

Biyoinformatik Veri Analizinde R ile Hiyerarşik Kümeleme

Prof. Dr.
Zeynel CEBECİ



© **PAPATYA YAYINCILIK EĞİTİM**
EĞİTİM BİLGİSAYAR SİS. SAN. VE TİC. A.Ş.

Ankara Cad. Prof. F. Kerim Gökay Vakfı İşhanı Girişi
No: 11/6 Çağaloğlu/İstanbul

Tel : (+90 212) 527 52 96 GSM: (+90 532) 311 31 10
Faks : (+90 212) 527 52 97
E-Posta : admin@papatyabilim.com.tr
Web : www.papatyabilim.com.tr

Biyoinformatik Veri Analizinde R ile Hiyerarşik Kümeleme – Prof. Dr. Zeynel CEBECİ

1. Basım Haziran 2019 Deneme basımı
2. Basım Şubat 2021 Gözden geçirilmiş basım

Yayına Hazırlayan : Cengiz UĞURKAYA (Ph. D)
Üretim : Necdet AVCI
Pazarlama Satış : Mustafa DEMİR
Satış : TDK Bilim ~ www.tdk.com.tr
Kapak Resmi : Zeynel CEBECİ
Kapak Tasarım : Papatya Kelebek Görsel Tasarım Atölyesi
Sayfa Düzenleme : Papatya Kelebek Görsel Tasarım Atölyesi
Basım : Özkaracan Matbaacılık (Sertifika No: 12228)
Evren Mah. Gülbahar Cad. No: 62 Güneşli/İstanbul

© Bu kitabın her türlü yayın hakkı yayınevine aittir. Yayınevinden yazılı izin alınmaksızın alıntı yapılamaz, kısmen veya tamamen hiçbir şekil ve teknikle ÇOĞALTILAMAZ, BASILAMAZ, YAYIMLANAMAZ. Kitabın, tamamı veya bir kısmının fotokopi makinası, ofset vs. gibi teknikle çoğaltılması, hem çoğaltan hem de bulunduranlar için yasadışı bir davranıştır.

Cebeci, Zeynel.

Biyoinformatik Veri Analizinde R ile Hiyerarşik Kümeleme / Zeynel Cebeci – İstanbul: Papatya
Yayıncılık Eğitim, 2021

xii, 300. ; 24 cm

Kaynakça ve dizin var.

ISBN 978-605-9594-44-8

Sertifika No: 11218

1. Veri Önışleme 2. Geçerlilik Testleri 3. Dendrogram Çizme 4. İleri Algoritmalar

I. Title

QA76.9 .D35 C64 2011

Nalan ve Çağatay'a

İçindekiler

Önsöz	xii
Kısaltmalar	xi
Bölüm 1. Kümeleme Analizine Giriş	13
Küme ve Kümeleme	13
Kümeleme Yöntemleri	15
Kümeleme Analizinin Aşamaları	18
Kümeleme Analizinin Uygulama Alanları	19
Hiyerarşik Kümeleme Yöntemleri	20
Bölüm 2. Birleştirici Kümeleme Yöntemleri	23
Örnek Veri ve Uzaklık Matrisinin Oluşturulması	26
Tek Bağlantı Yöntemi	29
Tam Bağlantı Yöntemi	34
Ortalama Bağlantı Yöntemi	39
Merkezci Yöntemler	43
Yöntemlerle İlgili Genel Özetleme	59
Bölüm 3. Ayırıcı Kümeleme Yöntemleri	61
Monotetik Ayırıcı Yöntemler	62
Politetik Ayırıcı Yöntemler	67
İleri Algoritmalar	73
Bölüm 4. Kümeleme Sonuçlarını (Dendrogram) Görselleştirme	75
Dendrogram (Ağaç Grafiği)	75
Afiş Grafik	76
Arapsaçı Grafiği	77
Korelasyon Grafikleri	79
Isı Haritaları	80
Bölüm 5. Küme Sayısının Belirlenmesi	81
Kofenetik Korelasyon Katsayısı	82
Birleşme Uzaklığındaki Değişim	84
Birleşme Katsayısı	84
<i>Calinski</i> ve <i>Harabasz</i> İndeksi	85

