRİMLERİ

Tablo - 1.1 Ürünler Tablosu.....	9
Tablo - 3.1 Algoritmadan Alınan İlk Değerler.....	44
Tablo- 3.2 Veri Tabanı Örneği.....	44
Tablo- 3.3 Algoritmanın Birinci Aşamasından Örnek Veriler.....	48
Tablo- 3.4 Algoritmanın İkinci Aşamasından Örnek Veriler.....	49
Tablo- 3.5 Değerlendirmede Kullanılan Veriler.....	50
Tablo- 3.6 SPSS Sonuçları.....	55

# 1. BÖLÜM

## 1.1 Veri Madencili i

Bilgisayarların ya amımıza daha çok girmesiyle birlikte, artık her yaptığımız işlem sayısal ortamda kayıt altına alınmaya başlandı. Marketlerde yaptığımız alışverişlerde aldığımız her bir ürün, hatta alıp bir süre sonra iade ettiğimiz bir ürün ve o ürünle birlikte aldığımız diğer ürünler bilgisayarlarda, veritabanlarında tutulmaya başlandı. Hastanelerde, belediyelerde veya ticaretle yaptığımız her işlem artık anında veritabanlarında yerini alıyor. Hatta bir mağazaya, alışveriş merkezine girerken ya da çıkarken, bazen de yolda yürürken kameraya çekilen görüntülerimiz bile bir veritabanı oluşturuyor. Bütün bunlar bir yılın halinde depolanırken içlerinde kim bilir ne gibi bilgiler gizlidir. Tüm bu veriler, veritabanlarında çıkarılmayı bekleyen değerli bir maden gibi durmaktadır. Bir bakıma etrafımız bir sürü veri varken bu veriler bilgiye dönüşmeyi beklemektedirler.

Veri madenciliği 1990'lerden beri, veri depolama araçları, barkod ve RFID teknolojilerine paralel olarak gelişmekte ve kullanım alanı yayılmakta olan bir konudur. Bu nedenle, kullanılan yer ve zamana göre çeşitli tanımları yapılmıştır. Çünkü her geçen gün daha da geliştiği için bugün yapılan bir tanım yarın yetersiz kalabilmektedir. En yaygın tanımlardan bir tanesi şöyle der: Veri madenciliği daha önceden bilinmeyen, geçerli ve uygulanabilir bilgilerin geniş veritabanlarından elde edilmesi ve bu bilgilerin işletme kararları verirken kullanılmasıdır.

Burada altının çizilmesi gereken noktalardan birincisi elde edilecek bilginin “önceden bilinmeyen” olmasıdır. Veri madenciliği sonunda ulaşılabilecek bilginin önceden bilinmiyor olmasından kasıt, elde edilecek sonucun tahmin edilmemesi anlamını taşımaktadır. Zaten tahmin edilebilen, beklenen sonuçlar için veri madenciliği kullanmak pek de ekonomik olmayacaktır. Ayrıca veri madenciliği tahmin edilen, öngörülen ya da başka yöntemlerle çıkarılmış sonuçların ispatını yapmak üzere kullanılacak bir araç değildir. Ayrıca, veri madenciliği daha önce hiç akla gelmemiş, düşünülmemiş sonuçları önümüze koymasıyla diğer

yöntemlerden farklılık gösterir. "Daha önce bilinmeyen" ya da tahmin edilemeyenle ilgili en ünlü örnek ise, artık klasikle mi , kulaktan kulağa anlatılan ve veri madenciliğinin "bilinmeyenini" çarpıcı bir şekilde önümüze koyan bira - çocuk bezi örneğidir:

Bir perakende mağazalar zincirinin yaptığı veri madenciliği araştırmasının sonuçlarına göre bira ile çocuk bezi satışları arasında, özellikle Cuma günleri, güçlü bir ilişki vardır. Çocuk bezi satın alan kişilerin büyük çoğunluğu aynı zamanda bira da satın almaktadırlar. Daha doğrusu, Cuma günleri çocukları için alıverişe çıkan babalar arada kendileri için de alıveriş yapmaktadırlar [Cabena, 1998].

Gartner Group tarafından yapılan bir değerlendirilimde ise veri madenciliği, istatistik ve matematik tekniklerle birlikte örüntü tanıma (pattern recognition) teknolojilerini kullanarak, depolama ortamlarında saklanmış bulunan verilerin elelenmesi ile anlamlı yeni korelasyon, örüntü ve ilişkilerin keşfedilmesi sürecidir [Akpınar, 2000]. Veri madenciliğini salt bir tanım olarak ele alırsak gerçekte var olan önemini tam olarak yansıtmamış oluruz. Salt bir tanımda veri madenciliği için bir tahmin aracı gibi yaklaşılabılır ya da veri madenciliği basit bir bilgisayar programı gibi görülebilir. Elbette ki sonuç olarak veri madenciliği kullanılacaksa bir bilgisayar yazılımı üzerinden kullanılacaktır; ancak veri madenciliğini tek bir programı anlamak ve onun arayüzünü kullanmak olarak ele almak bu bilim alanına büyük bir haksızlık olacaktır. Oysa kullanım alanlarına baktıkça durumun hiç de öyle olmadığını görülecektir. Bu nedenle öncelikle veri madenciliğinin kullanım alanlarına ve amaçlarını bilip daha sonra bu amaçlar için geliştirilmiş teknik algoritma ve yazılımları kullanmakta yarar vardır.

## **1.2 Veri Madenciliğinin Uygulama Alanları**

Veri madenciliği bankacılık, pazarlama, sigortacılık, sağlık gibi değişik alanlarda uygulanmaktadır. Veri madenciliğinin kullanılmasında sektör farkı gözlemlenmemekle beraber, geniş veri ambarlarının oluşturulmasına olanak veren, perakende satış, sigortacılık, sağlık gibi alanlarda kullanılması daha yaygın ve daha doğrudur.

## 1.3 Diğer Uygulamalar

Pazarlama ve risk yönetimi dışında veri madenciliği şu alanlarda da kullanılır:

### 1.3.1 Arayış

Telefon hatlarındaki parazitlenmeden dolayı oluşacak kayıpları ve buna bağlı olarak konu mada ortaya çıkan gürültüyü yok etme.

### 1.3.2 Biyoloji

DNA sıra (veri) analizi. İnsanda yaklaşık 100.000 gen vardır. Hastalıklara yol açan gen sıralama örneklerini binlerce gen arasından bulmak, tanımlamak oldukça zor bir iştir. Veri madenciliğiyle geliştirilen sıralama örnek analizi ve benzerlik arama yöntemleri DNA verisi üzerinde analiz yapmayı kolaylaştırır [Kut, Yılmaz].

### 1.3.3 Tıp

Bazı hastalıkların %100 kesin teşhisi mümkün olmamaktadır. Örneğin gebelik esnasında çocukta oluşabilecek herhangi bir down sendromu riskinin kesin tanısı dışı bulgularla sağlanamamaktadır. Buradaki dışı bulgulardan kasıt, anneden alınacak kan örneği, ultrason ile bebeğin görüntülenmesi, anne adayının yaşı, hamilelik ayı aldığı kilo vs gibi bulgulardır. Ancak bu bulguların hemen hiç biri hekime %100 tanı koyma olanağı vermez; %100 veya %100'e çok yakın bir tanı için anne karnından alınacak sıvının incelenmesi de gerekmektedir. Oysa bu işlemlerde 1/300 oranında bir düşük riski vardır. Dolayısıyla bu işleme girmeden önce hekimin anne karnındaki bebekte down sendromu olduğundan kuşlanması gerekmektedir. Bu amaçla yukarıda söz edilen dışı bulgular ve veri madenciliği teknikleri devreye girmektedir.

Daha önce bu işlem uygulanmış, dışı bulguları ve operasyon sonucu kaydedilmiş hasta adaylarına ait veritabanı, veri madenciliği algoritmaları tarafından incelenerek, bir makine öğrenmesi, sınıflandırma, karar ağacı vs. gerçekleştirilir. Daha sonra

gerçekle tirilen bu sisteme -örne in karar a acı- mevcut anne adayının bilgileri girilerek bebekteki risk oranı belirlenir. Bu oranın büyüklü üne ba lı olarak hekimin bir yarar risk analizi yapıp operasyona karar vermesi kolayla ır.

Tıp alanında bunun gibi ameliyat riski ta ıyan ancak, ameliyat öncesinde gerçekten ameliyat olması gerekti i tam olarak anla ılamayan hasta ve hastalıklar için de veri madencili i yöntemi kullanılır.

Ayrıca parmak izi tespiti, yüz eklinden kimlik tespiti, insan sesinin bilgisayar ve di er elektronik aygıtlarda komut olarak kullanılması konularında da kullanılan yapaya zeka teknikleri veri madencili i için de geçerlidir.

Amerika Bankası kendi ürünlerini kullanan banka mü terilerinin tespitinde veri madencili i kullanır ve mü teri ihtiyaçlarını kar ılamak için ürün ve servislerden olu an paketler sunar; Farmer Group irketinin "Risk Analiz Yönetim" paketi, sigorta oranı belirlenmesi, yatırım portföyü yönetimi, iyi ve kötü kredi riskleri ta ıyan irketlerin ve mü terilerinin belirlenmesinde veri madencili i kullanır; Twentieth Century Fox adlı film irketi ise fatura bilgilerini analiz ederek hangi aktörün, hangi filmin, hangi bölgede daha çok izlendi ini tespit ederek yeni film projelerini ba latmı ve bölge bazında gösterimler sunmu tur [Kut ve Yılmaz].

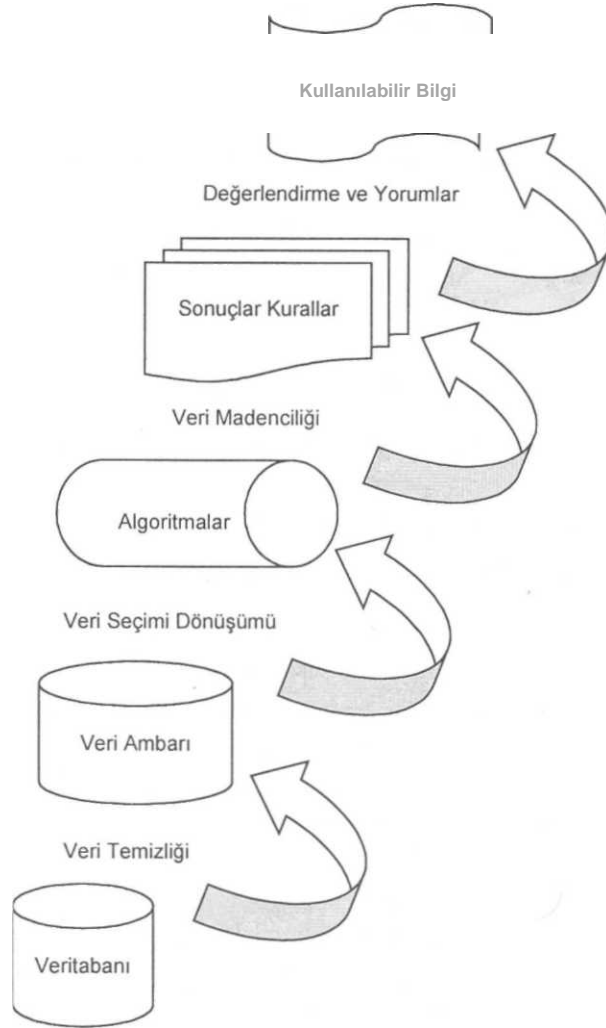
Görüldü ü gibi veri madencili i birçok ve birbirinden farklı konuyla ilgilenmekte, ba ka bir deyi le birbirinden farklı birçok konu veri madencili i yöntem ve teknikleri aracılı ı ile geli tirilmektedir. Bunun böyle olmasının ba lıca sebebi veri madencili i ve yapay zeka konularının -makine ö renmesi vs.- bir birlerine algoritma ve teknik olarak yakla malarıdır. Aslında veri madencili inin yapay zeka, makine ö renmesi adı altında geli tirilen algoritmaları kullanmakta ve ayrıca veri madencili i üzerine çalı anların geli tirdi i algoritmalar da yapay zeka alanına büyük yarar sa lamaktadır.

## 1.4 Veri Madencili i ğin Verilerin Hazırlanması

Veri madencili ğin ortaya ıkmasının ve gnmzde yaygın olarak kullanılıp bu konu zerine ara tırmalar yapılmasının en byk nedenlerinden bir tanesi, gnmzde byk veritabanlarının eri ilebilir olmasıdır. Bugn sper marketlerde yapılan alı veri lerden tutunda, di er kısım ve blmlerde alı an tm personelle ilgili her trl bilgi bilgisayarların belleklerinde tutulmaktadır. Ancak, bu veritabanlarındaki bilgilerin tamamının gerek ve do ru bilgiler oldu unu kimse %100 garanti edemez; ayrıca bu bilgilerin, mevcut haliyle yapaca ımız alı maya hizmet edece ği de kesin de ildir [Kimball]. Bu yzden, elimizdeki bilgilerin belirli i lemlerden geirilmesi gerekebilir.

Elimizdeki veritabanı bazı kayıtlar ynnden eksik olabilir. rne ğin, veritabanında kayıtlı birok ki ğinin medeni hali belirliyen, bu bilgi bazı kayıtlarda eksik olabilir; yani hi kayıt girilmemi olabilir, bu eksikli ği kayıp veriler (missing data) olarak isimlendirebiliriz. Bunun dı ğında, kayıtların bir kısmındaki bilgiler, a ırı u de erler ya da yanlış girilmi de erler olabilir; bunun en arpıcı rne ği bir ki ğinin do um tarihinin 1046 olarak girilmesi olabilir. Bu gibi bilgilere grlt ya da grltl veri denilir [Han-2001]. Bunun dı ğında, verilerin bir kısmı, gerekten yanlış , anlamsız bilgiler ierebilir; rne ğin rn kodları yanlış girilmi olabilir. Ayrıca, bir bilgi bir ka farklı yere gereksiz bir ekilde girilmi , aynı anlama gelebilecek birden fazla bilgi olabilir. rne ğin, kayıtlı ki ğilerin hem ya ları hem de do um tarihleri tutulmu sa bunlardan bir tanesi kesinlikle fazladır. Bazen mevcut de ği kenlerin birle mesi ve tek bir de ği ken gibi i leme ğirmesi mmkn olabilir; bu hem veri madencili ği alı ması esnasında bilgisayar alı ma zamanı karma ıklı ını azaltacak hem de elde edilecek sonuların gvenilirli ğini ve kalitesini arttıracaktır. Zaman karma ıklı ını nlemek iin tıpkı istatistik alı malarda yapıldı ğı gibi ana ktleden bir rnekleme alınarak elde edilen veri boyutu drlr. Ancak bu rneklemenin yapılabilmesi iin di er istatistik alı malarından farklı olarak elimizde ana ktle verilerinin tamamının olması gerekmektedir.

Veri madenciliği çalışmasının en başında yapılması gereken şey verilerin hazırlanmasıdır. Bu konuyu verilerin temizlenmesi ve verilerin yeniden yapılandırılması olarak iki başlık altında inceleyebiliriz.



ekil - 1.1 Veri madenciliği ve bilgi keşfi süreci

## 1.5 Veri Madenciliği Modelleri

Veri madenciliği konusu kullanıldıkları alanlara göre değişik modellere ayrılmaktadır. Bu modeller der tahmini model, veritabanı kümeleme modeli, bağlantı analizi ve fark sapmaları olarak dört ana başlık altında toplanabilir; bu modeller literatürde operasyon veya yöntem isimleriyle de anılmaktadır [Cabena, 1998]. Bu operasyonlar ya da yöntemler uygulamada değişik amaçlar için



kullanılırken birçok teknik ve algoritmalarından yararlanılmaktadır; kullanılan teknik ve algoritmalar genel olarak tahminleyici, tanımlayıcı veya her iki yaklaşımı da içerebilirler.

Herhangi bir amaç için, örneğin hedef müşteri kitlesinin seçimi gibi tek bir modelin kullanımını söz konusu olmayabilir; diğer modeller de farklı teknik ve algoritmalarla aynı amaç için kullanılabilirler.

Veri madenciliğinde kullanılan algoritma, teknik ve modeller sonuçta birer bilgisayar yazılımıdır. Bu yazılımlar her ne kadar matematiksel ve algoritmik olarak birbirlerinden ayrılırsa da bazı ortak özellikleri vardır; bu özelliklerden bir tanesi yazılımların öğrenme işlemidir. Yazılımlar kendilerine girilen verilerden bakıp inceleyerek diğer algoritmalarla bu verilerden bazı sonuçlar ve kurallar çıkarırlar. Bu inceleme işlemine ise "öğrenme" adı verilir. Daha sonra, bu çıkarımlar verilerin diğer kısmına uygulanarak sınanırlar. Yazılım bir bakıma kendi kendini sınav yaparak ne kadar öğrenimini sınar. Bulduğu sonuca göre de gerekirse çıkarımlarını (kurallar vs) yeniler. Yenilenen bu yeni bulguların ayrı bir işlemle doğrulanması gerekmektedir. Bu işlemi doğrulama denilir. Daha sonra aynı öğrenme olup olmadığını da kontrol edilmesi gerekir. Aynı öğrenme algoritmasının çıkardığı kuralların sadece üzerinde çalıştığı veriler için geçerli olması, dışarıdan başka veriler geldiğinde, üretilen kuralların geçersiz olması durumudur. Bu durumda, üzerinde çalışılan veriler üzerinde, denendiğinde tam sonuç veren, çıkarım ve kurallar, aynı performansla diğer veriler üzerinde gösteremezler.

Veri madenciliğinde sınıflandırma, kümeleme, bağlantı analizi ve dolandırıcılık tespiti gibi konularda kullanılmak üzere birçok algoritma geliştirilmiştir. Bu algoritmaların bazıları sadece sınıflandırma ya da kümeleme gibi konuları ilgilendirirken bazıları ise diğer versiyonlarla birden fazla konuda kullanılabilirler. Örneğin, genetik algoritmalar, yapay sinir ağları gerek sınıflandırma ve gerekse de kümeleme modellerinde kullanılabilirler; oysa Apriori algoritması sadece birliktelik kurallarının belirlenmesinde kullanılan bir algoritmadır.

### 1.5.1 De er Tahmini Modeli

De er tahmini ya da tahminsel modeldeki ö renme daha çok bir insanın ö renmesine benzemektedir. İnsan tüm yaşamı boyunca çevresini sürekli gözleyerek bir şeyler öğrenir. Aslında çevremizdeki her şey gerçek bir veritabanından başka bir şey değildir. Tahminsel model de kendisine verilen veritabanını inceleyerek, bu veritabanındaki temel unsurları birbirine benzeterek tanımlamaya, onları isimlendirmeye ve sınıflamaya çalışmaktadır. Tıpkı bir çocuğun kadın ve erkek cinsiyetlerini sınıflandırması gibi. Çocuk için ilk önce cinsiyet kavramı yani bir sınıflandırma yoktur; daha sonra anne-baba, teyze, hala, amca, kendinden büyük ve küçük erkek ve kız çocuklarını görür. Aslında tüm bunlar çocuk için bir veritabanıdır. Bu veritabanını inceleyen çocuk kadınla erkek arasındaki temel farkları belirler daha sonra kendisine hiç tanımadığı kız çocuğu gösterildiğinde bir önceki deneyimine/ öğrenmesine dayanarak bunun kız olduğuna karar verir. Aslında yaptığı tamamen bir sınıflandırma, genelleme yapma işlemidir; daha doğrusu bir tahmindir. Bu şekildeki bir öğrenmeye de denetimli öğrenme adı verilir [Cabena, 1998].

