

Kümeleme Analizi

Muratberk Ada, Fikret Altunay, Merve Civelek, Sevecen Kaplan, Pınar Koç
Danışman: Dr. Kumru Didem Atalay

ÖZET

Kümeleme analizi grupları kesin olarak bilinmeyen, birimleri, değişkenleri birbiriyle benzer alt kümelere (grup, sınıf) ayırmaya yardımcı olan çok değişkenli istatistiksel analiz yöntemlerinden biridir. Kümeleme analizinin temel amacı birimleri sahip oldukları karakteristik özellikleri temel alarak gruplandırmaktır. Kümeleme analizi son yıllarda gündemde olan analiz yöntemlerinden biridir. Bu yöntem özellikle bilim ve iş alanında birçok durumda uygulanabilen, en kolay yorumlanabilen ve en etkili olan yöntem olma özelliğini taşır. Bu nedenle hemen hemen tüm bilim alanlarında bu yöntemden yararlanılmaktadır.

Son yıllarda tıp alanında klinik olayların zenginliğine bağlı olarak, karmaşık vakaların incelenmesinde daha etkin ve yorumlanması kolay bir yöntem olan kümeleme analizi sıklıkla kullanılmaktadır. Bu çalışmada kümeleme analizi ayrıntılı bir şekilde açıklanarak, tıp alanında bu analizi kullanarak yapılan çalışmalar incelenmiş ve buna ilişkin örnekler verilmiştir.

Anahtar kelimeler: Çok değişkenli istatistiksel analiz, benzerlik ölçütleri, farklılık ölçütleri, ağaç diyagramı.

GİRİŞ

Kümeleme analizi farklı yapıdaki verilerin küme yapısını ve küme sayısını araştırır. Bu analiz gruplama yapısını bulurken küme içindeki gözlemlerin aynı yapıda, kümeler arasındaki gözlemlerin ise farklı yapılarda olmasını amaçlar. Gözlemler için bu ayrıştırmalar benzerlik ve farklılık ölçütleri kullanılarak yapılır. Bu ölçütler genellikle uzaklık ölçütlerine dayandırılarak bulunur. Bazı durumlarda korelasyon da kullanılır. Kümeler için sınıflandırma yapılırken gözlemlere ait grafiklerden de yararlanılabilir. İki değişken olması durumunda gözlemlere ait saçılım grafiği oluşturulabilir. İki değişkenden fazla veri olması durumunda ise temel bileşenler analizi kullanılarak iki değişkene indirgenerek grafik çizilebilir.

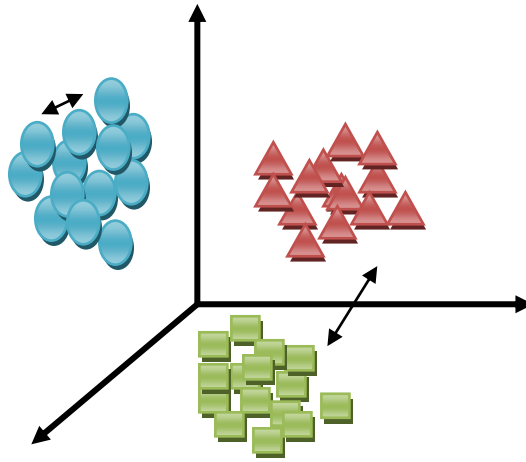
Kümeleme analizi, hemen hemen tüm bilim alanlarında kullanılan bir yöntemdir. Tıp, biyoloji, psikoloji, sosyoloji, arkeoloji gibi belirsizlik koşullarının ve karmaşık oluşumların bulunduğu bilim alanlarında ise daha yoğun olarak kullanılan bir yöntemdir. Örneğin, tıp alanında; hastalıkların sınıflandırılması, psikiyatride; paranoya, şizofreni gibi semptomların doğru sınıflandırılması (teşhis profilleri oluşturması), laboratuvar bulguları ile klinik bulguların oluşturduğu veri matrislerinden hastalık alt gruplamalarının ya da yeni semptomların tanımlanması gibi amaçlarla kümeleme analizinden yararlanılmaktadır (12).