Sahte T2 İndeksi	85
Kübik Kümeleme Kriteri	86
Hata Kareler Ortalamasının Karekökü	86
Ortak Standart Sapma	87
R2 İstatistiği	87
Yarı Kısmi R2 İstatistiği	88
Diğer Ölçütler	89
Bölüm 6. R'nin Kurulması ve Çalıştırılması	91
R'nin İndirilmesi	92
R'nin Kurulması	95
R'de Çalışma	98
R Paketleri	101
R Betikleri	105
R ile Programlama	107
R Fonksiyonları	110
R'den Çıkma	115
Bölüm 7. Veri Kümeleri ve Veri Önışleme	117
Veri Kümesi Oluşturma	117
İçsel Veri Üretimi	120
Dosyalardan Veri Okuma	121
Bölüm 8. Veri Kümeleri ve Önışleme	135
Kayıp Değerlerin İşlenmesi	135
Veri Dönüştürme	148
Standartlaştırma	150
Bölüm 9. R'de Birleştirci Kümeleme Analizi	151
Uzaklık Matrisinin Oluşturulması	151
Kümeleme İşlemi	151
Dendrogram Çizme	154
Kümeler ve Küme Elemanlarının Belirlenmesi	157
Dendrogram Kesme	163
Küme Serpilme Grafiği	165
Esnek Beta ile Kümeleme	167
Kümeleme Sonuçlarının Karşılaştırılması	168
Karışıklık Matrisi	172
Kutu-Bıyık Grafikleri	173
Birleşme/Ayrılma Yüksekliği Grafikleri	174

Grafiklerin Saklanması	175
Kofenetik Korelasyonlar	178
Performans Testi	181
Bölüm 10. R'de Ayırıcı Kümeleme Analizi	185
Diana ile Ayırıcı Kümeleme Analizi	185
Mona ile Ayırıcı Kümeleme Analizi	192
Bölüm 11. Kümeleme Geçerlilik Testleri	195
Bölüm 12. İleri Görselleştirme Araçları	201
Hclust nesnesinin diğer nesnelere dönüştürülmesi	205
Dendrogram Nesneleri	205
Phylo Nesneleri ve Filogenetik Ağaçlar	209
Dendextend ile Grafik Analiz	217
Dendrogram Nesnesi Özelliklerinin Ayarlanması	218
Arapça Grafiği	224
Dendrogramların Farklılıklarını Bulma	230
Korelasyon Matrisi ve Grafiği	231
Fowlkes-Mallows İndeksi	234
Bk Grafiği	236
Farklı Dendrogram Sunumları	237
Isı haritaları	242
Bölüm 13. Mikrodizelerde Kümeleme Analizi	251
Akciğer Kanseri Verikümesi ile Analiz	251
Küçük Yuvarlak Mavi Hücre Tümörleri Verikümesi ile Analiz	261
Akut Lenfositik Lösemi Verikümesi ile Analiz	269
Kaynakça	279
İngilizce Türkçe Terimler Kılavuzu	287
Yazarımız ~ Biyografi	289
Dizin	291

Önsöz

Hiyerarşik kümeleme, veri madenciliğinde örüntü tanıma, makine öğrenmesi, pazarlama ve müşteri yönetimi, gen keşfi, ilaç tasarımı gibi birçok fen, sosyal ve yaşambilimleri alanında en yaygın kullanılan keşifçi ya da açınsal istatistiksel analizdir. Çoğu kez, analiz öncesi belli bir küme sayısı (k parametresi) ve küme merkezleri için başlatma değerleri istememesi gibi önemli avantajlar sunması nedeniyle K-ortalamlar ve K-ortancalar gibi yine yaygın olarak kullanılan hiyerarşik olmayan yöntemlere göre tercih edilmektedir. Analiz için ihtiyaç duyulan yegâne şey veri öğeleri arasındaki benzerlik ölçüsü olduğundan uygulanması kolay bulunmaktadır.

R, istatistik ve grafik analiz için güçlü bir hesaplama ortamı ve programlama dilidir. R projesi kapsamında GNU lisansı ile dağıtılan açık kaynak ve özgür yazılım olması, yeni algoritmaları geliştirmek ve test etmek için etkin ve kolay bir programlama ortamı sunması nedeniyle istatistik, matematik, bilgisayar bilimleri ve biyoinformatik gibi alanlarda hızla popüler olmasını sağlamıştır. Dünya genelinde hemen her kuramsal ve uygulamalı alanda çalışan araştırmacılar, veri analizcileri ve öğrenciler tarafından yaygın şekilde kullanılmaktadır. Günümüzde R için çok çeşitli alanlarda çalışan araştırmacılar tarafından geliştirilmiş 16000'e yaklaşan sayıda R paketi bulunmakta olup CRAN, BioConductor ve Github üzerinden dağıtılmaktadır.

Bu kitap veri madenciliğinde ya da bilgi keşfinde önemli bir açınsal istatistik aracı olarak hiyerarşik kümeleme analizine yöntemler, R ile uygulamalar ve bazı gerçek veri kümeleri üzerinde analiz örnekleriyle kapsamlı bir bakış sağlamaktadır.

Bu kitap, daha önce R ile çalışmamış olanların birkaç gün içinde ileri düzeyde kümeleme analizi yapmalarını sağlayacak bir yaklaşımla yazılmıştır. Kitap, hem veri madenciliği ve istatistik konulu dersler için bir uygulama rehberi hem de biyoteknoloji ve biyoinformatik bilim dallarında çalışan araştırmacılar için bir başvuru eseri olacak şekilde tasarlanmıştır.