Denetimli öğrenmede, öğreniciye (öğrenici bir algoritmadır) nesnelere ve nesnelere özellikleri ve yine bu nesnelere tanımlanmış, ilerideki amaçta tahmini istenecek olan özellikleri verilir. Veri madenciliğindeki nesnelere günlük yaşamdaki nesnelere gibidir; ancak, buradaki nesnelere veritabanındaki her bir kayıttır. Şekil-2.3'deki ürünler tablosu nesnelere oluşmaktadır. Örneğin, veritabanındaki su nesnesinin özellikleri veritabanına girilmiştir. Normal yaşamda yukarıda belirtilen, bir çocuğun öğrenme işleminde, gördüğü nesnelere özelliklerini değerlendirilir. Veri madenciliğinde de yine, tasarlanan algoritma ya da program veritabanına girilmiş olan nesnelere özelliklerini değerlendirir. Günlük yaşamda, çocuk su ile bardak nesnesinin birbirinden farklı nesnelere olduğunu sürekli görerek, birden fazla defa bu nesnelere maruz kalarak öğrenir. Hatta bulaşıkta bu iki nesneyi iaretlerle büyüklere tanıtmaya çalışır; bu nesnelere tanır, fakat isimlerini bilmez. Çünkü onun için nesnelere özelliklerinin önemi vardır; isimlerinin değil! Daha sonra, bu nesnelere sabun, deterjan, su vs gibi isimlerini de öğrenerek bu isimleri nesnelere özellikleriyle özdeşleştirir. Bundan sonra herhangi bir nesnenin bir veya birkaç özelliği verilince bu nesnelere isimlerini tahmin eder. Veri madenciliği algoritmaları

da veritabanındaki ya da daha do ru bir deyi le olu turulmu veri ambarındaki nesnelerin özelliklerini, nesnelerin isimleriyle ili kilendirerek bu nesnelerin birbirinden farklı ya da benzer, aynı sınıftan nesneler olduklarını bulur ve ö renir. Daha sonra, kendisine verilen de i ik özellikleri de erlendirerek bu özelli e sahip olan nesnenin ismini tahmin eder.

Adı	Renk	Miktar	Kategori	Tip
Su	Berrak	0,5lt	Gıda	sıvı
Deterjan	Ye il	330 gr.	Temizlik	jel
Sabun	Kırmızı	150 gr.	Temizlik	katı
Bardak	Muhtelif	10 Pak.	Piknik	Plastik

**Tablo - 1.1 Ürünler tablosu**

Denetimli ö renmenin tersi duruma ise denetimsiz ö renme denilir. Denetimsiz ö renmede nesnelerin özellikleri verilirken tahmin için kullanılacak herhangi bir parametre verilmez; yani nesnelerin isimleri verilmez. Örne in, yukarıda verilen ilk örnekte, iki cinsiyet arasındaki farkı ö renen çocu a, erkek-kadın parametreleri, di er örnek için su, deterjan, bardak vs de erleri denetimsiz ö renme için verilmez.

Tahmini yakla ım, uygulamada, hedef mü teri kitlesinin seçimi, pazar sepeti analizi, çapraz satı , mü teri ili kileri yönetimi gibi alanlarda kullanılır. Bu yakla ımdan yola çıkarak geli tirilen en temel iki teknik, karar a açları ve sınıflandırmadır. Ayrıca genetik algoritmalar, yapay sinir a ları, Bayes yöntemi gibi yöntemler de bu amaç için kullanılırlar. Bu model çerçevesinde geli tirilen teknik ve algoritmalar ileride ayrıntılı bir ekilde ele alınmaktadır.

### 1.5.2 Ba lantı Analizi

Tahmini modelde kullanılan yazılım kendisine verilen veritabanının bir bütün olarak dü ünür ve ö renmesini de bu bütünü temel alarak gerçekleştirir. Oysa ba lantı analizinde veritabanındaki her bir kayıt veya kayıtlar grubu arasında bir ba lantı, ili ki yaratılmaya çalışılır. Ba lantı analizi bir veritabanındaki kayıtlar ya da bir graf üzerindeki dü ümler arasında çok rastlanan kuralları ortaya çıkarır.

Ba lantı analizinin en çok kullanıldığı alanlardan bazıları şunlardır: Çapraz satı , stok fiyat hareketleri ve hedef müşteri kitlesinin belirlenmesi [Cabena, 1998]. Ba lantı analizi dört ana başlık altında incelenebilir. Bunlar birliktelik kuralları, örüntü tanıma, ardı ık zaman örüntüleri ve benzer zaman ke fidir.

### 1.5.3 Birliktelik Kuralları

Birliktelik kuralı' belirli türlerdeki veri ili kilerini tanımlayan bir modeldir. Bu yönden de tanımlayıcı bir modeldir. Herhangi bir ürün alındı ında bu ürünün yanında bir ba ka ürünün de satın alınması bir birliktelik kuralı verir. Ürünler ve bu ürünlerin birlikte alınmaları söz konusu olunca, hemen anlaşı lca ı gibi birliktelik kuralları daha çok perakendecilik sektöründe faaliyet gösteren işletmelerde uygulanmaktadır. Örne in bir süpermarkette yapılan alı veri lerin incelenip hangi ürünün hangi ürünle birlikte satın alındı ının belirlenmesi birliktelik kurallarını ilgilendirir. Bunun dı ında örne in ileti im a larında meydana gelen hataların belirlenmesinde de kullanılabilir [Dunham, 2003].

Bilindi i gibi veriler bilgisayar üzerinden sayısal olarak iletilirler; her bir verinin ta ıdı ı bir sayısal de er vardır. Bu sayısal de erler arasında ili ki kurularak hangi sayısal de erden sonra hangi sayısal de erin gelebileceğinin otomatik olarak belirlenip söz konusu sayısal de erin kaybolması, okunamaması gibi durumlarda bu sayının yerine gelebilecek sayısal de er (tahmini) otomatik olarak konularak ileti imin kopmaması sağlanır. Örne in, verilen mesaj ".gün onu göremedim" olsun, "gün" sözcü ünün önünde 2 karakter kaybolmu tur. Bu durumda yazılım, daha önceden inceledi i veritabanından çıkarttı ı birliktelik kurallarıyla bu eksikli i

kar ıla tıracak ve olası olarak "bugün, o gün" gibi olasılıkları bulacak ve otomatik olarak en muhtemel seçene i "..gün" yerine koyacaktır. Unutulmamalıdır ki! Hiçbir a amada insan müdahalesi yoktur. Yani "bugün, o gün" olasılıkları tamamen veritabanına girilmi belki de milyonlarca kayıt üzerinde, yine yazılımın yaptığı ı bir ke ifle bulunmaktadır. Bu yüzden de veri madencili i bilgi ke fi olarak da anılır.

Görüldü ü gibi birliktelik kuralları aslında olaylar arasındaki probabilistik korelasyonu tanımlar. Olaylar arasındaki korelasyon ise sık sık beraber gözlenen olaylardır. Herhangi bir veritabanında birliktelik kurallarının tanımlanması veritabanı bilgi ke fi sürecinin ilk adımıdır. Veritabanındaki herhangi bir X'in aynı zamanda Y'yi içermesi bir birlikteliktir. Bu durum çok bilinen bir örnekle öyle açıklanabilir: "Bira içeren %30 alı veri in, %2'si aynı zamanda çocuk bezi de içermektedir." Burada %30 güven seviyesini, %2 ise bu güven seviyesine olan deste i belirtmektedir [Joshi].

Bu ili kileri ortaya çıkartmaya çalı an bir analistin bir takım kurallar elde edebilmesi için minimum kabul edilebilir destek ve güven de erlerini belirlemesi gerekir. Birçok birliktelik kuralları algoritması "üret ve test et stratejisi" kullanır. Örne in Apriori algoritması k. döngüsünde, aday k-dizilerini belirler ve bunların güven seviyeleri de veritabanının taranmasıyla yapılır. Bulunan birçok güven seviyesi ise ço u zaman önceden belirlenmi olan güven seviyesinden dü ük çıkaca ı için bu güven seviyesi dü ük çıkan aday diziler ret edilirken, muhtemelen daha az sayıda olan di er adaylar kural üretmek için kullanılacaktır: Bunun dı nda kullanıcıdan, yani veri madencili i analizini yapacak ki iden ayrıca bir destek de eri alınır [Orhunbilge, 1999]. Destek de eri elde edilen bilgilerin mevcut benzer bilgiler içindeki oranını temsil eder. Minimum güven, minimum destek seviyeleri Apriori ve bu konuyla ilgili di er algoritmalarla birlikte ileride ayrıntılı olarak ele alınmaktadır.

### 1.5.4 Örüntü Tanıma

Örüntü tanıma, daha önce belirlenmiş bir model diyebileceğimiz çok boyutlu bir örüntünün veritabanındaki benzerlerini ya da 'en benzerini' aramaktır [Dunham, 2003]. Herhangi bir yazılı metni tanımak ya da o metnin çok benzerini bulmak örüntü tanımanın konusuna girer. Bunun dışında parmak izi, ses, yüz tanıma, kan hücrelerinin kararlaştırılması, el yazılarının tespiti gibi alanlarda da uygulanır. Dolayısıyla örüntüden kasıt el, yüz resim, çizim ve ses gibi varlıkların sayısal ortamda sergiledikleri ekildir.

Aslında, örüntü tanımada yapılan işlem bir çeşit sınıflandırmadır. Elimizde, ulaşılmaması gereken bir örnek vardır ve biz bu örneğin benzerini veya mümkünse aynısını aramaktayız. Bir algoritmayla bu örneğe benzeyenleri bir araya topluyor, daha sonra bunları en benzerden, en azadana sıralıyoruz.

Futbol maçlarında fanatizmi ve şiddeti sona erdirmek ya da en azından azaltmak için alınan önlemlerden bir tanesi de daha önceden olay çıkartmış belirlenmiş kimselerin statlara alınmamasıdır. Ancak kimlik tespitine rağmen, baskının kimliğini kullanarak ya da kalabalıktan istifade ederek, bu kişiler statlara girebilmektedirler. Bu kişilerin eldeki foto rafları statta çekim yapan bir kameradan alınan seyirci foto rafları ile sürekli olarak karşılaştırılarak, bu kişilerin yakalanması sağlanmaktadır. Burada kullanılan çeşitli örüntü tanıma algoritmaları vardır. Hemen hepsi 0 ve 1'lerden oluşan çeşitli matrisleri karşılaştırarak bir benzerlik aramaktadır.

Günümüzde, klavye ve fare gibi giriş cihazları yerlerini, doğal insan sesine bırakmaya başlıyorlar; bilgisayara herhangi bir dosyayı aç denildiğinde bilgisayar doğrudan gidip o dosyayı açıyor. Aslında bu işi yaparken, bilgisayar söylenen cümleyi anlamsal olarak algılamıyor. Sizin ağzınızdan çıkan sesle daha önceden kaydedilmiş sözcüklerin sayısal ortamdaki örüntülerini/seslerini karşılaştırarak örüntü tanıma işlemi yapılıyor ve bulunan örüntünün karşısındaki komut yerine getiriliyor[Joch, 1994].

Ba ka bir örnek basketbol oyunundan verilebilir. Topun, hangi oyuncu tarafından hangi yolları izleyerek hareket etti i belirlenerek, daha sonra belirli bir oyuncunun, söz geli i 3 ya da 5 saniye sonra ne yapaca ının önceden belirlenmesi bir örüntü tanıma uygulamasıdır.

Örüntü tanıma için geli tirilmi birçok algoritma vardır. "K-en yakın kom u algoritması", "do rusal sınıflayıcı", "üstsel sınıflayıcı" bunlardan sadece birkaç tanesidir. Bu algoritmalardan do rusal olanlar yalnızca do rusal i levleri yerine getirirken, do rusal olmayanlar her türlü i levi yerine getirebilir.

### 1.5.5 Ardı ık Zaman Örüntüleri

Yukarıda örüntü (pattern) sözcü ünün, herhangi bir çizim, ses, resim, parmak izi vs gibi bir ekil oldu undan söz edilmi ti. Bu örneklere ek olarak, bir kimsenin yaptı ı i ler de örüntü olarak tanımlanabilir. Örne in bir mü terinin süt, peynir ve ekmek satın alması bir örüntüdür. Bu noktadan hareket edilerek bir mü terinin birinci gün A ürünü, onu izleyen gün veya günlerden birinde B ürünü ve daha sonraki bir günde de C ürünü alması ise yine bir örüntü olu turacaktır. Ancak bu sefer birbirini izleyen, yani zaman içinde ardı ık olan bir örüntü olu turacaktır.

Veri madencili i, perakendecilik sektöründen gelen yo un ilgiyle geli imini sürdürmektedir. Barkod teknolojisindeki geli me, perakendecilik sektöründe faaliyet gösteren i letmelerin büyük miktarda ve sık denilebilecek aralıklarda olu an bilgileri, süratle depolamasına olanak sa lamı tır. Satı miktarlarını gösteren bu verilere sepet verisi diyebiliriz[Agrawal, 1995]. Bu veriler genel olarak i lem tarihi, satılan ürünler, i lem sıra numarası gibi bilgileri içerir. lem yapılan yerde, burası tipik olarak bir süpermarkettir; üye kartı ve klüp kartı diyebilece imiz kartlar kullanılıyorsa bu bilgiler arasında mü teri numarası da bulunabilir. Bu kayıtlar incelenerek, ardı ık zaman örüntüleri belirlenebilir. E er mü terilerin bir kısmı, önce bir CD oynatıcısı, daha sonraki bir zamanda ABC isimdeki bir filmi, ardından XYZ isimdeki bir filmi satın alıyorsa bu bir ardı ık örüntü olu turacaktır. Ba ka mü teriler ABC filmi ile XYZ filmi arasındaki zaman diliminde, farklı filmler ya da farklı ürünler alsalar bile onlarda bu örüntüye dahildirler.

### 1.5.6 Kümeleme Analizi

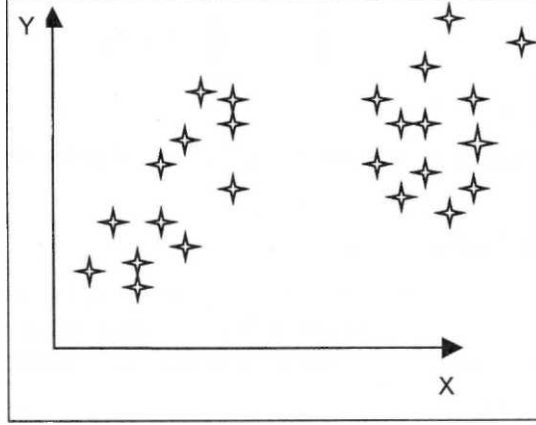
Kümeleme analizi (clustering) veri madenciliğinin en önemli alanlarından birisidir; amacı, nesnelere birbirlerine olan benzerliklerine göre gruplara ayırmaktır. Elde bulunan veriler incelenerek birbirlerine benzeyenler bir kümeye, benzemeyenler ise bir başka kümeye toplanırlar.

Verilerin kümeleme analizine göre modellenmesinde matematik, istatistik, makine öğrenimi ve yapay zeka gibi birçok alandan yararlanılır. Makine öğrenimi açısından, her bir küme gizli bir örüntüyü temsil eder ve uygulanan öğrenme ise bir denetimsiz öğrenmedir. Bu açıdan bakıldığında ise kümelemeyi, "gizli örüntülerin ortaya çıkartılması için uygulanan bir denetimsiz öğrenme yaklaşımı" olarak tanımlanabilir [Berkin].

Kümelemedeki öğrenmenin denetimsiz öğrenme olmasının nedeni önceden belirlenmiş sınıfların olmayışıdır. Önceden sınıflar belirli olsaydı, zaten bu kümeleme değil, bir sınıflandırma modeli adını alacaktı. Önceden sınıflar belirli iken, yani kadın ve erkek diye iki ayrı sınıf varken yapılan (algoritmik) öğrenmeye denetimli öğrenme; herhangi bir sınıf ismi verilmeden yapılan öğrenmeye ise denetimsiz öğrenme denilir. Örneğin, veritabanındaki kayıtlarda her kaydın yanında kadın veya erkek bilgisi yazılıyor olsun, bu durumda veritabanı üzerinde yapılan herhangi bir (kadın veya erkek olduğu) kural çıkarma işlemi denetimli öğrenmedir. Ancak aynı veritabanında, kayıtların yanında kadın mı erkek mi olduğu bilgisi yok iken yapılan kural çıkarma işlemi denetimsiz öğrenmedir. Bu işlem aynı zamanda veritabanını (iki) kümeye ayırma, yani kümeleme işlemidir. Burada kadın/erkek gibi bir etiket ya da sınıf olmayacağı için kümeleme kayıtlar arasındaki benzerlik veya mesafe ölçütüne göre yapılır.



ki verinin benzerli inden kasıt ise aralarındaki mesafenin ölçülmesi ve de erlendirilmesidir [Nanopoulos, 2001]. Bu de erlendirme, veritabanındaki di er verilere kıyasla iki verinin ne kadar yakın ya da benzer oldukları açısından yapılabilece i gibi önceden belirlenmi kısıtlar e ik de erleri çerçevesinde de yapılabilir.



ekil - 1.2 Euclid uzayında veriler

Örne in X ve Y özelliklerine sahip olan veriler toplulu u dü ünelim. Bu verileri X ve Y de erlerine göre iki boyutlu bir Euclid uzayında ekil-2.7 deki gibi gösterilebilir.

ekil - 1.2 de görülen noktaların birbirlerine olan mesafesi esas alındı nda toplamda iki kümenin olu tu u açıktır. Ancak, de i ken sayısı arttıkça bu kümeleri görmek iki boyutlu uzayda oldu u gibi kolay olmayacaktır. Örne in,  $A=\{1, 1, 2, 2, 5\}$  dizisine,  $B=\{1, 2, 3, 4, 2\}$  dizisi mi, yoksa  $C=\{1, 3, 5, 1, 3\}$  dizisi mi yakındır? Bunun yanıtını ekil 2.7 deki gibi bir grafikte görerek ya da göstererek vermek kolay olmayacaktır. Bu nedenle, diziler arasındaki mesafe matematiksel olarak ölçülmelidir. Yani mesafe (A, B) ile mesafe (A, C) bilinmelidir.

Öyleyse, bir grup noktanın küme olup olmadı mının belirlenmesi için bir mesafe ölçümüne ihtiyaç vardır. ki nokta arasındaki uzaklı ı mes (x, y) olarak gösterilirse, mes (x, y) bize x ile y arasındaki uzaklı ı verecektir. Genel olarak a a ıdaki gibi kabul edilir [Ulman].

$\text{mes}(x, y) = 0$  ise iki nokta arasındaki uzaklık sıfırdır.

$\text{mes}(x, y) = \text{mes}(y, x)$  ise mesafe simetriktir.

$\text{mes}(x, y) = \text{mes}(x, z) + \text{mes}(z, y)$  ise bir üçgen e itsizli i vardır.

## 1.6 Sınıflandırma Teknikleri ve Algoritmaları

Sınıflandırma en çok bilinen veri madencili i tekniklerinden birisidir; resim, örüntü tanıma, hastalık tanıları, dolandırıcılık tespiti, kalite kontrol çalı maları ve pazarlama konuları sınıflandırma tekniklerinin bolca kullanıldı ı alanlardır. Sınıflandırma tahminleyici bir modeldir; havanın bir sonraki gün nasıl olaca ı ya da bir kutuda ne kadar mavi top oldu unun tahmin edilmesi aslında bir sınıflandırma i lemidir [Dunham, 2003].