KÜMELEME ANALİZİNİN TANIMI

Kümeleme analizi, gruplanmamış bir veri topluluğunu benzerliklerine göre sınıflamaktadır. Bu sınıflama sayesinde veriyi araştırmacıya uygun, işe yarar özetleyici

bilgiler biçimine dönüştürür. Daha açık bir ifadeyle, doğal grupları kesin olarak bilinmeyen, birimleri değişkenleri ya da birim ve değişkenleri birbirleri ile benzer alt kümelere ayırmaya yardımcı olan yöntemler topluluğudur (12, 13). Kümeleme analizi birim ve bu birimlere ait değişkenlerin sınıflamaları hakkında kesin bilginin bulunmadığı bir popülasyondan alınan n tane birimin, p tane değişkene ilişkin gözlem sonuçları ile ilgilenir. Kümelemede homojen nesnelere birbiri ile birleştirilerek heterojen gruplar oluşturulur ve birimler hiyerarşik bir düzene sokulur. Sınıflandırma yapmak gözlem sonuçlarının çok az bir kayıpla bir araya toplanmasını sağlar (8).

Kümeleme analizi, temel amacı nesnelere (birimleri) sahip oldukları karakteristik özellikleri temel alarak gruplamak olan çok değişkenli teknikler grubudur. Kümeleme analizi, nesnelere küme içerisinde çok benzer biçimde, kümeler arasında farklı olacak biçimde kümeler. Kümeleme işlemi başarılı olursa, bir geometrik çizim yapıldığında nesnelere küme içerisinde birbirine çok yakın, kümeler ise birbirinden uzak olacaktır (5).



Şekil 1. Küme içi ve kümeler arası uzaklıklar

Aynı gruptaki verilerin benzerlikleri, farklı gruptaki verilerin ayrılıkları ne kadar başarılı bir kümeleme yapıldığının ölçütüdür. Gruplar istatistiksel olarak incelendiğinde tamamen birbirlerinden farklı ve kendi içlerinde eş dağılmış olarak oluşturulur. Kümeleme analizinde bireylerin veya nesnelere benzerlik ve farklılık oranlarına göre belirlenmiş gruplar veya sınıflar yer alır (3, 10).

Kümeleme analizinde verilerin normal dağılımlı olması gerektiği varsayımı olmakla birlikte normallik varsayımı prensipte kalmakta, uzaklık değerlerinin normalliği yeterli görülmektedir. Ayrıca kümeleme analizinde kovaryans matrisine ilişkin herhangi bir varsayım bulunmamaktadır (13).

Kümeleme analizi, temel olarak dört değişik amaç için uygulanır. Bu amaçlar aşağıdaki gibi sıralanabilir:

- a) n sayıda birimi, nesneyi, oluşumu, p değişkene göre saptanan özelliklerine göre olabildiğince kendi içinde türdeş ve kendi aralarında farklı alt gruplara ayırmak,
- b) p sayıda değişkeni n sayıda birimde saptanan değerlere göre ortak özellikleri açıkladığı varsayılan alt kümeler ayırmak ve ortak faktör yapıları ortaya koymak,

- c) Hem birimleri hem de deęişkenleri birlikte ele alarak ortak n birime p deęişkene gören ortak özellikli alt kümelere ayırmak,
- d) Birimleri, p deęişkene göre saptanan deęerlere göre, izledikleri biyolojik ve tipolojik sınıflamayı ortaya koymak.

Kümeleme analizinin aşamaları

1. Birimler arasında var olan benzerliğin belirlenebilmesi için kullanılacak ölçülerin ve deęişkenlerin belirlenmesi,
2. Birimler arasındaki benzerliklerin belirlenmesinden sonra birimlerin kümelenmesi,
3. Oluşturulan kümelerin uygun olup olmadığının belirlenmesi,
4. Kümelerin uygun olarak elde edildiği varsayımı altında bunun istatistik geçerliliğinin ortaya konması,

şeklinde sıralanabilir (2).

Kümeleme analizinin pratik problemler için birçok uygulaması mevcuttur. Bu uygulamalar kullanım amaçlarına göre iki grupta incelenebilir. Bunlar anlama için kümeleme ve fayda için kümelemedir. Anlama için kümelemede oluşturulan sınıflar ya da nesne grupları konuyu daha iyi anlamayı, olayı bütün hatlarıyla kavramayı kolaylaştırır. Fayda için kümeleme ise nesne hakkında elde edilen karakteristik bilgilerle nesnenin ait olduğu kümenin özellikleri hakkında bilgi edinmeyi sağlar (10).