Yararlı olması dileğiyle,

Zeynel Cebeci
Adana, 2020

Kitap Hakkında

Kitabımızın ilk bölümünde küme, kümeleme terminolojisi ve kümeleme yöntemleri taksonomisi anlatılmaktadır.

İkinci bölümde birleştirici kümeleme yöntemlerinin teorik temelleri küçük bir veri kümesinde uygulamalı olarak tüm ayrıntıları ile sunulmaktadır.

Üçüncü bölüm ise ayrıcı kümeleme yöntemlerini açıklamakta ve örneklemektedir. Monotetik ve politetik ayrıcı yöntemleri verilmiştir ve ileri algoritmalar ele alınmıştır.

Dördüncü bölümde kümeleme sonuçlarının görselleştirilmesi için kullanılan temel grafik türleri örneklerle açıklanmıştır.

Beşinci bölümde küme sayısının belirlenmesinde kullanılan ölçütler ve teknikler tanıtılmıştır. Çeşitli örnekler verilerek konunun daha kolay anlaşılmasına gayret edilmiştir.

Altıncı bölümde R istatistiksel hesaplama ve grafik analiz yazılımının indirilmesi ve kurulması, R'de çalışma, R dilinin tanıtımı ve R fonksiyonlarına giriş yapılmıştır.

Veri okuma ve ön işleme örnekleri ise yedinci ve sekizinci bölümlerde sunulmuştur. Dokuzuncu bölüm ayrıntılı olarak R ile birleştirici kümeleme; onuncu bölüm ise R ile ayrıcı kümeleme analizini açıklamaktadır.

On birinci bölümde kümeleme geçerlilik testleri konusu örneklerle anlatılmıştır.

On ikinci bölümde ise ileri düzeyde görselleştirme araçları ve grafik türleri sunulmaktadır.

On üçüncü bölüm gen ifade profillerinin analizi ve gen keşfi çalışmalarında kullanılan teknik ve yöntemlerin üç gerçek veri kümesine uygulanması ve yorumlanmasını kapsamaktadır. Mikrodizilerde kümeleme analizi anlatılmıştır.

Kısaltmalar

Kısaltma	Terim
AU	Approximately Unbiased
BP	Bootstrap Probability
CA	Clustering Algorithm
CCC	Cubic Clustering Criterion
CPCC	Cophenetic Correlation Coefficient
HAC	Hierarchical Agglomerative Clustering
HDC	Hierarchical Divisive Clustering
MCAR	Missing Completely At Random
MNAR	Missing Not At Random
RMSSTD	Root Mean Square Standard Deviation
VRC	Variance Ratio Criterion
WPGMA	Weighted Pair-Group Method using Arithmetic Averages
WPGMC	Weighted Pair-Group Method using the Centroid Average
UPGMA	Unweighted Pair-Group Average

Kitaptaki Örneklerin Uygulanması

Kitapta herhangi bir konu ile ilgili örnekler aşağıda görüldüğü gibi bir kutu içinde Courier yazı tipinde verilmiştir.

```
> xdianak <- diana(x, metric = "euclidean",  
+   stand = F)
```

Komut örneklerinde:

- a) Satır başlarındaki > simgesi R tarafından otomatik olarak görüntülenen komut satırı göstergesidir. Örnekler çalıştırılırken yazılmaması gerekir.
- b) Satır başlarındaki + simgesi R tarafından otomatik olarak görüntülenen komut satırı devam simgesidir. Girilen komutlar tek bir satıra sığmadığında R izleyen satıra geçer ve satırın başına önceki satırdaki komutun devamı olduğunu gösteren + işareti koyar. Kitaptaki örnekler çalıştırılırken yazılmaması gerekir.
- c) Komutlar çalıştırdıktan sonra sonuçları görüntülenir. Bazı sonuçlar [1], [7] gibi köşeli parantezler içinde sayılarla başlayan şekilde gösterilir. Bunlar sonuçlar bir vektör olduğunda satır başındaki elemanın indis numarasıdır. Sadece bilgi amaçlıdır. Dikkate alınmaması gerekir.
- d) Kitapta gerek bazı komutlar ve gerekse bazı sonuçları vurgulamak amacıyla koyu/kalın yazılmıştır. Komut çalıştırırken dikkate alınmamalıdır.
- e) Kitabın elektronik sürümünde (e-kitap) R komut örnekleri önce fare ile işaretlenip blok halinde kopyalanır. Daha sonra R çalışma ortamında sağ fare düğmesi tıklanarak *Sadece komutları yapıştır* seçeneği çalıştırılır. Bu işlem kopyalanan komut blokunda bulunan > ve + simgelerini ve diğer metinleri otomatik olarak temizleyerek sadece R komutlarının çalıştırılmasını sağlar.
- f) Kitabın elektronik sürümünde (e-kitap) R komut örnekleri > ve + simgesi olmadan kopyalanıp çalıştırılmak istenirse R çalışma ortamında sağ fare düğmesi tıklanarak *Yapıştır* seçeneği seçilmelidir. Bu işlem kopyalanan komut blokundaki kodların hemen çalıştırılmasını sağlar.