Sınıflandırma i leminde elimizdeki sınıf veya istatistiksel de imiyle ba ımlı de i ken hem sınıfsal hem de sürekli de er ta ıyabilir; bu anlamda regresyon ve çok terimli regresyona yakla maktadır [Akpınar, 2000]. Ancak, veri madencili i çerçevesinde bu istatistiksel yöntemlerin dı ında sınıflandırma i leminde "Bayesyen sınıflandırma algoritması", "karar a açlarına dayalı algoritmalar", "yapay sinir a ları" [Lipmann, i987] temelli algoritmalar ve "k-en yakın kom u algoritması" gibi birçok teknik ve algoritma geli tirilmi tir.

### 1.6.1 Karar A açları

Karar a açları sınıflandırma problemlerinde en çok kullanılan algoritmalarından birisidir. Di er yöntemlerle kıyaslandı ında karar a açlarının yapılandırılması ve anla ılması daha kolaydır denilebilir [Agrawal, 1993]. Bu teknikte sınıflandırma için bir a aç olu turulur; daha sonra, veritabanındaki her bir kayıt bu a aca uygulanır ve çıkan sonuca göre de bu kayıt sınıflandırılır. Temel olarak iki adımdan olu tu u

söylenbilir: Birincisi a acın kurulması, ikincisi de verilerin teker teker a aca uygulanarak sınıflandırmanın gerçekleştirilmesi şeklindedir.

Karar a acına dayalı olarak geliştirilen algoritmalar genel olarak a a ıda verilen kaba kod çerçevesinde çalışır:

D: Ö renme veritabanı

T: Kurulacak a aç

T = 0 // başlangıçta a aç boş küme

Dallara ayırma kriterlerini belirle

T = kök dü üümü belirle

T = dallara ayrılma kurallarına göre kök dü üümü dallara ayır;

Her bir dal için

do

Bu dü üüme gelecek de i kenisi belirle

if (durma koşulluna ula ıldı)

    Yaprak ekle ve dur

else

loop

Verilen kaba kodda sözü edilen durma/sonlandırma kriterini açıklamakta yarar vardır. Karar a açılırken eldeki veritabanının bir kısmı öğrenme işlemi için kullanılarak açılır; bu arada veritabanının bir kısmı da açılırken a açılırken test etmek için kullanılır. A açılırken kurulan sistemin çalışması belirlenir. Eğer a açılırken belirlenen düzeyde çalışmıyorsa dallanma durdurulur ve sınıflandırma tamamlanır. Programdaki durdurma kriteri a açma hassasiyetini de ortaya koyar. Geçerli durdurulan bir a açma daha fazla dallanacak ve a açma daha geniş olacaktır, çalışma süresi uzayacaktır. Bunun karışıklığında ise daha duyarlı sonuç verecektir. Erken durdurulan a açma ise her ne kadar daha hızlı çalışsa da tam öğrenmenin gerçekleşmesi her zaman tartışılmalıdır [Dunham, 2003].

A açma sırasında yapılan işlemlerden bir tanesi de budama işlemidir. Budama a açma sonucunu etkilemeyen ve sınıflandırmaya herhangi bir katkısı olmayan dalların a açma sonucunu azaltmasıdır. Bir bakıma gereksiz ayrıntıların sonuçtan çıkartılması işlemidir. A açma birçok düğüm ve dal oluşursa, a açma alt dalları ve yapraklarına ulaşan veri sayısı da azalacaktır; bu da a açma hassasiyetini azaltacaktır [Cabena, 1998]. Budamanın gerçekleştirilmesi için kullanılan algoritmanın işlemi budamanın hızı açısından önemlidir; ancak daha önemli bir unsur ise budamanın hangi ölçüte göre yapılacağını belirlemesidir. Yani biraz önce sözünü ettiğimiz “gereksiz ayrıntıların” ne olduğunu belirlemesidir.

Kullanılan algoritmaların çoğunda varsayılan (default) değer olarak %5 - %30 arası değerlerden düşük anlamlılık gösteren değerler budanırken, bu anlamlılığın belirlenmesi kullanıcıya bırakılmaktadır. Budama, gerek a açma kurulumu esnasında gerekse de kurulduktan sonra yapılabilir.

Karar a açma dayalı olarak geliştirilen birçok algoritma vardır. Bu algoritmalar birbirlerinden kök, düğüm ve dallanma kriteri seçimlerinde izledikleri yol açısından ayrılırlar.

## 1.6.2 Statistiki e Dayalı Algoritmalar

Veri madencili inde verilerin önceden verilen sınıflara göre ayrılması, aslında, gelecekte elde edilecek sonuçların tahminidir; yani sınıfların tahminidir. Regresyon, lojistik regresyon, zaman serileri analizi ve Bayesyen yaklaşım gibi istatistiksel yöntemler kullanılarak bu sınıflandırma i lemleri gerçekleştirilebilir.

### 1.6.2.1 Bayesyen Sınıflandırma

Bayesyen sınıflandırma tekni i, elde var olan, hali hazırda sınıflanmış verileri kullanarak yeni bir verinin mevcut sınıflardan herhangi birine girme olasılı mını hesaplayan bir yöntemdir. Bayesyen kuralına dayalı geli tirilmiş algoritma ve sınıflandırma teknikleri bu adla anılır.

### 1.6.2.2 Regresyon

Regresyon analizi herhangi bir de i kenin bir veya daha fazla ba ka de i kenler arasındaki ili kinin matematiksel bir denklem ekinde yazılmasıdır, yazılan bu denkleme regresyon denklemi adı verilir.

Regresyon, sınıflandırma için a a ıdaki iki yaklaşım çerçevesinde kullanılır [Dunham, 2003].

Bölme: Veriler sınıfa ba lı olarak çe itli bölgelere ayrılır.

Tahmin: Çıktı de erinin hesaplanması için formüller üretilir.

Bir ba ımlı de i kenin tek bir ba ımsız de i kenle açıklanabildi i durumlarda kullanılan regresyona "basit regresyon analizi" denilirken ba ımlı de i kenin birden fazla ba ımsız de i kenle açıklandı ı durumlarda kullanılan regresyona ise "çoklu regresyon analizi" denilir. Bunun dı ında kullanılan fonksiyonun, yani olu turulan denklemin türüne göre de bir ayırım yapılacak olunursa "do rusal" ve "do rusal

olmayan" regresyon analizi olarak yine ikiye ayrılabilir. Bu durumda, yukarıda söz edilen her iki ayırımın kombinasyonları da olacaktır ve basit doğrusal, basit doğrusal olmayan vs gibi ayrımlar da yapılabilecektir.

En küçük kareler yöntemiyle elde edilen doğrusal bir regresyon denklemi aşağıdaki gibi verilebilir:

$$y = a + b x + e$$

Burada  $a$  doğrusal fonksiyonun sabitidir;  $b$  ise doğrusal fonksiyonun eğimidir. Yani regresyon katsayısı olarak da adlandırılan  $x$ 'teki bir birimlik değişimin,  $y$  üzerinde, yine  $y$  cinsinden yaratacağı değişimi gösteren bir katsayıdır [Orhunbilge, 1999].  $y$  bağımlı değişkeni yani tahmin edilecek değeri temsil etmektedir. Veri madenciliği açısından ise  $y$  sınıfları temsil etmektedir,  $y$ 'nin aldığı değere göre verilen  $x$  değerinin hangi sınıfa dahil olacağı tahmin edilecektir,  $x$  ise bağımsız değişkeni temsil etmektedir; daha doğrusu sınıflanacak verinin giriş değeridir. Bu giriş değerleri belirli bir regresyon denkleminde yerine konularak bir  $y$  değeri hesaplanacak ve bu hesaplanan  $y$  değerine göre elimizdeki  $x$  değerlerinin hangi sınıfı temsil ettiğini tahmin edilmeye olacaktır.

Birden fazla  $x$  değeri için ise çoklu doğrusal regresyon denklemi aşağıdaki gibi olacaktır.

$$y = a + b_1x_1 + b_2x_2 + \dots + b_ix_i + e$$

Regresyon analizinde daha üst dereceden fonksiyonlar kullanılmaz. Bunun nedeni elde edilen verilere çok bağımlı sonuçlar elde edilecek olmasıdır. Veri madenciliğinde elde edilen verilere ayrı bağımlı sonuçlar elde edilmesine ayrı öğrenme denilir. Ayrı öğrenme durumunda elde edilen sonuçlar öğrenme için kullanılan verilere uygulandıında çok iyi sonuçlar verirken, dışarıdan başka veriler üzerinde kullanıldıında daha az hassas sonuçlar elde edilmesine yol açar [Cabena, 1998].

### 1.6.3 Mesafeye Dayalı Sınıflandırma Algoritmaları

Sınıflandırma yapılırken eldeki verilerin birbirlerine olan uzaklığı veya benzerliği kullanılarak sınıflamanın gerçekleştirilmesi bir başka sınıflandırma tekniğidir. Veriler arasındaki mesafe ölçülürken en çok kullanılan mesafe Euclid mesafesidir. Mesafeye dayalı algoritmalarından en bilineni "K-en yakın komşu algoritması"dır. Bu algoritmanın literatürde geliştirilmiş birçok türevi vardır. Genel olarak tümü aynı davranışı kullanarak işlem yaparlar.

#### 1.6.3.1 K- En Yakın Komşu

En yaygın algoritmalarından birisidir. İngilizce kaynaklarda K-Nearest Neighbor ya da KNN şeklinde ifade edilir. Sınıflandırma yapılırken veritabanındaki her bir kayıtdaki diğer kayıtlarla olan uzaklığı hesaplanır. Ancak, bir kayıt için diğer kayıtlardan sadece k adedi göz önüne alınır. Algoritmanın isminden de anlaşılacağı gibi bu k adet kayıt, başka bir deyişle veritabanındaki nokta, mesafesi hesaplanan noktadaki diğer kayıtlara nazaran en yakın olan kayıtlardır. Bu yöntem coğrafi bilgi sistemlerinde çok kullanılır; belirlenen bir noktaya en yakın şehir, istasyon vs belirlenmesi aslında k-en yakın komşu algoritmasının temelini oluşturur [Beyer, 1999].

Algoritmada k değeri önceden seçilir; değerinin yüksek olması birbirlerine benzemeyen noktaların bir araya toplanmasına, çok küçük seçilmesi ise birbirine benzemeyeni, yani aynı sınıfın noktaları oldukları halde, bazı noktaların ayrı sınıflara konmasına ya da o tür noktalar için ayrı sınıfların açılmasına neden olur. Tipik k değerleri 3.5 ve 7'dir [Khan, 2002].

Algoritmanın çalışma ilkesi aşağıda verilen adımlarla özetlenebilir:

Uygun bir mesafe ölçüm uzayı belirle. (Euclid mesafe ölçümü en çok kullanılmaktadır.)  
 Birbirine en yakın k adet noktayı belirle. Belirlenen grubun en çok rastlandığı sınıfı belirle. Bu gruba belirlenen sınıfın ismini ata.

Bugüne kadar, k-en yakın kom u algoritmasının daha etkin ve daha hızlı çalı an birçok uyarlaması tasarlanmı tır. Sınıflandırmanın dı nda kümeleme i lemlerinde de kullanılabilir.

## 1.7 Yapay Sinir A ları

Yapay sinir a ları biyolojik sinir a larından esinlenerek geli tirilmi bir bilgi i leme sistemidir. Yapay sinir a ları, yapay sinir hücrelerinin birbirleriyle çe itli ekilde ba lanmasından olu ur ve genel olarak katmanlar ekinde düzenlenir. En belirgin özellikleri birbirlerine ba lı nöronlar, ba lantılar arasındaki a ırlıkların belirlenmesi ve ate leme fonksiyonudur. Geçmi i 1942 yılına kadar gitmektedir.

## 1.8 Birliktelik Kuralları ve li ki Analizi

Veritabanlarındaki bilgi miktarı arttıkça birçok kurum ve kurulu sahip oldukları bilgiler arasındaki ili kileri ortaya çıkarma çabası içerisine girmi tir. Böylesi yı nlar halindeki bilgiler arasındaki ili kiler kurumlar için altın de erinde sonuçlar do urabilecek kararlarının alınmasında önemli rol oynamaktadır. li ki analizi veritabanındaki bir dizi bilgi ya da kaydın di er kayıtlarla olan ba lantısını açıklayan i lemler dizisidir. Yani bir kayıt varken, herhangi bir ba ka kaydın var olma olasılı ı nedir? Ya da bu iki kayıt varken, di er bir üçüncü hatta dördüncü kaydın veritabanına girme olasılı ı nedir? li ki analizi bu tür soruların yanıtım verir ve verilerin birlikte olan kurallarını ortaya çıkarır.

li ki analizi satı -pazarlamadan, ürün katalog tasarımlarına kadar birçok alanda kullanılmaktadır. Örne in, herhangi bir ürün satın alırken, bu ürünün yanında ba ka bir ürün ya da ürünlerin satın alınması, bu ürünler arasındaki ba lantıyı ortaya koyar. Bu ba lantıların ortaya çıkartılması ve bunun bir kural olarak ortaya konması ise birliktelik kuralları yani ili ki analizi konusuna girer. Literatürde bu tür çalı malara "pazar sepeti analizi" denilir. Pazar sepeti analizi mü terilerin alı veri alı kanlıklarının veritabanındaki bilgiler aracılı ıyla ortaya çıkartılması i lemidir. Bu alı veri alı kanlıklarının ortaya çıkartılması alı -veri merkezindeki ürünlerin



yerle tirilmesi, marketin alanının tasarımı ve markette sergilenecek ve satılacak olan ürünlerin belirlenmesinde yardımcı olur.

### 1.8.1 AIS Algoritması

AIS algoritması, geni nesne kümeleri üretmek için geli tirilmi bir algoritmadır. 1993 yılında geli tirilmi tir. Veritabanındaki isimlerin, yani nesne isimlerinin A 'dan Z'ye sıralanması kısıtını ta ır.

AIS algoritması veritabanını birçok kez tarar ve her tarama esnasında tüm i lemleri okur. İlk tarama esansında veritabanındaki nesnelere, teker teker sayarak hangilerinin geni nesnelere oldu unu belirler. Bunlardan geni olanlar aday nesne kümeleri olarak i aretlenirler. Bir i lem tarandıktan sonra, bir önceki taramada geni oldukları belirlenen nesne kümeleriyle, o i lemin nesnelere arasındaki ortak nesne kümeleri belirlenir. Belirlenen bu ortak nesne kümeleri i lemde mevcut olan di er nesnelere birle tirilerek yeni aday kümeler olu turulur. Herhangi bir 1 nesne kümesi, bir i lemdeki nesnelere birle ip aday kümelere birini olu turabilmesi için, birle ece i nesnenin hem geni olması hem de harf sırası açısından nesne kümesi içindeki tüm nesnelere sonra geliyor olması gerekir.

AIS algoritması bu adımı gerçekte tirebilmek için, bir budama tekni i kullanır. Budama tekni inin özünde, aday kümeler içindeki gereksiz kümelerin silinmesi vardır. Bu adımdan sonra, her aday kümesinin deste i hesaplanır. Destek seviyeleri minimum destek seviyesine e it veya bu seviyeden büyük çıkanlar geni nesne kümesi olarak i aretlenirler. Bir sonraki taramada bu geni i areti ta ıyan kümeler, yukarıda anlatıldı ı gibi bir sonraki aday kümelerin belirlenmesi için kullanılır [Agrawal, 1993].

### 1.8.2 SETM Algoritması

Bu algoritmada,  $L_k$  geni nesne kümesinin her bir elemanı iki parametreden oluşur; bunların bir tanesi nesnenin ismiyken diğeri bu nesneyi ayırt etmeye yarayan bir özellik numarasıdır. Algoritma içinde bu numara TID (Transaction Identification) olarak kullanılır. Benzer şekilde her bir aday nesne kümeleri,

$C_k, \langle \text{TID}, \text{Nesne Kümesi İsmi} \rangle$

formatında tutulur. [Houtsma ve Swami, 1993]

AIS algoritmasında olduğu gibi SETM algoritması da veritabanını birçok kez tarama. İlk tarama esnasında veritabanındaki nesnelere, teker teker sayarak hangilerinin genel nesnelere olduğunu belirler. Sonraki taramalarda, bir önceki geçişte genel olarak elde edilen nesne kümelerini kullanarak aday kümeleri belirler. Farklı olarak,

SETM algoritması aday kümelerle birlikte üzerinde çalışılan veritabanının TID bilgisini de tutar. Bundan sonra, aday nesne kümeleri nesne ismine göre sıraya dizilir ve küçük nesne kümesi silinir. Her veritabanı TID numarasına göre sıralanmışsa, bir sonraki tarama esnasında herhangi bir veritabanındaki genel nesne kümeleri  $L_k$ 'nin TID

numarasına göre sıralanmasıyla elde edilir. Bu şekilde veritabanı bir kaç kez tarama. Artık başka herhangi bir genel nesne kümesi bulunamadığında algoritma sonlandırılır.

SETM algoritmasında TID bilgisinin de tutulması, algoritmanın yer karmaşıklığını arttıracaktır, bu dezavantajın dışında başka bir eksi nokta ise, aday nesne kümesinin desteği hesaplanırken  $C_k$  sıralanmış halde değildir, bunun için nesne kümelerinin bir kez daha sıraya dizilmesi gerekecektir. Bu da zaman karmaşıklığını arttıran bir unsurdur.

### 1.8.3 Apriori Algoritması

Apriori algoritması bağlantı analizlerinin yapılabildiği kurallarının ortaya çıkartılması konusunda en çok bilinen ve kullanılan algoritmadır. Geni nesne kümelerinin ortaya çıkartılması için kullanılır.

Geni nesne kümelerini ortaya çıkartan algoritmalar eldeki tüm verileri birçok kez tararlar. İlk taramada, her bir nesnenin destek seviyesi, hesaplanarak kullanıcı tarafından bağlantıda girilen minimum destek seviyesi ile karşılaştırılır ve her bir nesnenin geniş olup olmadığına bakılır. Bundan sonraki her tarama bir önceki taramada geniş olarak tespit edilmiş nesnelere bağlantı ve geniş nesne kümeleri oluşturulur. Bu geniş nesne kümelerine aday nesne kümeleri denir. Taramanın sonunda ise hangi aday nesne kümesinin gerçekten geniş olduğu kontrol edilir. Daha önce de belirtildiği gibi bir nesne kümesinin geniş olarak adlandırılabilmesi için o nesne kümesinin kullanıcı tarafından verilen minimum destek seviyesinin üzerine bir destek seviyesine sahip olması gerekir. Bir sonraki taramada, yine bir önceki taramada geniş olarak seçilen nesne kümelerinden bağlantı ve veritabanının sonuna kadar bu nesne kümelerinin destekleri hesaplanır. Bu işlem, bağlantı yeni geniş nesne kümeleri bulunamayana kadar sürer. [Agrawal ve Spirant, 1994]

Apriori algoritması daha önceden ortaya atılmış olan AIS ve SETM algoritmalarından her bir geçişte aday nesne kümelerinin sayılması ve bu aday kümelerinin üretilmesiyle ayrılır. Hem AIS algoritmasında hem de SETM algoritmasında, tarama esnasında, veriler okunurken aday nesne kümeleri üretilir. Bir işlem (T) (transaction) okunduktan sonra, geniş nesne kümelerinin bu işlemlerde olup olmadığına da bakılır. Yeni aday nesne kümelerinin üretilmesi ise işlemlerdeki diğer nesnelere elde edilen geniş nesne kümelerinin birleştirilmesiyle üretilir [Agrawal, 1993]. Tabii bu da, gereksiz yere, aslında küçük nesne kümesi olan birçok aday nesne kümesinin sanki geniş nesne kümesiyymi gibi üretilmesi ve sayılması sonucunu doğurur. Bu da algoritmanın zaman karmaşıklığını artırır.