Bölüm 1

Kümeleme Analizine Giriş

Küme ve Kümeleme

Kümeleme bireyler, nesnelere, olaylar ve durumlar gibi soyut ve somut veri öğelerini benzerliklerine göre gruplara ayırmak işlemidir. Bir başka deyişle öğeleri özelliklerine göre gruplandırma işidir. Günlük yaşamda kümelemenin çok sayıda örneğine rastlamak mümkündür. Örneğin, süpermarketlerdeki tatlılar, baharatlar, süt ve süt ürünleri, et ve et ürünleri, temizlik ürünleri vb. ürünler reyon olarak adlandırılan yüzlerce farklı bölümdeki raflarda sergilenirler. Botanik bahçelerinde bitkiler kayalık bitkileri, güller, tıbbi bitkiler ve ibreliler gibi farklı alanlarda gruplanarak yetiştirilirler. Kütüphaneler veya kitapçılardaki kitaplar konularına göre farklı raflarda dizilirler. Belli bir coğrafya bitki örtüsü veya hayvan biyoçeşitliliği bakımından farklı alt bölgelere veya ekosistemlere bölünürler. Sosyologlar toplumu farklı demografik bölümlere ayırarak incelerler. Pazarlama faaliyetlerinde müşteriler bazı özellikler açısından gruplara ayrılarak değerlendirilir. Günlük yaşamın bir parçası olarak her gün kullanılan kara taşıtları segment olarak tanımlanan farklı gruplara ayrılırlar.

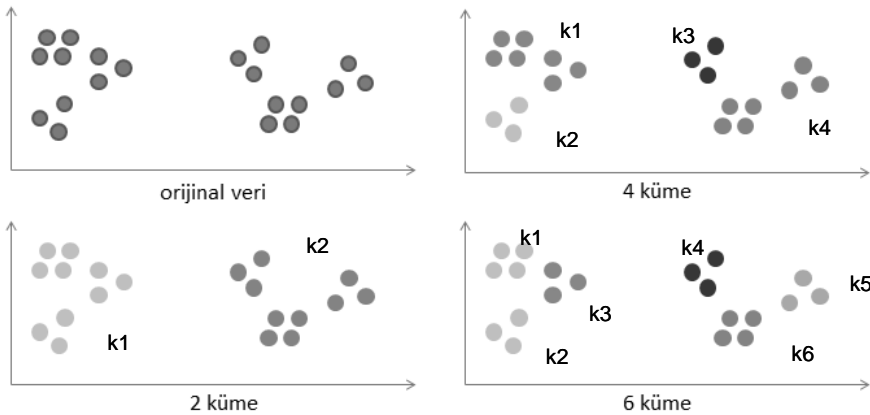
Yukarıdaki örneklerin hepsinde öğeler belli özelliklerine göre ait oldukları küme (sınıf veya grup) içine atanırlar. Öğelerin sınıfa veya gruba atanmasında yani kümelemede gerekli özellikler veya kriterler atama işlemi öncesinde bellidir. Yapılan işlem bir sınıflama işlemidir. Sınıflama analizi ile öğeler ait oldukları sınıf etiketleriyle etiketlenmekte ya da işaretlenmektedirler. Bu bir önbilgi dâhilinde gerçekleştirilir ve yani eğitilmiş bir sınıflama işlemidir. Bu nedenle makine öğrenmesi gibi bazı bilim dallarında *eğitilmiş öğrenme* olarak da adlandırılmaktadır.

Kümeleme analizi incelenen veri kümesindeki öğeleri gruplamak için kullanılan yöntem ve teknikleri kapsamaktadır. Ancak kümeleme analizinde öğelerin atacağı sınıf üyelikleri hakkında daha önce elde edilmiş herhangi bir önbilgi yoktur. Öğeler ve içinde yer aldıkları gruplar ya da kümeler daha önceki önbilgi ile değil, doğrudan analiz ile saptanmaya çalışılmaktadır. Kümeleme analizi tipik şekilde verilen özellikleri kullanarak ve neden olduklarını araştırmaksızın sadece sınıfsal yapıları keşfetmeyi hedefler. Bu yüzden kümeleme analizi normalde büyük bir problemin küçük bir ön çözümünü amaçlar. Sonuç olarak verideki yapıları ya da örüntüyü görmek ve daha sonraki analizlerde kullanmak üzere yapılan keşifsel bir analizdir. Bu nedenle, bu tür bir sınıflama *eğitilmiş öğrenme* işlemidir. Bazı disiplinlerde, örneğin biyolojide *sınıflama analizi* veya *sayısal taksonomi* olarak da adlandırılmaktadır.

Yukarıdaki açıklamalardan da anlaşılacağı üzere, kümeleme analizi örüntü tanıma ve önceden atanmış etiketlerle sınıflama kurallarını araştıran ayırma analizi ve karar analizi gibi yöntemlerden farklıdır. Ayırma analizi, öğelerin ait olduğu küme (ya da grup veya sınıf) üyelikleri hakkında önbilgilere gereksinim duyar. Bu bilgiler sınıflama kuralı oluşturmak için kullanılır. Kümeleme analizi sınıfları bulmada bir ön sonuç sunmak için kullanışlı olabilirken diğer analizlerde çok daha fazlası söz konusudur. Örneğin, veri kümesi öğelerini temsil edecek özelliklerin neler olduğuna karar vermek örüntü tanıma gibi uygulama alanlarında temel amaçlardan biridir.

Kümeleme analizi tek bir yöntem veya teknik değildir. Aksine veri öğelerinin kendileri ve ilişkilerini açıklamak amacıyla, onları küme olarak adlandırılan gruplara ayırmak için kullanılan yüzlerce farklı algoritma veya yaklaşımları kapsayan çok değişkenli analiz yöntemlerinden oluşmaktadır. Aslında verilerdeki küme veya sınıfları bulmak sanatı olarak da tanımlanmaktadır (Kaufman ve Rousseeuw, 1990). Kümeleme birbirine benzeyen öğelerin yani birbirine yakın olanların aynı gruba; farklı olanların ise diğer gruplara yerleştirilmesiyle gerçekleştirilir. Amaç, kendi içlerinde mümkün olduğunca birbirine benzer (homojen) öğelere sahip olan ancak birbirinden mümkün olduğunca uzak grupların elde edilmesini sağlamaktır. Bir başka deyişle, aynı kümedeki öğelerin olabildiğince birbirine benzeyen öğeler olmaları; diğer kümelerdekilerden farklı olmaları kümeleme analizinde temel amaçtır.

Kümelerdeki öğelerin benzerliği yüksek ve kümeler arasında farklılık ne denli fazla ise oluşturulan kümeler belirgin biçimde birbirinden ayrık olurlar ki, bu kümeleme işleminin iyi ve kaliteli yapıldığını gösterir. Ancak elde edilecek küme sayısı ve bunların verideki hangi öğeleri kapsayacağı konusunda etkin bir istatistiksel test söz konusu değildir. Bir öğenin birden fazla kümeye üye olmasına imkan veren bulanık kümeleme bir istisna olmak üzere, kümeleme analizinde genel olarak öğelerin birbiriyle çakışmayan kümelere ayrılması araştırılır.



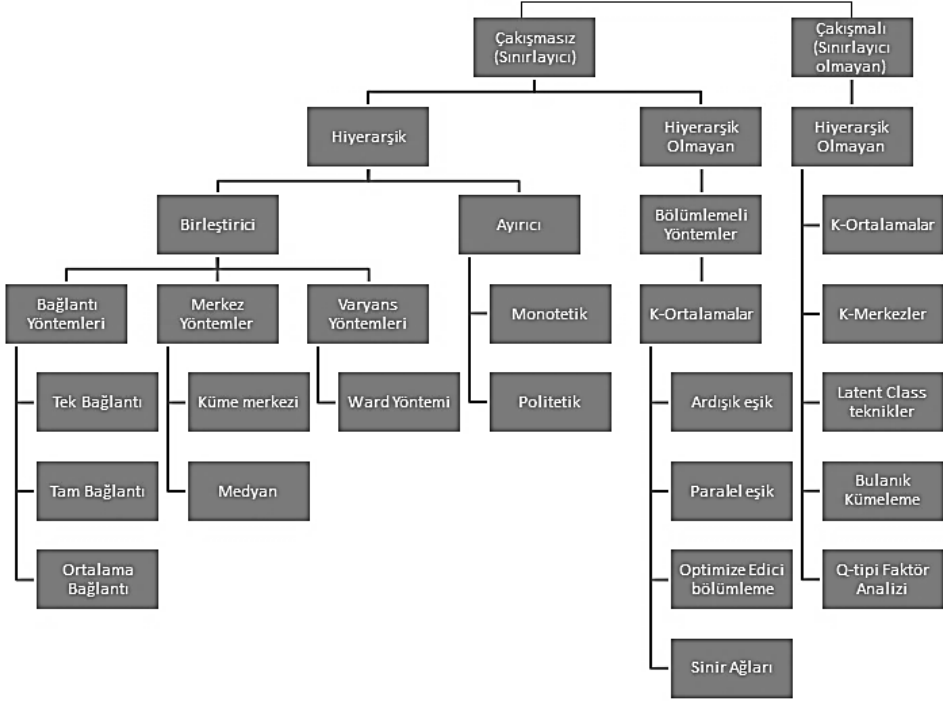
Şekil 1. Aynı veriden elde edilmesi muhtemel kümeler