Apriori algoritması ise aday nesnelere üretirken veritabanındaki işlemleri hiç işine sokmadan, yalnızca bir önceki taramada geniş olduğu tespit edilmiş nesne

kümelerini kullanarak oluşturulur. Apriori algoritması genel bir nesne kümesinin herhangi bir alt kümesinin de genel olacağına dayanır. Böylece  $k$  adet nesneden oluşan bir nesne kümesi,  $k-1$  adet nesneye sahip genel nesne kümelerinin birleştirilmesi ve alt kümeleri genel olmayanların silinmesiyle elde edilebilir. Bu birleştirme ve silme işlemi sonunda daha az sayıda aday nesne kümeleri olacaktır.

Agrawal ve Srikant tarafından geliştirilen Apriori algoritması 1994 yılında 20. VLDB (Very Large Database Endowment) konferansında sunulmuştur. Bu bildiride, Agrawal ve Srikant algoritmanın çalışma ayrıntılarını ve algoritmanın kaba kodunu açıklanmıştır [Agrawal ve Srikant, 1994].

- Verilerin ilk taraması esnasında, genel nesne kümelerinin tespiti için, tüm nesnelere sayılır.
- Bir sonraki tarama,  $k$ 'nci tarama olsun, iki adımdan oluşur.
- Apriori-join fonksiyonu kullanılarak,  $(k-1)$ 'inci taramada elde edilen,  $L_{k-1}$ , nesne kümeleriyle,  $C_k$  aday nesne kümeleri oluşturulur.
- Sonra veritabanı taranarak,  $C_k$  'daki adayların desteği sayılır.
- Hızlı bir sayım için, verilen bir  $l$  değeri için,  $C_k$  'yı oluşturan adayların çok iyi belirlenmesi gerekir.

#### 1.8.4 AprioriTid Algoritması

Önceki ayrıtta belirtildiği gibi, algoritmalar desteği hesaplamak için tüm veritabanını tarar; ancak her adımda veritabanının tamamının taramasına gerek olmayabilir. Bu yaklaşımla Agrawal, Apriori algoritmasıyla birlikte AprioriTid algoritmasını da sunmuştur.

AprioriTid algoritması da taramadan önce aday nesne kümelerini belirlemek için ekil-4.2 de görülen apriori-gen fonksiyonunu kullanır. Apriori-gen en büyük farkı ilk geçi ten sonra veritabanının destek seviyesini bulmak için taranmamasıdır. Bu i için  $C_k$  kullanılır. SETM algoritmasında oldu u gibi  $C_k$  'nin her elemanı  $\langle TID,$

$\{X_k\}\rangle$  formundadır. Burada  $X_k$ , TID numaralı i leminde bulunan potansiyel geni k nesne kümesidir,  $k = 1$  iken  $C$ , veritabanına kar ılık gelir. Bununla beraber her nesne, nesne kümesiyle yer de i tir,  $k > 1$  oldu u durumlarda  $C_k$  algoritmanın onuncu adımında oldu u gibi üretilir, t i lemindeki  $C_k$  bir elemanı  $\langle TID, c \rangle$  eklindedir. Burada  $c$ , t i lemindeki  $C_k$  ya ait bir aday elemanıdır,  $\{c \in C_k \setminus c\}$ . E er bir i lemin,

herhangi bir k nesne kümesi adayı yoksa bu durumda  $C_k$  nm bu i lem için herhangi bir girdisi, elemanı olmayacaktır. Daha do rusu bu i lemin TID numarasını ta ımıyor olacaktır. Böylece  $C_k$  'daki girdi sayısı, özellikle bu k de erleri için, veritabanındaki i lem sayısından daha küçük olabilir. Bunun dı nda yine büyük k de erleri için her girdi kendisine kar ılık gelen i lemden daha küçük olabilir. Çünkü o i leminde çok az sayıda aday barınmıyor olabilir. Ancak, küçük k de erleri için bunun tersi olacaktır; yani girdiler kendilerine kar ılık gelen i lemlerden daha büyük olabileceklerdir [Agrawal ve Srikant, 1994].

## 2. BÖLÜM

### 2.1 Kümeleme Analizi

Kümeleme analizi, sınıflandırmada olduğu gibi sahip olunan verileri gruplara ayırma işlemidir. Kümeleme yabancı kaynaklarda clustering ya da segmentation olarak adlandırılmaktadır. Sınıflandırma işleminde, sınıflar önceden belirli iken kümelemede sınıflar önceden belirli değildir. Verilerin hangi gruplara/kümelere, hatta kaç de i ik gruba ayrılacağı eldeki verilerin birbirlerine olan benzerliğine göre belirlenir. Belirlenen her bir gruba da küme ismi verilir. Kümeleme analizi biyoloji, tıp, antropoloji, pazarlama, ekonomi ve telekomünikasyon gibi birçok ve birbirinden çok farklı alanlarda kullanılmaktadır [Dunham, 2003].

Örneğin elimizdeki bir perakende mağazasına ait veritabanında, müşterilerimizin sadece isim ve yaşlarının tutulduğunu varsayalım. Bu durumda, müşterilerimizi kümelere ayırmak istersek, onları yaşlarına göre ayırmamız doğru olacaktır. Dolayısıyla yaşları nispeten birbirlerine yakın olanlar aynı kümede toplanacaktır. Yaşları 20, 22, 26, 27, 40, 45, 46, 47, 49 ve 49 olan 10 müşteri varsa, 20-27 yaşları arasında olanlar birbirlerine, veritabanındaki diğer kişilere göre daha yakın olduğundan bir kümede toplanırken, yaşları 40-49 arası olanlar ise başka kümede toplanacaktır. Oysa veritabanındaki yaşları 19, 20, 21, 21, 21, 26, 26, 26, 27, 27 ve 28 ekinde olsaydı, 19 -21 arası bir kümede, 26 -28 arası da bir küme olması daha anlamlı olacaktı. Bu durumda, bir önceki örnekte 20 yaşındaki bir müşteriyle, 27 yaşındaki bir müşteri aynı kümede tutulurken, ikinci örnekte bunlar ayrı kümelere bulunacaklardır. Buradan da görüldüğü gibi veriler sadece tanıdıkları özelliklere göre değil, diğer verilerle olan benzerliğe ve veritabanındaki diğer verilerin durumuna göre de kümelere ayrılıyorlar. Bu durumda kümeleme sonuçları dinamik olabilir. Bu da kümelemeyi sınıflandırmadan ayıran bir başka özelliktir.

Kümelemenin matematiksel tanımı ise şu şekilde yapılabilir:

Elimizde  $D = \{ X_1, X_2, X_3, \dots, X_n \}$ ,  $n = 1, 2, \dots, m$  veritabanı olsun, her bir  $X_n$  bir kaydı temsil etsin.  $X = \{ x_1, x_2, x_3, \dots, x_i \}$ ,  $i = 1, 2, \dots, m$  her bir  $x_i$ , ad, soyad, yaş ve

gelir gibi özellikler olsun. Kümelemedeki amaç  $D$  veritabanını,  $j$  adet  $K$  kümesine bölmek ve  $K_j \subseteq D$  koşulunun sağlanmasıdır.

Daha önce belirtildiği gibi kümeleme analizinde, sınıflandırmadan farklı olarak belirlenecek kümelerin özellikleri önceden bilinmemektedir ve üstelik ortaya çıkacak küme sayısı da belirli değildir. Ancak, algoritmaların zaman karmaşıklığını, alınacak sonuçların kullanılabilirliğini artırmak için, literatürdeki algoritmaların bir kısmı ya küme sayısını ya da her bir kümede bulunacak eleman veya bu elemanlar arasındaki minimum - maksimum benzerlik uzaklık ölçütünü kullanıcıdan ister.

### 2.1.1 Benzerlik ve Uzaklık

Veritabanındaki veriler kümelere ayrılırken, benzerlik ve uzaklık kavramlarından yararlanır. Bu, veritabanındaki her bir kaydın diğer bir kayıtlarla olan benzerliğini ya da her bir kaydın veritabanındaki diğer kayıtlardan olan uzaklığı gibi oluşturulan gerçek ve aday kümeler arasındaki mesafe ve benzerliğini de içerir. Sözgelimi, veriler birbirlerine olan uzaklığına göre başlangıçta 8 ayrı kümeye ayrıldılarsa, bu 8 ayrı kümenin gerçekten farklı özelliklere sahip birer küme olup olmadığının da belirlenmesi gerekecektir. Bu durumda, oluşturulmuş bu kümeler arasındaki mesafe / benzerlik de ölçülmelidir. Birbirlerinden pek de farklı olmayan kümeler birleştirilerek tek bir küme haline dönüştürülebilir. Bu işlem, tüm veriler taranıp kümeler ortaya çıktıktan sonra yapılabileceği gibi veritabanının taranması ve veriler arasındaki benzerlik ve mesafenin ölçümü esnasında da yapılabilir. Bu nedenle kümeler arasındaki mesafenin ölçülmesiyle iki veya daha fazla kümenin birleştirilmesi söz konusu olduğu gibi aynı zamanda, bir kümeden birden fazla küme üretilmesi de söz konusu olacaktır. Bunun için de sürekli olarak kümelerin büyüklüğü ve çapı ölçülmelidir.

$D$  olarak göstereceğimiz bir veritabanındaki

$$D = \{X_1, X_2, X_3, \dots, X_n\}, n = 1, 2 \dots m,$$

$$X_m \text{ ile } X_j \text{ arasındaki mesafe, } mes(X_m, X_j)$$

Euclid uzayında, u ekilde hesaplanır:

$$\text{ben}(X_m, X_j)_{\text{DICE}} = \frac{2 \sum_{i=1}^n x_{mi} x_{ji}}{\sum_{i=1}^n x_{mi}^2 + \sum_{i=1}^n x_{ji}^2}$$

bu mesafeye Euclid mesafesi denir [E ecio lu, 2001].

Benzerlik kavramı ise mesafenin tersi bir anlam içerir ve iki veri arasındaki yakınlı 1

$$\text{ben}(X_m, X_j) = \frac{1}{1 + \text{mes}(X_m, X_j)} \quad \text{eklinde ifade edilebilir.}$$

Bununla beraber, benzerlik ölçümü için çe itli yöntemler vardır [Rasınussen, 1992] [Bacher]. Bunlardan birincisi Dice benzerlik ölçümüdür.  $\text{ben}(X_n, X_j)$ ,  $X_n$  ile  $X_j$  arasındaki benzerlik Dice benzerlik ölçüm yöntemiyle u ekilde hesaplanır.

$$\text{ben}(X_m, X_j)_{\text{DICE}} = \frac{2 \sum_{i=1}^n x_{mi} x_{ji}}{\sum_{i=1}^n x_{mi}^2 + \sum_{i=1}^n x_{ji}^2}$$

Aynı benzerlik Jaccard ile u ekilde hesaplanır;

$$\text{ben}(X_m, X_j)_{\text{JACCARD}} = \frac{\sum_{i=1}^n x_{mi} x_{ji}}{\sum_{i=1}^n x_{mi}^2 + \sum_{i=1}^n x_{ji}^2 - \sum_{i=1}^n x_{mi} x_{ji}}$$



Ba ka bir benzerlik ölçümü ise Cosine'dir. Cosine yöntemi ile a a ıdaki gibi hesaplanır;

$$\text{ben}(X_m, X_j)_{\text{COSINE}} = \frac{\sum_{i=1}^n x_{mi} x_{ji}}{\sqrt{\sum_{i=1}^n x_{mi}^2 \sum_{i=1}^n x_{ji}^2}}$$

Bu yöntemlerin dışında, benzerlik ölçümünde kullanılan bir ba ka yöntem ise Overlap'tir.

$$\text{ben}(X_m, X_j)_{\text{OVERLAP}} = \frac{\sum_{i=1}^n x_{mi} x_{ji}}{\min\left(\sum_{i=1}^n x_{mi}^2, \sum_{i=1}^n x_{ji}^2\right)}$$

### 2.1.2 Kümeleme Analizinin Sınıflandırılması

Literatürde birçok kümeleme algoritmasının adı geçmektedir. Algoritmalar birbirlerinden, kümelemenin olu turulu ekleline göre ayrıldıkları gibi kullanılan yeri türüne, yapılacak olan çalı manın amacına göre de farklılıklar gösterirler [Han ve Kamber, 2001]. Kümeleme algoritmaları, genel olarak hiyerar ik ye bölümlmeli olarak ikiye ayrılırken, bu konuda yapılmı bir literatür taraması bu algoritmaların daha alt bölümlere ayrılabilce ini göstermektedir [Berkhin].

- Hiyerar ik Yöntemler
- Topla ım (Agglomerative) Kümeleme Algoritmaları
- Bölünür (Dizsiye) Kümeleme Algoritmaları
- Bölümlmeli (Partitioning) Yöntemler

- Yer de ğ i tiren Algoritmalar
  
- Olasılıksal Algoritmalar
  
- K-Medoid Yöntemler
  
- K-Means Yöntemler
  
- Yo ğ unlu a Dayalı Algoritmalar
  
- Yo ğ unlu a Dayalı Ba ğ lantılı Kümeleme
  
- Yo ğ unluk fonksiyonlu kümeleme
  
- Grid Temelli Yöntemler
  
- Kategorik Verinin Yinelenmesine Dayanan Yöntemler
  
- Kısıtlara Dayanan Yöntemler
  
- Makine Ö ğ renmesi Alanında Kullanılan Yöntemler
  
- Gradient ğ nme ye Yapay Sinir A ğ ları
  
- Ölçeklenebilir Kümeleme Yöntemleri

### **2.1.3 Hiyerar ğ ik Yöntemler**

Hiyerar ğ ik yöntemler bir küme a ğ acı yaratır. Bu küme a ğ acındaki her bir dü ğ üm o ğ ullara sahiptir; a ğ aç yapraklarla son bulur. A ğ a ğ ıdan yukarıya, toplama kümeleme

algoritmaları ve yukarıdan a a ıya bölünür kümeleme algoritmaları olarak iki grupta toplanabilir.

Topla ım kümeleme algoritmaları, ba langıçta veritabanındaki her bir noktayı bir küme olarak görür. Bu kümeleri birle tire birle tire birbirinden ayrı kümeler olu turur.

Bölünür kümeleme algoritmaları ise ba langıçta veritabanındaki tüm noktaları tek bir kümeymi gibi görür. Veritabanını taradıkça, birbirine benzemeyen noktaları kümeden dı arı atarak önceden verilmi , k kadar kümeye da ıtır. Hiyerar ik kümeleme algoritmaları benzerlik ve mesafe ölçütlerini kullandıkları için kullanılması kolay, hemen hemen her türlü veri türüne uygun ve esnek algoritmalarıdır; ancak özellikle bölünür algoritmalar için k küme sayısının verilmesi bir dezavantajdır.

Ba ka bir dezavantaj ise bu kategorideki algoritmalar bir kümeyi olu turduktan sonra, yapıları gere i olu turulan bu kümeyi bir daha kontrol etmezler. Bu algoritmalar yukarıda sözü edilen ba lantı ölçümlerini, kümeler arası mesafe ölçümü kullanırlar. Birçok algoritma  $N \times N$  bir mesafe ya da benzerlik matrisi çıkartarak, kümelemeyi bu matrise dayanarak yapar. Bu matris, her bir elemanın di er elemanlarla olan benzerli ini veya aralarındaki mesafeyi gösteren matristir.

Kümeleme analizinde, algoritmaların zaman ve yer karma ıklı ını en çok artıran unsur da bu mesafe/benzerlik matrisidir. Kümeleme analizi için geli tirilen algoritmalar u ya da bu ekilde söz konusu mesafe/benzerlik matrisini, özellikle yer karma ıklı ını azaltmak için bellekte tutmadan ya da kısaltarak tutarak kümeleme i lemlerini yerine getirirler. Örne in, belirli bir e ik de erinin altındaki benzerliklerin atlanması, kullanılan bilgisayarın belle ini alan karma ıklı ı açısından rahatlatacaktır [Olson, 1993].

### 2.1.3.1 SLINK Algoritması ve Tek Ba lantı Tekni i

Slink algoritması tek ba lantı ya da en yakın kom u tekni ini kullanır [Sibson, 1973]. Bu teknikte, kümeler arasındaki mesafe ölçülürken, iki küme içinde birbirine en yakın iki elemanın uzaklığı ya da ba ka bir deyi le iki kümeyi en yakın kılan elemanların mesafesi kümeler arası mesafe olarak kabul edilir. Algoritmanın zaman karma ıklı  $O(n^2)$ 'dir.

Öncelikle eldeki verilerin, mesafe/benzerlik matrisi çıkartılır; bu matrisi bir a aç haline dönü türür. ebeke modellerinden en küçük maliyetli a aç çıkartılarak, verilen e ik de erine göre kümeler olu turulur.

### 2.1.3.2 CURE Algoritması

Kümele i lemi yapılırken, olu turulan kümelerin kalitesini en çok etkileyen faktör, ana veri toplulu u içinde di er verilerden uzakta bulunan ve sayıları az olup aslında hiç bir kümeyle ait olmaması gereken uç verilerdir. CURE (Clustering Using Reprisentatives- Temsilciler Kullanarak Kümeleme) algoritması bu uç verilerin olu turulan kümelerin kalitesini etkilememesi dü üncesiyle 1998 yılında geli tirilmi bir algoritmadır. Küresel bir geometrik ekil ta ımayan veri gruplarının kümelene mesi için oldukça elveri li bir algoritmadır.

CURE algoritması öncelikle her girdiyi sanki ayrı bir kümeymi gibi ele alır ve her adımda bu küme temsilcilerinin birbirlerine olan yakınlıklarına göre ya birle tirir ya da ayrı kümeler olarak tutar. Öncelikle her bir küme için c adet iyi da ıtılmı temsilci nokta seçilir. Seçilen bu noktalar kümelerin fiziksel eklini geometrik özelli ini ortaya koyar. Daha sonra bu da ıtılmı noktalar bir a katsayısıyla kümenin ortasına, merkezine do ru kaydırılır. Da ıtılmı olan noktalar bu kaydırma i leminden sonra artık o kümenin temsilcileri olarak kabul edilirler. Bundan sonra iki küme arasındaki uzaklık, her biri bir kümeyle ait olan en yakın temsilci çifti arasındaki uzaklıktır.

Temsilcilerin bir  $\alpha$  katsayısıyla kümenin merkezine kaydırılması kümedeki yüzey anomalilerini tolere etti  $\epsilon$  gibi uç verilerin etkisini de azaltır. çünkü uç veriler tipik bir  $\epsilon$  ekilde merkezden uzakta yer alırlar ve sonuç olarak da bu veriler merkeze do ru daha fazla hareket etmi olacaktırlar. Bu uç verilerin uzun mesafeli hareketleri farklı iki kümenin birle tirilmesini önleyecektir [Guha, 1998]. Kullanılan  $\alpha$  katsayısı aynı zamanda, olu an kümelerin  $\epsilon$  eklini belirlemede de kullanılabilir.  $\epsilon$ 'nın alaca ı de er 0-1 arasındadır. Küçük de erli  $\epsilon$  da ılımı noktaların çok az yer de i tirmesine neden olurken kümelerin de ekilsel olarak uzunla masına yol açar.  $\epsilon$  de erinin büyük olması ise da ılımı noktaları küme merkezine oldukça yakla tıraca ı için daha toplu halde kümeler olu acaktır.  $\epsilon = 1$  durumunda ise CURE algoritması merkezi temelli algoritmalara yakla acaktır.