ok sayıda kmeleme ynteminin mevcudiyetine raėmen kmeleme iřleminin geerliliėi, yani analiz sonucunda elde edilen kme sayıları ve kmelenme řekli bakımından eřitli glkler sz konusu olabilmektedir. řekil 1’de 20 ėe iin muhtelif kmeleme sonuları grlmektedir. řeklin sol stndeki grafiėe bakıldıėında orijinal verinin iki kmeden oluřtuėu/oluřacaėı dřnlebilir. Ancak daha dikkatli bakıldıėında bu iki kmeden herbirinin er alt kmeden oluřtuėu dřnlerek toplam 6 kme olduėu sonucuna varılabilir (řekil 1’in saė alt kesindeki gsterildiėi gibi). Ancak bu deėerlendirme insan grme sisteminin bir kusuru da sayılabilir. Orijinal veriyi drt kmeye ayırmak da akla gelebilir varılabilir (řekil 1’in saė st kesindeki gsterildiėi gibi). Ancak bu kmelemeler iin makul bir neden belirtmek ve yorumlamak fazlasıyla gttr. Bu nedenle kmeleme analizinde, kmeleme iřleminin geerliliėi zerinde durulması gereken ok nemli bir husus olup yoėun alıřılan arařtırma konuları arasında yer almaktadır. Diėer yandan kmeleme iřleminin geerliliėi, yani ka kme oluřturulacaėına karar vermek iin zerinde alıřılan veriyi iyi incelemek ve anlamak gerekli olup alan uzmanlıėı ve deneyim de gerektirir.

Kmeleme Yntemleri

İlk olarak Tryon (1939) tarafından kullanıldıėından beri zerinde en ok alıřılan ok deėiřkenli istatistik analiz yntemlerinden biri olarak kmeleme analizi iin geliřtirilmiř ok sayıda yntem, teknik ve algoritma sz konusudur. Downs ve Barnard (1995) kmeleme yntemlerinin en genel řekliyle *hiyerarřik* ve *hiyerarřik olmayan* yntemler olmak zere iki gruptan oluřtuėunu belirterek ok basit bir sınıflama vermiřlerdir. Kmeleme yntemleri zerinde birok arařtırıcı tarafından yapılan inceleme ve derlemelerde eřitli taksonomiler nerilmiř olmasına karřın ideal bir taksonomi olmadıėını syleyebiliriz. nk bir kmeleme yntemi farklı ltlere gre farklı sayıda kmeleme yntemi sınıfında yer alabilir. Bu nedenle kmeleme yntemleri iin belli bir sınıflama/taksonomi vermek yerine Sneath ve Sokal’ın yaptıėı gibi herhangi bir kmeleme yntemini ařaėıdakilere uygunluk bakımından eřitli kategorilere dhil etmek mmkndr (Sneath & Sokal, 1973).

- *Hiyerarři: Hiyerarřik yntemler / Hiyerarřik olmayan*
- *Kmeleme yn: Birleřtirici / Ayırıcı*
- *Geiř durumu: Tek geiřli / ok geiřli*
- *Eřanlılık: Ardıřık / Eřanlı, Tek deėiřkenli geiř / ok deėiřkenli geiř*
- *lek: Yerel / Genel*
- *Yineleme: Direkt / Yinelemeli (İteratif)*
- *Aėırlık: Aėırlıklı / Aėırlıksız*
- *Uyumsallık: Adaptif / Adaptif olmayan*
- *Katılık: Katı / Yumuřak*
- *akıřma / Sınırlama: akıřmalı / akıřmasız, Sınırlayıcı / Sınırlayıcı olmayan*



Şekil 2. Kümeleme yöntemleri taksonimisi

Yukarıda anlatıldığı gibi kümeleme yöntemlerine ait belli bir taksonomi önermek zor olsa da başlangıç itibariyle kullanılacak basit taksonomilerden biri aşağıda sunulmaktadır:

- Hiyerarşik yöntemler
 - Birleştirici yöntemler
 - Bağlantı yöntemleri
 - Tek bağlantı
 - Tam bağlantı
 - Ortalama bağlantı
 - Merkezci yöntemler
 - Küme merkezleri
 - Ortanca
 - Varyans yöntemleri
 - Ward
 - Diğer Yöntemler
 - BIRCH, CURE, CHAMALEON, OPTICS

- Ayırıcı yntemler
 - Monotetik ayırıcı yntemler
 - MONA algoritması
 - Politetik ayırıcı yntemler
 - DIANA algoritması
- Hiyerarřik Olmayan Yntemler
 - Blmleyici yntemler
 - Tek geiřli yntemler
 - Leader algoritması
 - Uzaklık Tabanlı Yntemler
 - Katı algoritmalar
 - K-ortalamalar
 - K-medoid (PAM)
 - K-mod
 - SOM
 - CLARANS
 - Yumuřak algoritmalar
 - Bulanık C-ortalamalar (FCM)
 - Olabilirlikli C-ortalamalar (PCM)
 - Bulanık ve Olabilirlikli Karma algoritmalar (FPCM, PFCM)
 - Yoęunluk tabanlı yntemler
 - En yakın komřu tabanlı yntemler
 - DBSCAN, OPTICS, DENCLUE, CLIQUE...
 - Izgara tabanlı yntemler
 - BANG, STING, CLIQUE, WaveCluster, PROCLUS, ORCLUS, OPTIGRID, ...
 - Model tabanlı yntemler
 - MLE
 - Beklenti oklaması (EM)
 - Grafik tabanlı modeller
- Karma Yntemler