CURE algoritmasının en kötü durumdaki zaman karma ıklı ı  $O(n^2 \log n)$ 'dir. Bununla beraber, Guha 1997'de hazırlamı oldu u bir teknik raporda, veri noktalarının boyutu küçük oldu unda bu zaman karma ıklı ının  $O(n^2)$  oldu unu söyler [Guha, 1997]. Algoritmanın alan/yer karma ıklı ı ise  $O(n)$ 'dir.

Bunun dı nda, algoritmanın bellekte daha az yer i gal etmesini sa lamak için tüm veriler üzerine algoritma ko turulmadan önce, ana kümeden belirli bir miktarda örnek alınarak CURE algoritması bu örnek küme üzerinde uygulanır. Rastgele yapılan bu örnekleme olu turulacak kümelerin kalitesini arttırmaktadır.

### 2.1.3.3 CHAMELEON Algoritması

Chameleon algoritması ilk olarak 1999 yılında Karypis ve arkadaşları tarafından geli tirilmi bir algoritmadır [Karypis, 1999]. Chameleon algoritması, iki küme arasındaki benzerli i dinamik bir model kullanarak belirler. Di er algoritmalardan farklı olarak iki alt kümenin birbirine olan benzerli i ve yakınlı ı bu iki kümeden her birinin kendi iç benzerlikleri ve yakınlıkları ile kıyaslanarak belirlenir ve bu kar ıla tırma sonucunda bu iki alt küme birbirine yakınsa birle tirilir. Bu sayede daha kaliteli ve homojen kümeler yaratılmı olunur. Benzerlik / mesafe matrisinin olu turulabildi i tüm veri türleri ve veri kümeleri için uygulanabilecek bir algoritmadır.

Daha önce incelenen SLINK, CURE ve ileride incelenecek olan DBSCAN ve K-means vs. gibi algoritmalar bazı veri kümelerinin tespit edilmesinde yanılabilirler.

### 2.1.3.4 BIRCH

BIRCH, çok büyük veritabanlarının kümelenebilmesi için geliştirilmiş bir algoritmadır. Ayrıca gürültülü verilerin kontrol edilmesi için bu alanda öne sürülen ilk algoritmadır [Zhang, 1996].

Çok büyük veritabanlarının kümelere ayrılması için tutulması gereken  $N \times N$  matrisini bilgisayar belleği açısından maliyeti,  $ba \cdot ka$  bir deyimle alan karmaşıklığı çok yüksek olacaktır. BIRCH algoritması bilgisayar belleğinde daha az yer kaplayan bir teknikle sahip hiyerarşik yapıda bir kümeleme algoritmasıdır. Temel olarak, kümelemenin yapılabilmesi için bir ağaç oluşturulur ve gerekli tüm bilgilere haiz bu ağaç taranarak kümeleme işlemleri gerçekleştirilir. Kümeleme, her bir düümde mesafe ölçümleri için gerekli bilgilerin tutulduğu ağaç üzerinde gerçekleştirilir.

Kümeleme özellikleri diye adlandırılan ve kümeler hakkındaki bazı bilgileri içeren “bir ağaçtan” yararlanılması BIRCH algoritmasının en tipik özelliğidir. Ancak, sadece sayısal veriler üzerinde kullanılabilir.

Yukarıda sözü edilen bu ağaç CF ağacı olarak adlandırılır ve bu ağaç üç bilgiyi tutar ve algoritmanın kodlanması esnasında bu bilgileri sürekli olarak günceller. Bu bilgiler:

N: Kümede bulunan nokta sayısı

LS: Kümedeki noktaların değerlerinin toplamı

SS: Kümedeki noktaların değerlerinin karelerinin toplamıdır.

CF a acı her bir dü ümün alabilece i dü üm sayısı belirli olan dengeli bir a aştır. Her bir dü üm kendisine ba lı olan alt dü ümlerle ilgili CF de erlerini (N, LS, SS) tutar. A acın her bir yaprak dü ümü ise bir kümeyi temsil eder. Dendrogramda kullanılan tekni in tam tersi olarak CF a acı yukarıdan a a ı do ru çalı ır yani toplama algoritması de il, hiyerar ik fakat bölünür bir kümeleme algoritmasıdır. A aca yeni noktalar eklendikçe CF a acı yaratılmı olur; her bir nokta kendisine en yakın olan yapra a ba lanır ve yaprakların büyüklü ü (T e ik de eri) daha önceden algoritmaya verilmelidir. Noktalar eklene eklene büyüyen yaprak T e ik de erini a arsa, a aça dengeleme veya bölme i lemi yapılır. Buradaki T e ik de eri aslında yapra ın çapıdır. Büyük T de erleri küçük a aaların olu masına neden olurken, T de eri azaldıkça a acın büyüklü ü de artacaktır. Büyük a aalar bellekte daha fazla yer tutaca ı için T de eriyle oynanarak bilgisayarın belle indeki yer karma ıklı ı ayarlanabilir. Bu i lemler yapılırken, veritabanının birden fazla defa okunup taranmasına gerek yoktur; dolayısıyla algoritmanın zaman karma ıklı ı  $O(n)$ 'dir.

CF a acının olu turulması;

Veritabanımız  $D = \{t_1, t_2, \dots, t_i\}$ ,  $i=1, 2, \dots, n$  olsun

$t_i$  elemanının eklenece i yapra ı bul

E ik de eri a ılmamı sa, bu yapra a (küme)  $t_i$  ekle

CF özelliklerini yeniden hesapla

aksi takdirde;

yeterli yer varsa  $t_i$  'yi ayrı bir yaprak olarak ata

yoksa

yapra ı ikiye böl ve  $t_i$  'yi uygun olana ekle.

CF a acının olu turulmasından sonra tam kümeleme i lemine geçilir. CF a acında bazı kısıtlar olu aca ından elde edilen kümeler do al kümeler olmayabilir. Bu nedenle eldeki kümeler (tercihen merkezci yakla ima uygun) ba ka bir algoritmayla birle tirilir. Gerekirse bu i lemden sonra ikinci bir kümeleme daha yapılabilir.

### **2.1.4 Bölümlemeli Yöntemler**

Bölümlemeli yöntemlerde  $n$  adet nokta önceden verilen  $k$  küme sayısına ( $k < n$ ) göre kümelere ayrılır. Hiyerar ik yöntemlerin tersine kullanıcı tarafından verilen bazı kriterlere uygun kümeler yaratılırken, yaratılacak küme sayısı önceden belirlidir. Kullanıcı algoritmaya kümeler arasındaki minimum/maksimum mesafeyi ve kümelerin iç benzerlik kriterlerini de vermek zorundadır [Giudici, 2004].

Bölümlemeli algoritmalar genel olarak hiyerar ik algoritmalarından daha hızlı çalı ırlar; çünkü hiyerar ik algoritmalarındaki gibi bir benzerlik/mesafe matrisi kullanmak zorunda de illerdir. Bundan dolayı da büyük veritabanlarının kümeleneğinde hiyerar ik yöntemlere göre daha uygundur. Bununla beraber önceden verilen kritere uygun birden fazla sonuç çıkarmak mümkün olabilir. Bu durumda algoritmanın gerçekten en uygun çözümü bulup bulamadı ı ise hiç bir zaman bilinmeyecektir [Dunham, 2003]. Bunun ö renilebilmesi için verilerin da ıtılarak, sıra ve yerleri de i tirilerek, algoritmanın tekrar ko turulması gerekecek ve çıkan sonuçların birbirleriyle kıyaslanması gerekecektir. Bu da zaman maliyetini oldukça artıracaktır.

#### **2.1.4.1 K- Ortalama (K-Means) Algoritması**

K-Ortalama algoritması sürekli olarak kümelerin yenilendi i ve en uygun çözüme ula ana kadar devam eden döngüsel bir algoritmadır. Bölümlemeli algoritmaların tipik özelliklerini ta ır. Bu alandaki benzer algoritmaların ço u ya K-Ortalama algoritmasından esinlenerek ya da bu algoritmanın geli tirilmesiyle ortaya çıkmı tır. Dolayısıyla bu algoritmanın anla ılması bundaki sonraki algoritmaların mantı mının kavranmasında önemli bir rol oynayacaktır [Han Jiawei ve Kamber Micheline, 2001].



İlk olarak 1967 yılında ortaya atılan [MacQueen, 1967] K-Ortalama algoritması eldeki verileri  $k$  adet kümede ve kümelerin ortalamalarına göre kümelere ayırır.  $k$  küme sayısı kullanıcı tarafından verilir. Burada kastedilen ortalama daha önce belirtilen küme merkezidir.

Girdiler:

$D = \{t_1, t_2, \dots, t_n\}$  // eldeki veritabanı

$K$  // verilen küme sayısı

Algoritma:

Keyfi olarak  $m_1, m_2, \dots, m_k$  ortalama belirle.

her bir  $t_i$ 'yi en yakın oldu  $u$   $m_i$ 'nin kümesine ata.

Kümelere ait  $m_1, m_2, \dots, m_k$  de erlerini yeniden hesapla.

Küme elemanlarında herhangi bir de  $i$  iklik yoksa dur.

ilk adımdan itibaren tekrar et.

Çıktı:

$K$  adet küme

### 2.1.4.2 PAM Algoritması

PAM algoritması  $k$  adet kümeyi bulmak için seçilen temsilcilerin etrafına ana kümedeki tüm elemanları toplayarak ve her defasında bu temsilcileri de  $i$  tirerek kümeleme  $i$  lemini tamamlar.

PAM algoritmasının temsilci olarak seçti  $i$  noktaya medoid denilir; dolayısıyla bu algoritma  $k$ -medoid algoritması olarak da anılır. Bu temsilci (medoid) seçiminden kasıt ise kümenin merkezine yakın mesafede bulunan noktanın belirlenmesidir.  $K$  adet küme için seçilen  $k$  adet temsilci belirlendikten sonra, veritabanındaki temsilci olmayan  $d_i$  er noktalar (veriler) kendilerine en çok benzeyen temsilcinin etrafında toplanır. Daha matematiksel bir ifadeyle,  $e$  er  $t_i$  bir temsilci ve  $t_j$  ise temsilci olmayan bir  $b_a$  ka nokta olsun,  $e$  er  $d(t_i, t_j) = \min_{t_e} d(t_j, t_e)$  ise  $t_j$ ,  $t_i$  tarafından temsil edilen kümeye aittir. Burada  $\min_{t_e}$  tüm temsilciler içindeki en küçü ü ifade ederken,  $d(t_a, t_b)$  ise  $t_a$  ve  $t_b$  noktaları arasındaki mesafe veya benze mezli i ifade etmektedir. Bu durumda bir kümenin kalitesi o kümedeki temsilciyle  $d_i$  er noktalar arasındaki ortalama mesafe ya da ortalama benzememe de eriyle ölçülebilir. PAM algoritmasında tüm benzerlik / mesafe ölçütleri kullanılabilir [Raymond ve Han, 1994].

### 2.1.4.3 CLARA Algoritması

CLARA algoritması bütün veritabanının tarayarak temsilci noktalar seçmek yerine, veritabanından rastgele bir örnek kümeyi alarak, PAM algoritmasını bu örnek küme üzerine uygular. Bu uygulama sonucunda olu acak olan kümelerin her birinin temsilcisi belirlenir. Daha sonra ana kümeyi olu turan veritabanından bir örnek küme daha seçilir. Bu esnada ilk temsilcilerin rastgele seçilmesi yerine bir önceki a amada belirlenmi temsilciler kullanılır. Bu da algoritma içinde temsilci de i imini azaltacak ve algoritma hem daha hızlı bir ekilde i leyecek hem de daha kaliteli sonuçlar verecektir. Bu tekrar örnekleme i leminin 5 defa yinelenmesi ve her defasında  $40 + 2k$  adet örnek seçilmesinin en iyi sonucu verdi i Kaufman ve Rousseeuw (1990) tarafından rapor edilmi tir.

PAM algoritmasıyla kıyaslandığında CLARA algoritmasının daha geni veri tabanlarında güvenli bir şekilde çalışabildiği belirlenmiştir; 10 kümede 1000 eleman gibi.

#### **2.1.4.4 CLARANS Algoritması**

CLARANS algoritması verilen  $n$  adet nesnenin temsilciler aracılığıyla ve bir ebeke diyagramından yararlanılarak  $k$  adet kümeye ayrılması şeklinde özetlenebilir.

$G_{n,k}$  ile temsil edilen bu ebeke diyagramında her bir düğüm  $\{O_1, O_2, \dots, O_k\}$ 'den oluşan  $k$  adet nesneyi temsil eder.  $O_1, O_2, \dots, O_k$  bir bakıma temsilcilerdir. ki düğüm birbirinden sadece bir nesnenin de i imiyle ayrılıyorsa bu iki düğüm komu olarak kabul edilir.

CLARANS algoritmasının iki parametresi vardır: maks-komu (maxneighbor) ve yerel-miktarı (numlocal). Maks-komu parametresi incelenecek komu sayısının üst limitini ifade ederken, yerel-miktarı ise elde edilecek yerel minimum nokta sayısının alt sınırını ifade etmektedir.

### 3. BÖLÜM

#### 3.1 Çalışmanın Konusu

Bu tez çalışmasında, eşitlik ayırıcı (equal-split) olarak adlandırılan parametre kullanılarak kategorisel verilere sahip bir veri tabanı üzerinde kümeleme algoritması gerçekleştirilmiştir.

Temel olarak, veri tabanındaki herhangi bir alan içindeki veri sayısı ve de i ken adedinden yararlanmaktadır. Bu parametre sayesinde hangi alan ya da alanların veri tabanını mümkün olan en az parçaya ayırabileceği tespit edilmektedir.

EP: Eşitlik-Ayrıcı Parametre (Equal-Split Parameter).

Herbir de i kenin merkeze uzaklığının mutlak toplamı,

N: Toplam veri sayısı,

NV: Her alan içindeki de i ken sayısı,

NV<sub>i</sub>: Her alan içindeki de i kenden kaç adet olduğu,

CA: De i ken merkezi,

$$CA = \frac{N}{NV}$$

$$EP = \sum_{i=1}^n |CA - NV_i|$$

#### 3.2 Çalışmanın Amacı

Bu çalışmanın ilk amacı, önerilen algoritmanın uygulanabilirliğini denetlemektir. İkinci amaç ise mevcut yöntemlerle arasındaki benzerlik ve farklılıklarını ortaya koymaktır. Karmaşık O(n) olan bu algoritma için aynı tür verilerle işlem yapabilen diğer algoritmalarla karşılaştırma yapılmıştır.

### 3.3 Çalı maya Konu Olan Algoritma

<p>Begin</p> <p>Calculate Area Count</p> <p>mEP:=-1</p> <p>Loop for All Areas</p> <p>  Do Begin</p> <p>    N:= Total Data Count (for current area)</p> <p>    NV:= Total Variable count (for current area)</p> <p>    <math>CA = \frac{N}{NV}</math></p> <p>    <math>EP = \sum_{i=1}^n  CA - NV_i </math></p> <p>    If mEP&gt;EP then mEP:=EP</p> <p>  End;</p> <p>End.</p>	<p>A1) Alan sayısını hesapla. Calculate Area Count</p> <p>A2) Sıradaki alana geç.</p> <p>A3) Toplam veri sayısını hesapla (N) .</p> <p>A4) Toplam de i ken sayısını bul (NV) .</p> <p>A5) Merkezi hesapla (<math>CA = \frac{N}{NV}</math>) .</p> <p>A6) Her bir de i kenin merkeze olan uzaklı ının mutlak toplamını hesapla</p> $EP = \sum_{i=1}^n  CA - NV_i  .$ <p>A7) En küçük de erli de i keni kök olarak seç.</p> <p>A8) De i ken de erlerinden birbirine e it olan varsa herhangi birini seç.</p> <p>A9) lemi tekrarlamak için A2. adıma dön.</p> <p>A10) Tüm veriler aynı oldu unda yapra a ula ılımı demektir.</p> <p>A11) Dur.</p>
---	---

Kümeleme i leminde ; Veriler olu turulan a ca yayılacaktır. Her dü ümdeki veri miktarı, verilerin merkezi, standard sapma, çap, yarıçap özelliklerinin kriter olarak seçilenleri ya da istenilenleri hesaplanır. Örne in kazanç (Gain) ya da Gini indeksi (Gini index) i lemlerinde oldu u gibi. E ik de erleri ba tan verildi inde budama i lemi yapılabilir. Örne in bir sonraki dü üme geçi te veri niceli inde anlamlı bir de i me yoksa dallanmaya gerek yok demektir. Bu durumda o dü üm (node), yaprak (leaf) haline dönü türülebilir. Yapraklarda olu an verilerin merkezi hesaplanıp merkezleri birbirine yakın olan yapraklar birle tirilebilir.

### 3.4 Uygulama

Kullanılan sentetik veri tabanlarından birisinin temsili gösterimi ve algoritmanın uygulanması sonucu elde edilen sonuçlar a a ıda gösterilecektir. Kullanılan veri tabanı 10000 veriden olu maktadır. Verilerde özellikle bizim bulunmasını istedi imiz alanlar ile ilgili de erler verilmi tir. Aynı zamanda veri tabanında mevcut olan alan ve verilerden bazıları yanıltıcı olması açısından özellikle yerle tirilmi tir.

### 3.4.1 Kullanılan Veri Tabanı

Kullanılan bu veri tabanında toplam veri sayısı 10000 'dir. Mevcut verilerin bazıları sayısal olmalarına rağmen kategorisel olarak sınıflandırılmıştır.