Kümeleme Analizinin Aşamaları

Kümeleme analizi, veri öğelerini belli kümelere ayıran bir algoritmanın işletilmesi veya hesaplamasının yapılmasından ibaret tek aşamalı bir işlem değil aşağıda listelenen aşamalardan oluşan çok aşamalı bir analizdir.

1. Önışleme işlemleri
 - a. Veri kümesinin okunması ve analize hazırlanması
 - b. Kayıp değerlerin tahmin edilmesi
 - c. Normalleştirme veya standartlaştırma işlemlerinin yapılması
2. Yakınlık ölçülerinin hesaplanması ve uzaklık matrisinin oluşturulması
3. Kümeleme yöntemi/algoritmasının belirlenmesi
4. Kümeleme analizinin yapılması
5. Kümeleme grafiklerinin çizilmesi
6. Kümeleme geçerliğinin kontrol edilmesi
7. Kümeleme profilinin oluşturulması
8. Kümeleme sonuçlarının yorumlanması

Kümeleme analizinde ilk aşama hemen her istatistik analiz veya veri madenciliği için de geçerli olan analiz öncesi işlemlerdir. Veri önışleme olarak adlandırılan bu aşamada analiz edilecek veri kümesinin okunması ve analize hazır hale getirilmesi hedeflenmektedir. Kümeleme işlemi uygulamadan önce veri kümesinin incelenmesi varsa kayıp veya eksik öğelerin tamamlanması veya tahmin yoluyla veriye eklenmesi gereklidir. Araştırmacının veri türlerini ve boyutunu incelemesi uygun analiz yöntem veya tekniğine karar vermek açısından son derece kritik bir konudur. Ayrıca nicel sürekli değişkenlere bazı hallerde standartlaştırma ve/veya normalleştirme işlemlerine de analiz öncesi gerek duyulabilir. Verideki aykırı değerler veya gürültülerin saptanması ve elimine edilmesi de gerekli olabilir.

Kümeleme analizinde uygun yöntem karar vermek de önemlidir. Karar sırasında gerek çalışma disiplini, gerek veri türü ve boyutu ve gerekse işlem hızı ve maliyeti belirleyici olmaktadır. Örneğin, biyolojik ya da ekolojik taksonomi çalışmalarında sıkça hiyerarşik kümeleme kullanılmaktadır. Öte yandan bilgi sistemlerinde veri madenciliği veya sosyal medya için metin madenciliği gibi uygulamalarda hiyerarşik olmayan bir kümeleme yöntemi; mümkün olduğunca hızlı işlem imkânı sağlayan bir yöntem veya algoritma amaçlanabilmektedir. Günümüzde veri madenciliği uygulamalarında büyük-veri (big data) söz konusudur ve mutlaka çok hızlı çözüm sağlayan yöntemlere gereksinim duyulmaktadır.

Küme geçerliliği, yapılan kümelemenin ne kadar uygun olduğunu ya da kümelemenin uyum iyiliğini test etmektir. Bunun için çeşitli çapraz doğrulamalar ve istatistiksel ölçütlere başvurulur. Küme profileme, ilgili boyutlarda nasıl farklılaştığını izah etmek için her bir kümenin karakteristiklerini açıklar. Bunun için ayırma analizi veya varyans analizi (ANOVA) gibi yöntemlere başvurulur.

Kmeleme Analizinin Uygulama Alanları

Kmeleme analizi, tarımdan astronomiye; gen biliminden bilgisayar bilimlerine hemen hemen tm temel ve uygulamalı disiplinlerde kullanılan bir analizdir. Tablo 1’de bilim dalları itibariyle kmeleme analizinin kullanıldıđı alanlardan bazıları listelenmektedir.