Alan Adı	Veri Sayısı (NV)	NV <sub>1</sub>	NV <sub>2</sub>	NV <sub>3</sub>	NV <sub>4</sub>	NV <sub>5</sub>	NV <sub>6</sub>
Marital Status	2	5673	4327				
Gender	2	4887	5113				
Children	6	5799	1323	735	809	755	579
Education	5	2578	1635	1979	2698	1110	
Occupation	5	1423	1730	1399	2953	2495	
Home Owner	2	3221	6779				
Cars	5	1802	2394	3993	989	822	
Commute Distance	5	3081	1630	1810	1520	1959	
Region	3	2928	5456	1616			
BikeBuyer	2	9000	1000				

Tablo - 3.1 Algoritmadan Alınan İlk Değerler

Marital Status	Gender	Children	Education	Occupation	Home Owner	Cars	Commute Distance	Region	Bike Buyer
Single	Male	0	Partial College	Clerical	No	1	0-1 Miles	Europe	Yes
Married	Female	0	Graduate Degree	Clerical	Yes	0	0-1 Miles	Europe	Yes
Single	Female	2	Bachelors	Skilled Manual	No	1	0-1 Miles	North America	Yes
Single	Female	0	High School	Professional	Yes	2	5-10 Miles	Pacific	Yes
Married	Female	1	Bachelors	Clerical	Yes	0	2-5 Miles	Europe	Yes
Single	Female	1	Partial College	Manual	No	0	0-1 Miles	Europe	Yes
Single	Male	0	Partial College	Skilled Manual	No	2	1-2 Miles	North America	Yes
Married	Female	0	Partial College	Clerical	Yes	1	1-2 Miles	North America	Yes
Single	Female	2	Bachelors	Skilled Manual	Yes	1	2-5 Miles	North America	Yes
Married	Female	0	Partial College	Clerical	Yes	1	1-2 Miles	North America	Yes
Single	Male	0	Partial College	Clerical	Yes	1	1-2 Miles	North America	Yes
Single	Female	2	Bachelors	Professional	Yes	2	5-10 Miles	Pacific	Yes
Single	Female	0	Graduate Degree	Professional	Yes	0	2-5 Miles	North America	Yes
Single	Female	0	Partial College	Professional	No	2	2-5 Miles	Europe	Yes
Single	Male	0	Graduate Degree	Clerical	Yes	0	0-1 Miles	Europe	Yes
Single	Female	0	Bachelors	Professional	Yes	1	5-10 Miles	Pacific	Yes
Married	Male	0	Bachelors	Professional	No	1	0-1 Miles	Pacific	Yes
Single	Female	0	Partial College	Skilled Manual	No	1	1-2 Miles	North America	Yes
Single	Male	0	Bachelors	Clerical	Yes	0	0-1 Miles	Europe	Yes
Single	Female	0	Bachelors	Professional	No	1	0-1 Miles	Pacific	Yes
Married	Female	1	Partial College	Professional	Yes	2	1-2 Miles	North America	Yes
Married	Female	0	Bachelors	Management	Yes	2	5-10 Miles	North America	Yes
Married	Male	0	Graduate Degree	Management	Yes	1	5-10 Miles	Pacific	Yes
Married	Male	0	Graduate Degree	Professional	Yes	0	0-1 Miles	North America	Yes
Married	Female	0	Bachelors	Clerical	No	0	0-1 Miles	Pacific	Yes
Single	Female	1	Partial High School	Clerical	Yes	2	5-10 Miles	Pacific	Yes
Single	Male	0	Graduate Degree	Professional	Yes	0	2-5 Miles	North America	Yes
Married	Male	5	Partial High School	Professional	Yes	4	10+ Miles	Pacific	Yes
Married	Female	0	Bachelors	Management	Yes	2	0-1 Miles	Pacific	Yes
Single	Female	0	High School	Skilled Manual	Yes	1	5-10 Miles	North America	Yes

Tablo - 3.2 Veri Tabanı Örneği

### 3.4.2 Delphi Programlama Dili ile Örnek Uygulama – 1

MARITALSTA	GENDER	YEARLYINCO	CHILDREN	EDUCATION	OCCUPATION	HOMEOWNER	CARS	COMMUTEDIS	REGION	AGE	BIKEBUYER
Single	Female	50000	0	Graduate Degree	Skilled Manual	Yes	0	1-2 Miles	North America	35	No
Married	Male	50000	0	Graduate Degree	Skilled Manual	Yes	0	1-2 Miles	North America	35	No
Single	Female	40000	0	Bachelors	Professional	No	1	0-1 Miles	North America	39	No
Single	Male	80000	3	Bachelors	Skilled Manual	Yes	0	2-5 Miles	North America	40	No
Married	Male	60000	0	Graduate Degree	Professional	No	0	0-1 Miles	North America	40	No
Married	Female	40000	0	High School	Skilled Manual	Yes	2	5-10 Miles	North America	29	No
Married	Female	20000	0	Graduate Degree	Clerical	Yes	0	0-1 Miles	Europe	45	No
Married	Female	40000	0	Partial College	Skilled Manual	Yes	1	5-10 Miles	North America	29	No
Single	Female	40000	0	Partial College	Skilled Manual	No	1	1-2 Miles	North America	30	No
Married	Male	50000	0	Partial College	Skilled Manual	Yes	1	5-10 Miles	North America	32	No
Single	Male	50000	0	Partial College	Skilled Manual	Yes	1	5-10 Miles	North America	32	No
Married	Female	50000	0	Partial College	Skilled Manual	Yes	1	5-10 Miles	North America	32	No

N Sayisi: 10000  
Kullanılan Alan Sayısı: 12  
En Küçük NV: 226

Buton1

NV-1	NV-2	NV-3	NV-4	NV-5	NV-6	NV-7	NV-8	NV-9	NV-10	Bütün NV adedleri	EP Listesi
Single	Male	0	Partial College	Clerical	No	1	0-1 Miles	Europe	Yes	Single = 4326	1346
Married	Female	5	Graduate Degree	Professional	Yes	0	10+ Miles	Pacific	No	Married = 5672	8238
		2	Bachelors	Management		3	2-5 Miles	North America		Male = 5112	6233
		1	High School	Skilled Manua		2	5-10 Miles			Female = 4886	2553
		4	Partial High School	Manual		4	1-2 Miles			0 = 5798	2897
		3								5 = 578	3598
										2 = 734	4775
										1 = 1322	2165
										4 = 754	4246
										3 = 888	333333333
										Partial College = 2697	8000
										Graduate Degree = 1634	

ekil - 3.1 Delphi Programı Uygulaması - 1

Algoritma çalı maya ba ladı nda veritabanındaki bütün alanları kullanarak N, NV, NV<sub>i</sub> de erlerini hesaplar. Tablo 3.2 de verilen örnek veritabanı ele alındı nda hesaplanan de erler Tablo 3.1 de verilmi tir. Bu durum öyle açıklanabilir; Tabloda 10000 adet veri bulunmaktadır. Böylece N sayısı ilk ba ta 10000 olarak belirlenmi olur. Tablodaki 1. alan "Marital Status" alanıdır. Bu alan "Married" ve "Single" olmak üzere iki adet veri türü içermektedir. Bu durumda "Marital Status" alanı için NV de eri 2 olarak hesaplanır. Aynı ekilde 2. alan olan "Gender" alanı "Male" ve "Female" olarak iki adet veri türü içermektedir. "Gender" alanı için veri türü sayısı NV de eri de 2 olarak hesaplanır. Bütün alanlardaki NV de erleri bu ekilde hesaplandıktan sonra NV<sub>i</sub> de erlerinin hesaplanmasına geçilir. Algoritma her bir alandaki NV de erlerinden o alanda kaç adet veri oldu unu hesaplayacaktır. Örne in Tablo 3.1 de de görülece i gibi "Education" alanında NV de eri 5'tir. Bu 5 de eri veri türü sayısı olarak "Bachelors", "Graduate Degree", "High School", "Partial College", "Partial High School", bilgilerinden elde edilmi tir. "Education" alanında, "Bachelors" verisinden 2578, "Graduate Degree" verisinden 1635, "High School" verisinden 1979, "Partial College" verisinden 2698, "Partial High School" verisinden 1110 adet bulunmaktadır. Buradaki alanlar sırasıyla NV<sub>1</sub>, NV<sub>2</sub>, NV<sub>3</sub>, NV<sub>4</sub>, NV<sub>5</sub> olarak kabul edilmektedir. Böylece Tablo 3.1'de de görülece i üzere "Education"

alanı için  $NV = 5$ , "Bachelors" olan  $NV_1 = 2578$ , "Graduate Degree" olan  $NV_2 = 1635$ , "High School" olan  $NV_3 = 1979$ , "Partial College" olan  $NV_4 = 2698$ , "Partial High School" olan  $NV_5 = 1110$  olarak hesaplanmıştır. Benzer şekilde tüm alanlar için ayrı ayrı bu veriler rekürsif (özyinelemeli) olarak hesaplanırken aynı zamanda EP değerleri de hesaplanmaktadır. Her bir alan için EP değerinin hesaplanması algoritmada yer aldığı gibi  $EP = \sum_{i=1}^n |CA - NV_i|$  formülünden yararlanılarak gerçekleştirilecektir.

Yukarıda verilen değerleri kullanarak "Education" alanı için EP değerinin hesaplanmasını gösterecek olursak;

$$CA^{\text{Education}} = \frac{N}{NV} = \frac{10000}{5} = 2000$$

$$EP^{\text{Education}} = \sum_{i=1}^n |CA - NV_i|$$

$$EP^{\text{Education}} = |2000 - 2578| + |2000 - 1635| + |2000 - 1979| + |2000 - 2698| + |2000 - 1110|$$

$EP^{\text{Education}} = 2552$  olarak hesaplanmıştır. Bu değer Tablo 3.1'de de gösterilmiştir.

Birinci aşamada sonuçta en küçük EP değerine sahip olan değeri kök olarak seçilir. Tablo 3.4'te de algoritmanın çalışması sırasında, her bir alanın her bir veri türü için aynı şekilde EP değerlerini nasıl hesapladığı örneklenmiştir. Örneğin bir sonraki dallanmada algoritma, "Marital Status" un "Married" ve "Single" alanları için ayrı ayrı aynı işlemleri gerçekleştirecektir. Bu aşamada bu şekilde açıklanabilir;

Tablo 3.1'e bakıldığında "Marital Status" alanının "Married" durumu için toplam veri sayısı 5673, "Single" için 4327 olarak hesaplanmıştır. Söz konusu olan "Married" değeri olduğunda göz önüne alınacak durumlar sadece "Marital Status" alanının "Married" değeri alındığında diğer alanlarda mevcut olan toplam veri sayısı ve verilerdir.  $N$ ,  $NV$ ,  $NV_i$  değerleri buna göre değerlendirilecektir. Bu durumda  $N=5673$  olarak değerlendirilecektir. Mesela bu aşamada "Children" alanı için yapılacak işlemleri gösterecek olursak;



"Marital Status" = "Married" iken

$$N = 5673$$

$$NV = 6$$

$$CA = \frac{5673}{6} = 945.5 \cong 945$$

$$\text{"Children" = "0" iken } NV_1 = 2948$$

$$\text{"Children" = "1" iken } NV_2 = 916$$

$$\text{"Children" = "2" iken } NV_3 = 339$$

$$\text{"Children" = "3" iken } NV_4 = 512$$

$$\text{"Children" = "4" iken } NV_5 = 542$$

$$\text{"Children" = "5" iken } NV_6 = 416$$

$$EP = |945 - 2948| + |945 - 916| + |945 - 339| + |945 - 512| + \\ |945 - 542| + |945 - 416|$$

$$EP = 4005 \text{ olur.}$$

"Marital Status" = "Single" iken

$$N = 4327$$

$$NV = 6$$

$$CA = \frac{4327}{6} = 721.16 \cong 721$$

$$\text{"Children" = "0" iken } NV_1 = 2851$$

$$\text{"Children" = "1" iken } NV_2 = 407$$

$$\text{"Children" = "2" iken } NV_3 = 396$$

$$\text{"Children" = "3" iken } NV_4 = 297$$

$$\text{"Children" = "4" iken } NV_5 = 213$$

$$\text{"Children" = "5" iken } NV_6 = 163$$

$$EP = |721 - 2851| + |721 - 407| + |721 - 396| + |721 - 297| + \\ |721 - 213| + |721 - 163|$$

$$EP = 4260 \text{ olur.}$$

Her a amada bir sonraki a amaya geçmeden en küçük de erli EP seçilmi olur. Böylece seçilen alan bir sonraki dü üm olacaktır. Bütün EP de erleri bu ekilde dallanma i lemi sayesinde hesaplanmı olur. Burada amaç veri tabanındaki alanları kullanarak veri tabanını en az sayıda fakat e it bölen alanları ortaya çıkarmaktır. E itlikten kasıt veri sayısı gözönüne alındı nda belirli bir alana yo unla manın kökte de il yapraklara do ru olmasıdır. Aksi halde veriler rastgele bölünmü olacaktır. Az sayıda olmasından kasıt ise mümkün oldukça iki parçaya ayırabilmektir. Elbette bu alanlardaki veri türü sayısı ile de orantılıdır. Böylece kök,

dü ümler ve yapraklar belirlenmi olmaktadır. Dikkat edilmesi gereken nokta, e er veri tabanında aynı tür veriden sadece bir adet bulunuyorsa ya da bir alanda kayıt sayısı kadar veri türü mevcut ise  $\sum_{i=1}^n |CA - NV_i|$  formülünden elde edilecek sonucun sıfır ( 0 ) olmasıdır. Formül mutlak de erler toplamına, bir bakıma e it uzaklık mesafesi hesaplamaya, dayandı ı için asla negatif de er üretmeyece inden elde edilecek en küçük de er sıfır oldu unda bu de er kök olarak seçilecektir. Bu durumda veritabanı e it parçalara ayılamayacaktır. E er bu tür verilere sahip alanlar var ise kök de eri olarak onlar seçilece inden algoritma do ru çalı mayacaktır. Bu duruma engel olmak için iki seçenek mevcuttur. Birincisi bu algoritmaya uygun verilere sahip veritabanlarının kullanılmasıdır. kinci yol ise algoritma programlanırken, EP de erlerinin seçim a masında iki probleme ait kriterleri kontrol etmektir. Bu i lem u ekilde gerçekleştirilebilir. EP de erlerinin en küçü ü seçilmeden en küçük olarak öngörülen EP de erine ait alandaki veri türü sayısı kontrol edilebilir. EP'nin ait oldu u veri türü sayısı 1 (NV=1) ise ya da veri türü sayısı veri sayısı ile aynı (NV=N) ise o EP de eri pas geçilip bir sonraki küçük de er en küçük olarak seçilebilir.

### 3.4.3 Excel 'de Elde Edilen Sonuçların ncelenmesi

	NV		EP	NV-1	NV-2	NV-3	NV-4	NV-5	NV-6
Marital Status	2	2	1346	673	673	*	*	*	*
Gender	2	1	226	113	113	*	*	*	*
Children	6	10	8265	4132	343.7	931.7	857.7	911.7	1088
Education	5	4	2552	578	365	21	698	890	*
Occupation	5	5	2896	577	270	601	953	495	*
Home Owner	2	6	3558	1779	1779	*	*	*	*
Cars	5	8	4774	198	394	1993	1011	1178	*
Commute Distance	5	3	2162	1081	370	190	480	41	*
Region	3	7	4245	405.3	2123	1717	*	*	*
BikeBuyer	2	9	8000	4000	4000	*	*	*	*

Tablo - 3.3 Algoritmanın Birinci A masından Örnek Veriler

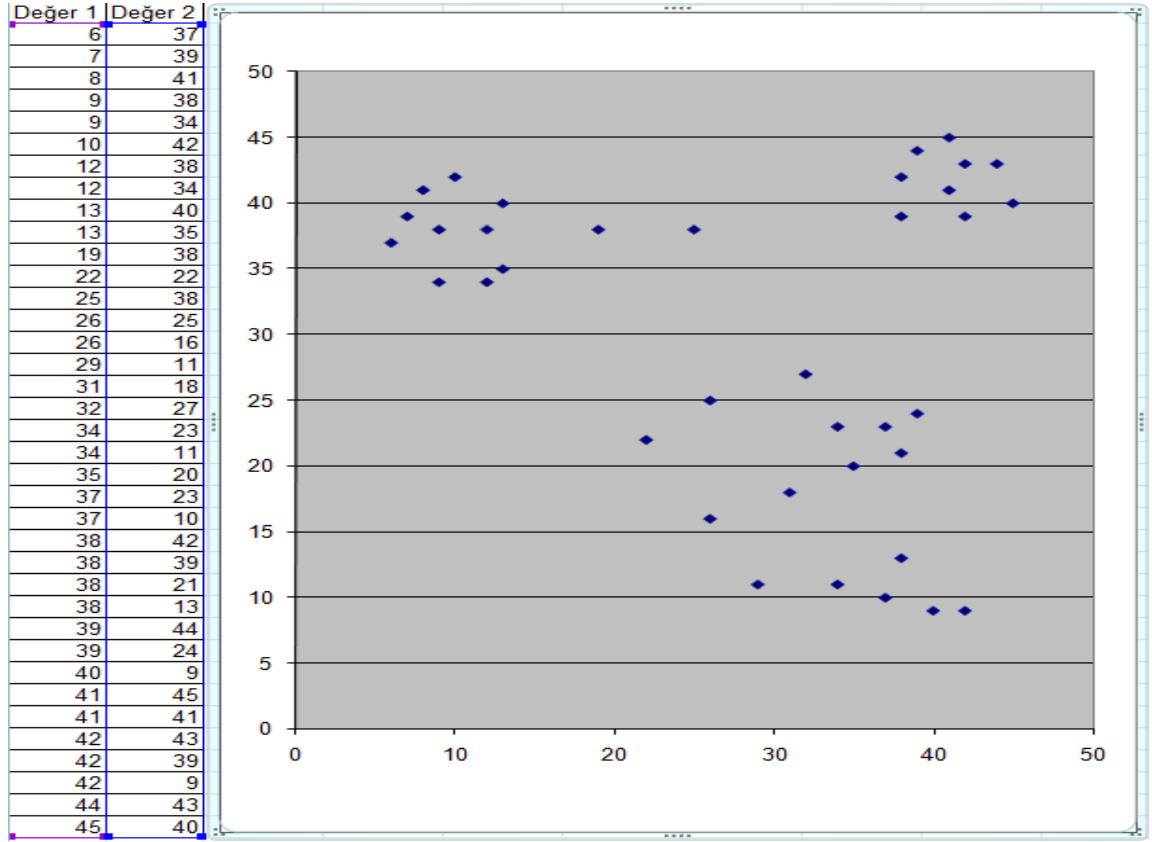
<b>Marital Status</b>	5673	EP	NV-1	NV-2	NV-3	NV-4	NV-5	NV-6
<b>Married</b>	<b>Gender</b>	367	183.5	183.5	*	*	*	*
	<b>Children</b>	4005	2003	29.5	606.5	433.5	403.5	529.5
	<b>Education</b>	1587.6	444.4	65.6	41.6	349.4	686.6	*
	<b>Occupation</b>	2176.4	474.6	19.4	613.6	722.4	346.4	*
	<b>Home Owner</b>	3487	1744	1744	*	*	*	*
	<b>Cars</b>	2130.4	26.4	17.4	1021	467.6	597.6	*
	<b>Commute Distance</b>	1230.8	615.4	105.6	228.6	266.6	14.6	*
	<b>Region</b>	2796	503	1398	895	*	*	*
	<b>BikeBuyer</b>	4673	2337	2337	*	*	*	*
<b>Marital Status</b>	4327	EP	NV-1	NV-2	NV-3	NV-4	NV-5	NV-6
<b>Single</b>	<b>Gender</b>	141	70.5	70.5	*	*	*	*
	<b>Children</b>	4259.67	2130	314.2	325.2	424.2	508.2	558.2
	<b>Education</b>	1005.6	133.6	299.4	20.6	348.6	203.4	*
	<b>Occupation</b>	783.6	102.4	289.4	12.6	230.6	148.6	*
	<b>Home Owner</b>	71	35.5	35.5	*	*	*	*
	<b>Cars</b>	2696.4	224.4	376.6	971.6	543.4	580.4	*
	<b>Commute Distance</b>	1008.4	465.6	264.4	38.6	213.4	26.4	*
	<b>Region</b>	1644.67	97.67	724.7	822.3	*	*	*
	<b>BikeBuyer</b>	3327	1664	1664	*	*	*	*

Tablo - 3.4 Algoritmanın ikinci A aşamasından Örnek Veriler

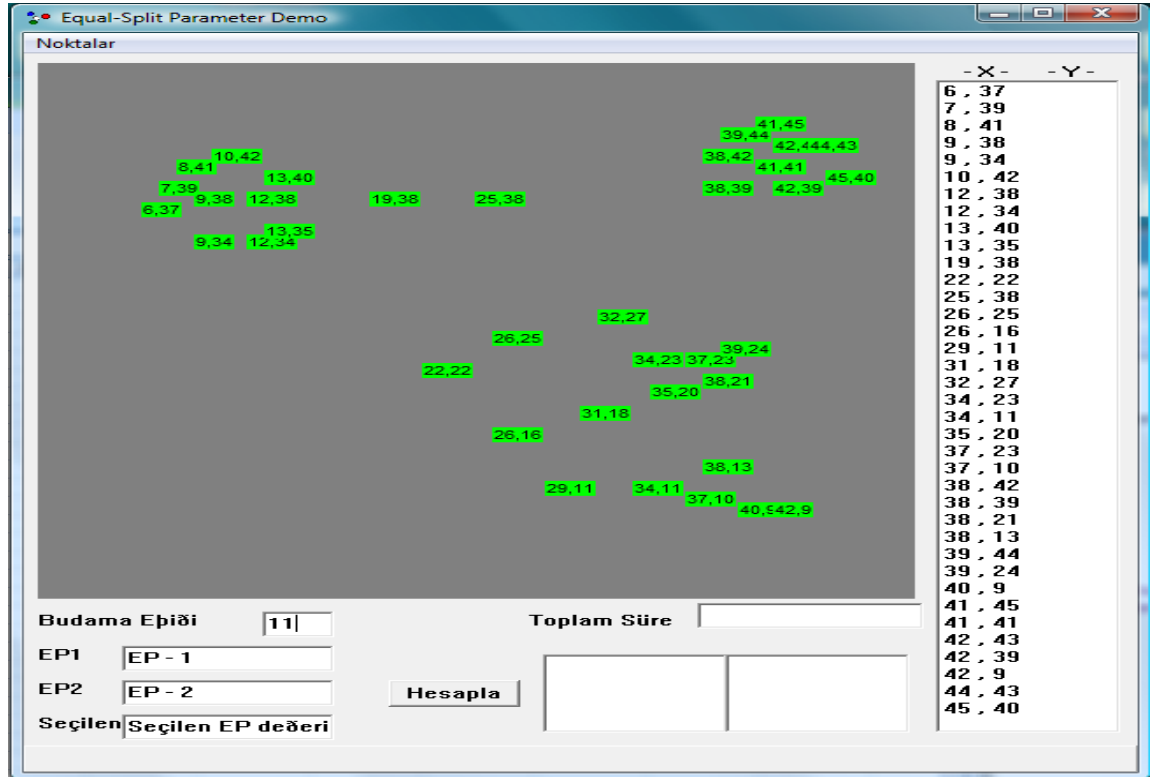
### 3.4.4 Algoritmanın ve Uygulamanın Sonuçları

Bu bölümde çalışması yapılan algoritmanın uygulama sonuçlarını de erlendirmek için ba ka bir örnek üzerinde elde edilen bilgiler de erlendirilecektir.