Tablo 1. Kmeleme analizi uygulama alanı rnekleri

Uygulama Alanı	Kullanım rnekleri
Ekoloji	Farklı ortamlardaki organizma topluluklarının meknsal ve zamansal karřılařtırmalarını yapmak.
Biyoloji	İliřkili ifade rntlerinden gen gruplarını saptamak.
Biyoinformatik	Evrim biyolojisi ve biyoinformatik alıřmalarında nemli bir konu olan gen aileleri iin homolog diziliřleri gruplandırmak.
Genetik	
Kimya	ok sayıda bileřiđi ok sayıda zelliđine gre gruplandırarak yapısal benzerliklerini ortaya koymak.
Eczacılık	İla keřif alıřmaları yapmak.
Tıp	PET taramalarında  boyutlu grntlerde kan ve doku trleri arasındaki farklılıđı, farklı grupları ortaya koymak. Hastalıklar iin tedavi yntemleri ve belirtileri kmeleyerek tıbbi taksonomiler oluřturmak.
Sosyal Ađlar	Sosyal ađ kullanıcı gruplarını belirlemek.
neri Sistemleri	Kullanıcılara kendilerine benzeyenlerden retilen grup zelliklerine gre rn ve hizmet nerilerinde bulunmak.
Bilgi eriřimi	Arama ve bilgiye eriřimde Web siteleri ve belgeleri gruplandırmak.
Grnt İřleme	Nesne tanıma, kenar bulma amalarıyla grnty farklı blmlere ayırmak (grnt blmleme).
Su analizi	Farklı alanları su zelliklerine gre gruplara ayırmak.
Eđitim	Okullar ve đrencileri zelliklerine gre gruplandırmak.
İklim/Meteoroloji	İklim rejimleri, hava durum analizi iin gruplandırmalar yapmak.
Tarım	Benzer arazi kullanım zelliklerine sahip tarımsal alanları tanımak ve sınıflandırmak.
Pazarlama	Ele alınan zelliklere gre farklı mřteri gruplarını belirlemek ve hedef pazarlama programları tasarlamak.
Kentsel planlama	Tr, deđer ve cođrafı konumuna gre konutları gruplandırmak ve gruplara iliřkin planlama alıřmaları yapmak.
Deprem alıřmaları	Fay hatları boyunca deprem merkez slerini gzlemek.
Bilgisayar/Yazılım Mhendisliđi	Grnt analizi/iřleme, Grnt blmleme, Makine đrenmesi, rnt tanıma, Arama motorları, Bilgi eriřimi, Veri madenciliđi, Ses madenciliđi, Metin madenciliđi, Sunucu/Web kmeleme.

Hiyerarşik Kümeleme Yöntemleri

Kümeleme analizinde en çok kullanılan yöntemlerden olan hiyerarşik kümeleme, kümeler biri diğerinin içinde yer alacak şekilde yuvalandırılarak oluşturulurlar. Ebeveyn küme olarak adlandırılan her bir üst küme, çocuk küme olarak adlandırılan alt kümeleri içerir. Ebeveyn kümeler böylece çocuk kümelerdeki elemanları da kapsarlar. Küme oluşumunda böylece alt-üst ilişkisine dayanan bir hiyerarşi söz konusu olduğundan bu tür kümeleme yöntemleri *hiyerarşik kümeleme* olarak adlandırılırlar. Hiyerarşik kümeleme yöntemleri kümeleri ve aynı zamanda hiyerarşiyi ortaya koyarlar.

Kümeleme yöntemleri öğeler arasındaki uzaklıkları kullanarak birbirine en yakın olanları bir araya getirerek bir küme ya da grup içine almayı amaçlarlar. Hiyerarşik kümeleme yöntemleri uzaklıkları kullanmalarına göre farklılık gösterirler:

- En küçük uzaklık tabanlı algoritmalar
- En küçük ortalama uzaklık tabanlı algoritmalar
- En küçük uzaklık kareleri toplamı tabanlı algoritmalar

Literatürde çeşitli hiyerarşik kümeleme yöntemleri ve teknikleri önerilmiş olup bunları aşağıdaki şekilde sınıflandırmak olasıdır:

- *Birleştirici yöntemler*
 - *Bağlantı yöntemleri*
 - *Tek bağlantı yöntemi*
 - *Tam bağlantı yöntemi*
 - *Ortalama bağlantı yöntemi*
 - *Tartısız ortalama bağlantı yöntemi (McQuitty yöntemi)*
 - *Tartılı ortalama bağlantı yöntemi*
 - *Merkezci yöntemler*
 - *Kitle merkezi (centroid) yöntemi*
 - *Ortanca (median) yöntemi (Gower yöntemi)*
 - *Varyans tabanlı yöntemler*
 - *Ward yöntemi*
 - *Ward II yöntemi*
 - *Esnek yöntemler*
 - *Esnek beta yöntemi*
 - *Genelleştirilmiş yöntem (Lance-Williams yöntemi)*

- *Ayırıcı yntemler*
 - *Monotetik ayırıcı yntemler*
 - *MONA algoritması*
 - *Politetik ayırıcı yntemler*
 - *DIANA algoritması*
- *İleri yntemler*
 - *BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)*
 - *CURE (Clustering Using Representative)*
 - *CHAMALEON (A Hierarchical Clustering Algorithm Using Dynamic Modeling)*
 - *OPTICS (Ordering Points to Identify the Clustering Structure)*

Kitabımızı Trkiye'nin İnternet Kitapçısı

TDK Bilim

www.tdk.com.tr



zerinden temin edebilirsiniz.