Delphi programlama dilinde gerçekleştirilen bu çalışmada, a a ıdaki tabloda yer alan verilerin kümelemesi yapılmaya çalışılmıştır.



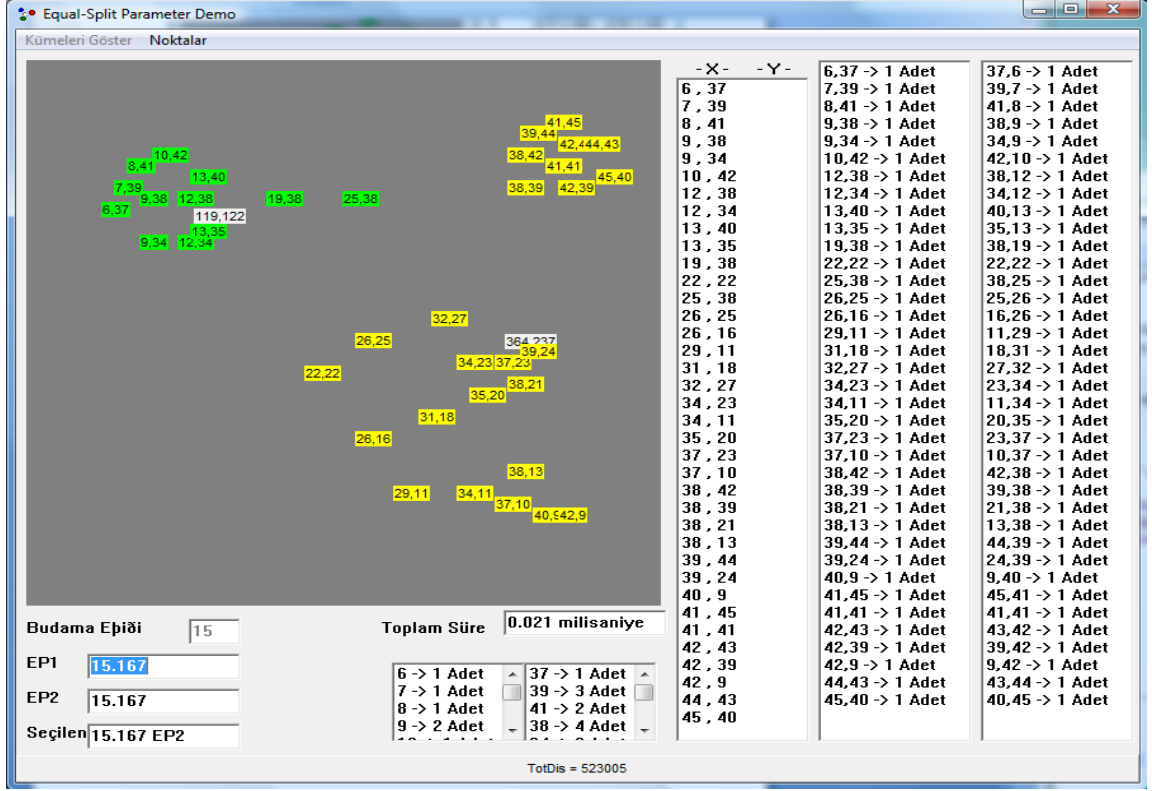
Tablo 3.5 De erlendirmede Kullanılan Veriler



ekil 3.2 Delphi Programı Uygulaması - 2

ekil 3.2 de Tablo 3.5 de verilerin delphi programında bir düzlemde gösterilmesinin gerçekleştiği görülmektedir. Program ilk çalıştırmada veriler tablo olarak görülmemektedir. Programın menüsündeki "Noktalar" seçeneği tıklandığında sağ taraftaki önceden belirlenmiş noktalar grafiksel olarak gösterilecektir. Programın temel amacı doğru bir biçimde ve mümkün olan en kısa sürede verilen noktaları kümelere ayırmaktır. "Noktalar" seçeneği tıklanarak veriler grafiksel olarak gösterildikten sonra bu değerleri kullanarak algoritma işletilecektir. Bunun için aşağıda bulunan "Hesapla" butonu kullanılacaktır. Bir önceki uygulamada anlatıldığı gibi verilerin kümelenebilmesi için aç hesaplamaları yapılacak ve ardından kümeleme işlemine geçilecektir. Kümeleme işleminin gerçekleşmesi için "Budama Eşiği" kullanılması gerekmektedir. Bunun sebebi kullanılan algoritmanın çok detaylı bir aç üretmesidir. Doğal olarak bu açın budama işlemine tabi tutulması gerekecektir. Bunun için küme merkezi, veri miktarı, standard sapma gibi özellikler kullanılabilir. Bu uygulamada budama için küme merkezi ve çap değerleri kullanılmıştır. Bu işlemlerin ne kadar süre aldığı hesaplanmış ve ekrana yazdırılmıştır. "Hesapla" butonuna tıklandığında Tablo 3.5 'teki Değer1 ve Değer2 alanında yer alan veriler için EP değerleri hesaplanır. Hesaplama sonucunda her veriden kaç adet değer olduğu da ekranda gösterilecektir. Diğer EP hesaplamaları için aynı işlemler tekrar edilip aç oluşturulacaktır. "Budama Eşiği" alanındaki veriye göre birbirine yakın noktalar aynı kümelere alınarak aç budanacak ve küme sayısı belirlenecektir. Bu işlemden sonra oluşan kümelerin gösterilmesi için menüden "Kümelere Göster" butonu tıklandığında kümeler gösterilecektir. Burada dikkat edilmesi gereken durum "Budama Eşiği" değerleridir. düşük budama değerleri küme sayısını arttırırken yüksek budama değerleri küme sayısını azaltacaktır. Uygulamada yaklaşık 3 adet budama aralığı tespit edilmiştir. "0 - 7" arasındaki değerler için 6 ya da daha fazla küme, "8 - 12" arasındaki değerler için 4 ya da daha az küme, "13 - 25" arasındaki değerler için 2 ya da daha az küme, olduğu görülmüştür. Bu durumlar aşağıdaki şekillerde gösterilmiştir. Bu örnek verilerle programın çalışması, verilen budama değerine göre, ortalama 22 milisaniyedir. Bu süreye ekrana çizim için geçen süre dahil edilmemiştir.





ekil 3.5 Budama E i inin 12'den Büyük Oldu u Durum

### 3.5 Çalı manın Önemi

Bu çalı ma geli mekte olan veri madencili i uygulamalarında yeni bir algoritma kullanarak farklı bir bakı açısı elde etmeyi hedeflemesiyle önem ta imaktadır. Olu turulan algoritma daha da geli tirilip di er algoritmalarla da birle tirilerek veri madencili i uygulamalarında etkin sonuçlara ula maya katkıda bulunacaktır.

### 3.6 Çalı manın Kısıtları

Bu tez çalı masının en önemli kısıtlarından birisi, kullanılan algoritmanın yeni olması ve bu sebepten belirli kriterlere uyan veri tabanlarında çalı abilmesidir. Seçilen veri tabanlarında bazı özel durumlara dikkat edilmesi gerekmektedir. Bu durumlardan birisi herhangi bir alanda var olan veri çe idinin her bir kayıt için

birbirinden farklı olmaması gereklili idir. Bir di er olumsuz durum ise herhangi bir alanda sadece bir adet veri çe idi bulunmasıdır ki her iki durumda da algoritma hatalı sonuçlar üretmektedir.

### 3.7 Ara tırmanın Modeli ve Hipotez

Bu tezde öne sürülen algoritma için geni bir literatür taraması yapılmı tır. Elde edilen ara tırma sonuçları bu çalı mada temel alınarak sonuç elde edilmeye çalı ılmı tır.

**Hipotez:** Kullan algoritma veri madencili i uygulamalarında kullanılabilecek ve mevcut algoritmalara yakın derecede etkili olacaktır.

### 3.8 Algoritmaya Ait Özellikler

Kullanılan algoritmanın veri setini yeniden taramasına gerek olmamasından dolayı karma ıklı ı  $O(n)$  'dir. Bu özelli i sayesinde daha fazla karma ıklık de erine sahip olan algoritmalarından daha hızlı sonuç vermektedir. Kar ıla tırmak gerekirse, BIRCH algoritmasında bir CF (Clustering-Feature) a acı yaratılır. CF a acı dallanma faktörü ve budama parametresi ile yükseklik dengeli bir a aç olur. BIRCH algoritmasında en kısa yoldan sınıflara ula mak için en uygun simge kök olarak seçilmi tir. Böylece algoritma veri tabanındaki en uygun alan ile bir sonraki dü üümü (node) belirlemektedir. Burada "en uygun" ifadesi ile, veri tabanını kabaca iki ya da daha fazla e parçaya bölmek kastedilmektedir. Buna en yakın alan en uygun alan olarak seçilir. Literatürde a acı olu turmak için farklı yöntemler vardır. Bunlardan biri entropi (entropy) kavramıdır. ID3 ve C4.5 algoritmaları dü üüm temsilcilerini bulmak için entropi yöntemini kullanırlar.

Kullanılan algoritma her dal ve ayrı olarak onun dü üümleri için EP'yi yeniden hesaplayan tekrarlamalı bir algoritmadır. Algoritma prosedürümüz do al kümeler olu turdu u için a aç budanmaya ihtiyaç duymaktadır. BIRCH algoritmasındaki gibi e ik de erlerini kullanarak budama yapılmakta. Yapraklardan ba layarak, kullanıcı tarafından verilen e ik de erlerinin altında kalan her yaprak en yakınındakine dahil oluyor. Böylece algoritma tarafından üretilen do al kümelerin sayısı azalıyor. Kümeler biçimlendirildikten sonra PAM, CLARA ya da K-Means gibi kümeleme algoritmaları kullanılarak yaprakların sayısı istenilen sayıya dü üürülebilir.



### 3.9 Hipotezin Testi

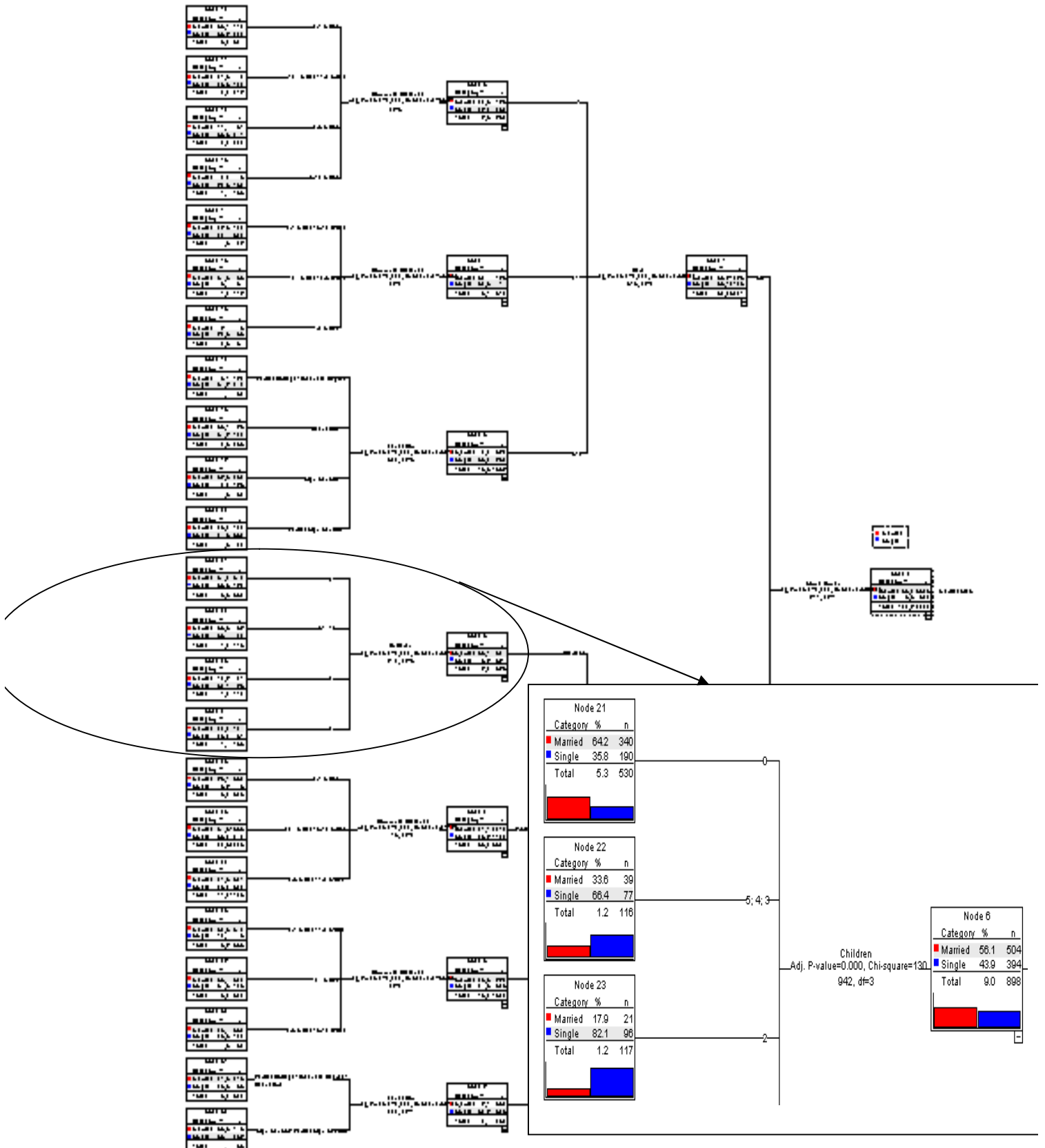
Bu ara tırmada öne sürülen hipotezi test etmek amacıyla, BIRCH ve CLARA algoritmalarının performansları kullanıldı. Algoritma veri setini yeniden taramaya ihtiyaç duymamaktadır. Kullanılan algoritmanın karma ıklı ının da  $O(n)$  olması önemli bir avantaj sa lamı tır.

### 3.10 Spss Programı le Elde Edilen Sonuçlar

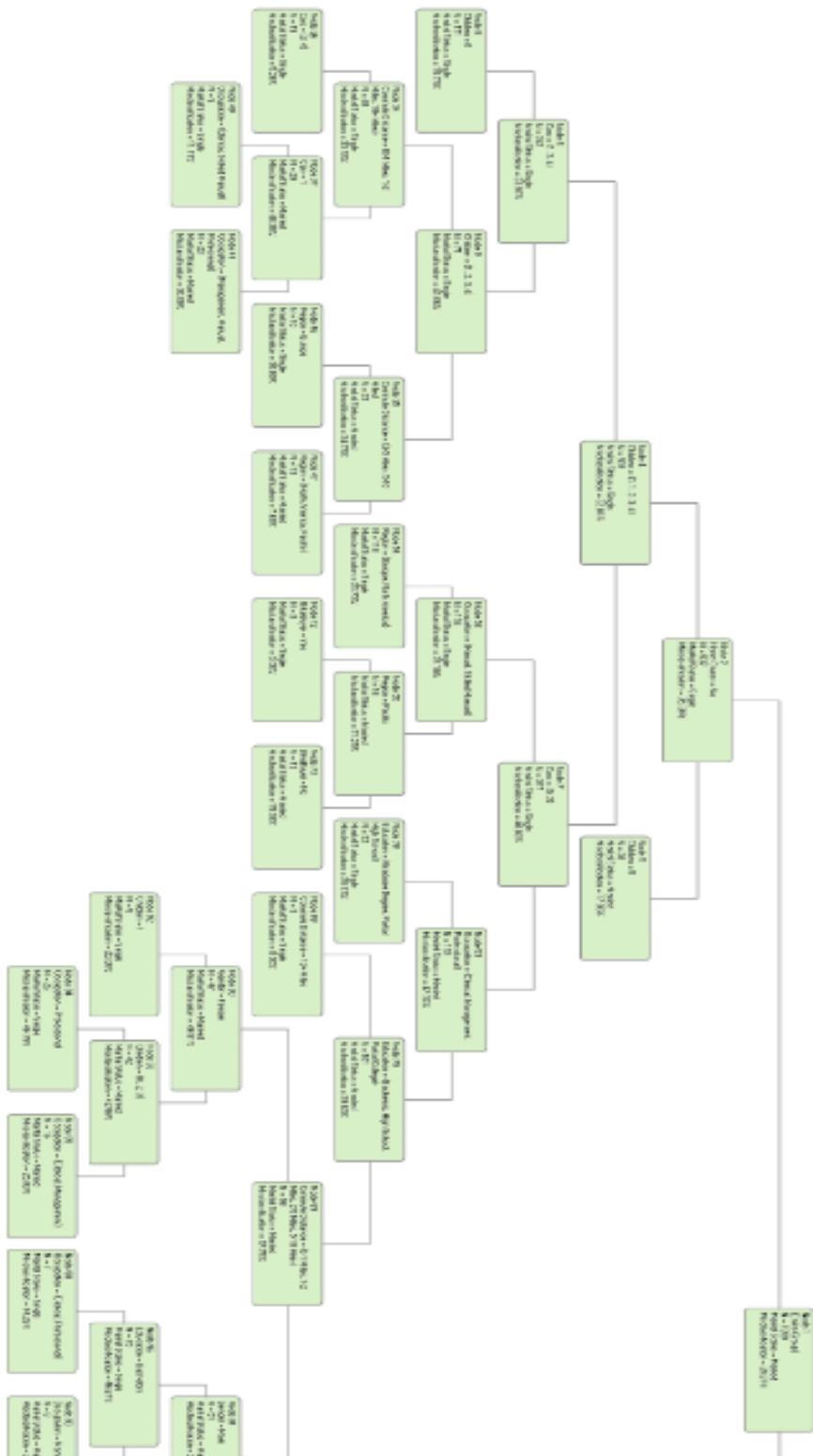
#### Model Summary

Specifications	Growing Method	CLARANS
	Dependent Variable	Marital Status
	Independent Variables	Gender, Children, Education, Occupation, Home Owner, Cars, Commute Distance, Region, BikeBuyer
	Validation	None
	Maximum Tree Depth	3
	Minimum Cases in Parent Node	100
	Minimum Cases in Child Node	50
Results	Independent Variables Included	Home Owner, Cars, Commute Distance, Education, Occupation, Children
	Number of Nodes	33
	Number of Terminal Nodes	23
	Depth	3

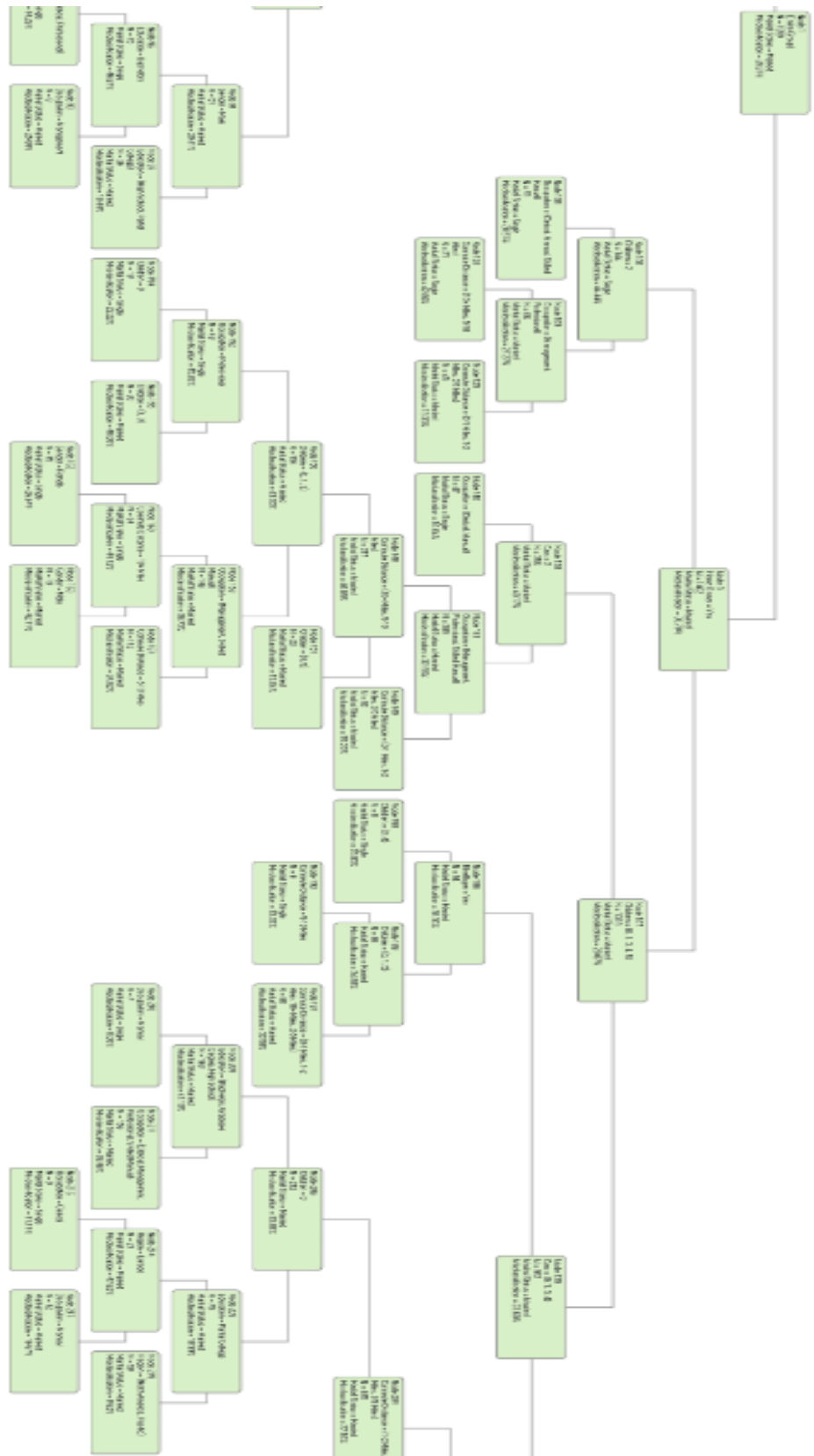
Tablo 3.6 SPSS Sonuçları



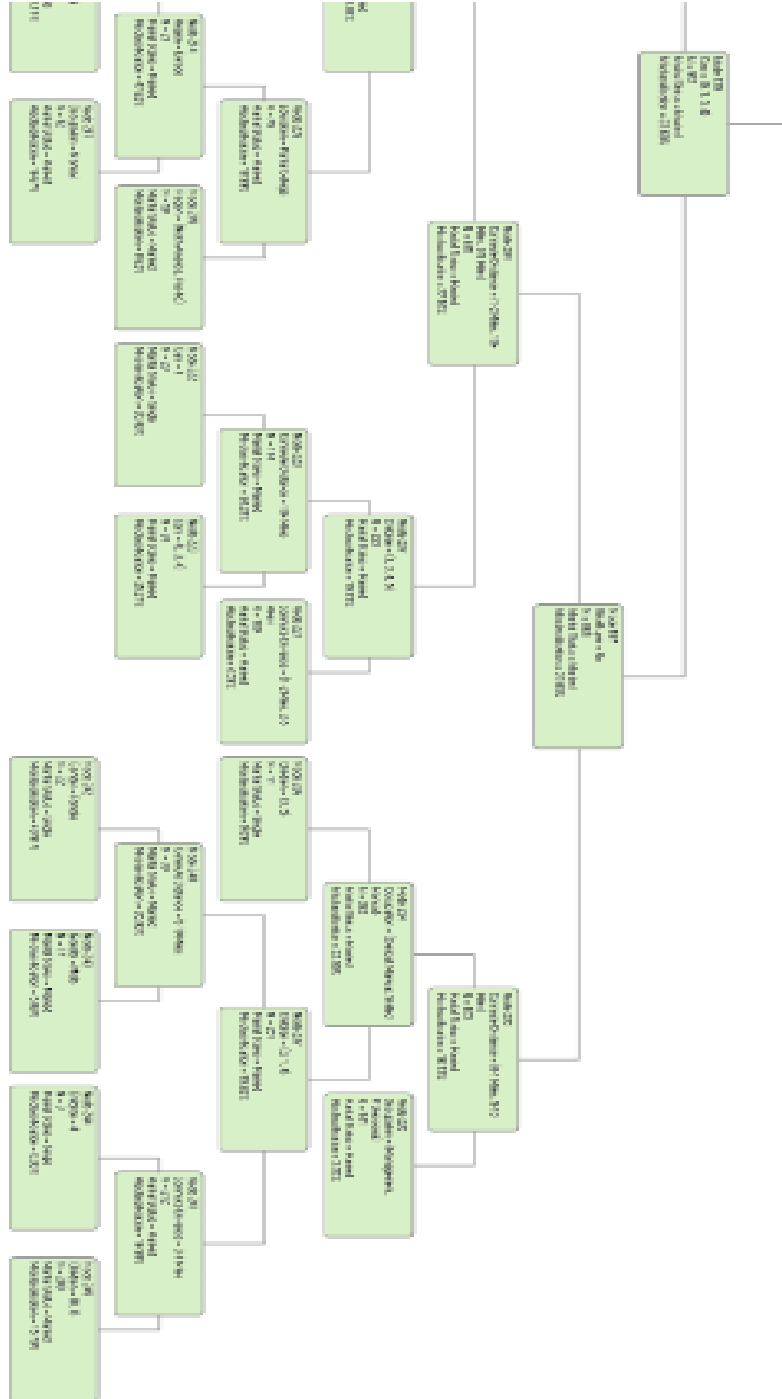
ekil - 3.6 SPSS A aci



ekil - 3.7 Algoritmadan Elde Edilen A aç - 1



ekil - 3.8 Algoritmadan Elde Edilen A aç - 2



ekil - 3.9 Algoritmadan Elde Edilen A aç - 3

## 4. BÖLÜM

### 4.1. SONUÇ VE ÖNERİLER

#### 4.1.1 Genel Sonuçlar

Bu bölümde, tez konusu kapsamında gerçekleştirilen uygulamanın sonuçları özetlenmektedir. Ayrıca algoritmanın kısıtları ve daha iyi hale getirilmesine yönelik çözümler de eleştirilmektedir.

Bu çalışmada, tasarlanan algoritmanın veri madenciliğinde kullanılanlar gibi kendine bir yer edeceğini gösterebilmek amacıyla gerçekleştirilmiştir. Bu algoritmanın, karmaşık ve daha fazla olan diğer algoritmalara oranla tercih edilip daha da geliştirilebilecek bir algoritma olduğu görülmektedir.

Veriyi kümelere ayırmak için ağaç metodunu Zhang, BIRCH algoritmasında kullanmıştır. Kümeleri biçimlendirmek için CF (Clustering - Feature) ağacı kullanılmış ve kümeleri tek bir tarama bazı tarıfsel istatıksel parametreleri kullanarak ve zıtlı tırma kümeleri oluşturdu. CF ağacında, tüm veri tabanını temsil eder; veri ayıklamanın bir dalı olan sınıflandırma için başka bir ağaç yöntemi kullanılabilir. veri tabanında yer alması olası, doğal kümeleri bulmak için bir karar ağacı kullanabiliyoruz. Bizim yöntemimizde robot, BIRCH yönteminde olduğu gibi tüm veri tabanını temsil etmiyor, buna rağmen tüm veri tabanı bölümünün en belirgin kısmıdır. Kısımların yaklaşıklık önemlerine-anlamlarına göre veri tabanını analiz ederek, doğal minimal kümelere ulaşıyoruz. Bu da oldukça iyi bir performans sağlamaktadır.

Daha önce bölüm 3.8 'de de değinildiği gibi algoritmanın bazı kısıtlılıkları bulunmaktadır. Geliştirme süreci devam eden bir algoritma olduğu için bu kısıtlar zamanla ortadan kalkabilecektir.

Çalışmada, farklı bir algoritma kullanılarak yeni bir yol elde edilmeye çalışılmış ve başarılı olabilecek bir sonuç elde edilmeye çalışılmıştır.

#### 4.1.2 Gelecek Ara tırmalar için Öneriler

Bizim yöntemimiz entropi ya da gini indeks yöntemine daha yakındır. Bu yöntemler de kök verisi kümelerine ayrılmı tır.

Ek olarak, verileri ve dalların arasında yer alan nodları ayırmak için kümeleri ve alt-kümeleri tanımladık. Kök de dahil olmak üzere, her nod veritabanını e ite en yakın parçalara ayırmak için hesaplanmı tır. Bu algoritmanın temel aldığı ölçüt budur. Böylece bir sonraki a amada kümeleme algoritmaları için daha iyi sonuç verecek bir algoritma elde edilmeye çalı ılmı tır. Bu algoritmanın, dikkate alındı ında ve di er sistemlere entegre edildi inde iyi sonuçlar verece i dü ünülmektedir.

## KAYNAKLAR

Agrawal Rakesh ve Shafer John c., Parallel Mining of Association Rules, IEEE Transactions on Knowledge and Data Engineering, Vol. 8, No. 6, 1996, s. 962-969.

Agrawal Rakesh ve Shafer John c., Parallel Mining of Association Rules, 1995.

Agrawal Rakesh ve Srikant Ramakrishnan, Fast Algorithms for Mining Association Rules, 20. VLDB Konferansı, İli, 1994, s.487-499.

Agrawal Rakesh ve Srikant Ramakrishnan, Mining Sequential Patterns, 11. Uluslararası Data Engineering Konferansı, 1995.

Agrawal Rakesh ve ark., Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications, ACM-SIGMOD Management of Data Konferansı, 1998, s. 94 - 105.

Agrawal Rakesh ve ark., Database Mining: A Performance Perspective, IEEE Transactions on Knowledge and Data Engineering, Aralık 1993, s.914-925.

Agrawal Rakesh, vd., Mining Association Rules between Sets of Items in Large Databases, 1993 ACM SIGMOD Konferansı Bildirisi, USA, 1993, s. 2.

Akpınar Haldun, Veritabanlarında Bilgi Keşfi ve Veri Madenciliği, İ.Ü. İktisadi İdari Bilimler Fakültesi Dergisi, C:29, sayı: İI Nisan 2000, s. 1-22.

Ankerst Michael ve ark., OPTICS : Ordering Points to Identify the Clustering Structure, ACM SIGMOD Management of Data Konferansı, Philadelphia, 1999, s. 49- 60.

Bacher Johann, Cluster Analysis, (çevrimiçi)  
[www.soziologie.wiso.unierlangen.de/koeln/script/chap3.pdf](http://www.soziologie.wiso.unierlangen.de/koeln/script/chap3.pdf) , (29/06/2008).

Berkhin Pavel, Survey of Clustering Data Mining Techniques, (çevrimiçi)  
<http://citeseer.nj.nec.com/berkhin02survey.html> , (7/4/2008).

Beyer Kevin ve ark., When Is 'Nearest Neighbor' Meaningful?, 7. Uluslararası Database Theory Konferansı (ICDT'99), srail, 1999, s. 217-235.



Cabena Peter ve ark., *Discovering Data Mining: From Concept to Implementation*, USA, International Business Machines Corporation, 1998, s.12.

Cheung David Wai-Lok, ve ark., *A Fast Distributed Algorithm for Mining Association Rules*, PDIS Bildirisi, 1996.

Colin Andrew ve Journal Dobbs, *Building Decision Trees with the ID3 Algorithm*, Haziran 1996.

Dempster Arthur. P. vd., *Maximum Likelihood from Incomplete Data via the EM Algorithm*, journal of the Royal Statistical Society Agglomerative, B serisi, Vol. 39, 1977, s. 1-38.

Dunham Margaret H., *Data Mining Introductory and Advanced Topics*, Prentice Hall, Pearson Education Inc., New Jersey, 2003, s. 8.

Efe Önder ve Kaynak Okyay, *Yapay Sinir Ağları ve Uygulamaları*, Bozaziçi Üniversitesi, 2004, s.10-12.

Ercio lu Ömer, *Parametric Approximation Algorithms for High-Dimensional Euclidean Similarity*, 5. Principles and Practice of Knowledge Discovery in Databases (PKDD'01) Konferansı, Freiburg, Almanya, Eylül 2001, s. 79-90.

Ester Martin ve ark., *A Density-Based Algorithm for Discovering Clusters in Large Databases with Noise*, 2. Uluslararası Knowledge Discovery and Data Mining Konferansı, 1996. Fausett Laurene, *Fundamentals of Neural Networks*, Prentice-Hall, 1994, s. 3.

Frenkel Brian, *Using Artificial Intelligence to Detect Fraud in Credit Cards*, (çevrimiçi) <http://tsel.cs.colorado.edu/~cs3202/papers2/BrianFrenkel.html> , (20/07/2008).

Giudici Paolo, *Applied Data Mining: Statistical Methods for Business and Industry*, Wiley, 2004, s. 83.

Guha Sudipto ve ark., *CURE: A Clustering Algorithm for Large Databases*, Teknik Rapor, Ben Laboratuvarları, Murray Hill, 1997.

Guha Sudipto ve ark., *CURE: An Efficient Clustering Algorithm for Large Databases*, ACM SIGMOD Konferansı, 1998, s. 73-84.

Han Jiawei ve Kamber Micheline, Data Mining Concepts and Techniques, Morgan Kaufman Publishers, Academic Press, 200 i, s. 106.

Joshi Karuna Pande, Analysis of Data Mining Algorithms, (çevrimiçi) [http://www.ernstedu.com/proj\\_rpt.htm](http://www.ernstedu.com/proj_rpt.htm), (04/04/2008)

Hinneburg Alexander ve Keim Daniel A., An Efficient Approach to Clustering in Large Multimedia Databases with Noise, Uluslararası Knowledge Discovery and Data Mining Konferansı (KDD'98), ABD, 1998, s. 58-65.

Kauffman L., Rousseeuw PJ, Finding Groups in Data: An Introduction to Cluster Analysis, John Wiley and Sons, 1990.

Khan Maleq ve ark., K-Nearest Neighbor Classification on Spatial Data Streams Using P-Trees, 6. Pasifik Asya Knowledge discovery and Data Mining Konferansı PAKKDD'02, Taiwan, 2002, s. 517-518.

Kimball Raph, DBMS Dealing With Dity Data (çevrimiçi) <http://www.dbms-mag.com/9609d14.html> (04/05/2008)

Kut Alp ve Yılmaz Sedat, Veri Madencilik Uygulamaları, nternet ders notu, (çevrimiçi) [courses.cs.deu.edu.tr/cse572/VeriMadencilikUygulamalar.doc](http://courses.cs.deu.edu.tr/cse572/VeriMadencilikUygulamalar.doc) (24/03/2008).

Lipmann Richard P., An Introduction to Computing with Neural Nets, IEEE ASSP Dergisi" Nisan 1987, s. 155- 162.

MacQueen, J., Some Methods for Classification and Analysis of Multivariate Observations, 5. Berkeley Matematiksel statistik ve Olasılık Sempozyumu, University of California Press, 1967, s. 281-297.

Manish Mehta vd., SLIQ: A Fast Scalable Classifier for Data Mining, S. Uluslararası Extending Database Technology Konferansı, Avignon, Fransa, Mart 1996.

Olson Clark F, Parallel Algorithms For Hierarchical Clustering, Teknik Rapor, 1993,

(Çevrimiçi) <http://citeseer.ist.psu.edu/17291.html> (30/06/2008). .

Orhunbilge Neyran, Uygulamalı Regresyon ve Korelasyon Analizi, Avcıol Basım\_Yayın, stanbul, 1999, s. 9.

Quinlan J. Ross, Induction of Decision Trees, Journal of Machine Learning, 1986, s. 81-106.

Quinlan, J. Ross, Simplifying Decision Trees", International Journal of Man-Machine Studies, sayı: 27,1987, s. 221-234.

Sibson R, An Optimally Efficient Algorithm for the Single Link Cluster Method, The Computer Journal, Vol. 16, Issue I, 1973.

Xiaowei Xu vd., A Distribution-Based Clustering Algorithm for Mining in Large Spatial Databases, ICDE, sayı: 14, 1998, s. 324-331.

Yamaç Aytekin, Temel Bileşenler Analizi ile Türkiye'nin AB Ülkeleri içindeki Yeri, Kara Harp Okulu Bilim Dergisi, 2002-2, (çevrimiçi), [www.kho.edu.tr/yayinlar/bilimdergisi/bilimder/doc/2002-2/4\\_bilder.doc](http://www.kho.edu.tr/yayinlar/bilimdergisi/bilimder/doc/2002-2/4_bilder.doc),(21/09/2008).

Yohannes Yisehac ve Webb Patrick, Classifications and Regression Trees, CART: A user Manuel For Identifying Indicators of Vulnerability to Famine and Chronic Food Insecurity, y.y, International Food Policy Research Institute, 1999, s. 15.

Zhang Tian, ve ark., BIRCH: An Efficient Data Clustering Method for Very Large Databases,