KÜMELEME ANALİZİNDE BENZERLİK VE FARKLILIK ÖLÇÜLERİ

Kümeleme aynı küme içerisindeki gözlemlerin birbirine benzer, diğer kümelerdeki gözlemlerden farklı olacak şekilde yapılmasıdır. Bu amaç için benzerlik ve farklılık kavramları kullanılır. Benzerlik iki nesne veya iki özellik arasındaki ilişkinin kuvveti olarak açıklanır. Bu nicel deęer alınan ölçüye veya veri tipine göre deęişik yollardan elde edilir. Farklılık ise, iki nesne arasındaki zıtlık ya da uyumsuzluğun bir ölçüsü olan farklılıkları ölçer. Benzerlik ve farklılık ölçümleri gözlemlerin birbirinden ayırt edilmesini sağlar ve bu sayede gözlemler gruplara ayrılır (3, 10).

Deęişken tipleri kesikli ve sürekli olmak üzere iki kategoride sınıflandırılabilir. Deęişkenin aldığı deęerlerin sayısı sonlu veya sayılabilir sonsuzlukta ise bu deęişkene kesiklidir denir, eđer deęişken birden çok aralıkta her deęeri alabiliyorsa bu deęişken sürekli dir.

Kümeleme analizinde deęişkenlerin ölçek türleri büyük önem taşır. Stevens ölçüm düzeylerini isimsel, sıralı, aralık ve oransal olmak üzere dört sınıfa ayırmıştır (10).

Birimlerin deęişkenlere göre birbirleri arasındaki uzaklıkları hesaplamak amacıyla çeşitli uzaklık ölçü birimleri ileri sürülmüştür. Bu ölçü birimleri veri matrisinde yer alan deęişkenlerin ölçü birimlerine göre farklılık gösterir. Eđer deęişkenler oransal ya da aralıklı ölçüyle elde edilmiş deęerler ise uzaklık ya da ilişki türü ölçülerden yararlanır. Eđer ikili gözlemlere göre ölçümler yapılmış ise birimler arasındaki benzerlik ve farklılık ölçülerinden yararlanır.

Birimlerin birbirleri ile olan benzerlik düzeyleri benzerlik (similarity, sim) matrisi ile gösterilir. sim matrisinin elemanları sim_{ij} ile gösterilir ve $sim_{ij} = 100(1 - d_{ij} / \max(d_{ij}))$ biçiminde hesaplanır. Birimlerin birbirinden farklılıkları (dissimilarity, diss) sim matrisinden yararlanılarak hesaplanır. Diss matrisinin elemanları $diss_{ij}$ ile gösterilir ve $diss_{ij} = 100 - sim_{ij}$ biçiminde hesaplanır (12).

Kümeleme analizinde birimler arasındaki uzaklıkların hesaplanmasında sıklıkla kullanılan ölçüler aşağıdaki gibi verilebilir:

Sürekli veriler için:

Her biri p tane sürekli (oransal ve aralıklı ölçekli) değişken içeren x_i ve x_j gözlem çifti arasındaki uzaklık $d_{ij} = d(x_i, x_j)$ olsun.

a) Minkowski Uzaklığı:

$$d_{\lambda}(x_i, x_j) = \left[\sum_{k=1}^p |x_{ik} - x_{jk}|^{\lambda} \right]^{1/\lambda}; \lambda \geq 1 \text{ için.}$$

b) Manhattan City-Block Uzaklığı ($\lambda = 1$ durumu):

$$d_1(x_i, x_j) = \sum_{k=1}^p |x_{ik} - x_{jk}|$$

c) Öklid Uzaklığı ($\lambda = 2$ durumu):

$$d_2(x_i, x_j) = \left[\sum_{k=1}^p |x_{ik} - x_{jk}|^2 \right]^{1/2}$$

d) Ölçekli Öklid Uzaklığı:

$$d_2(x_i, x_j) = \left[\sum_{k=1}^p w_k^2 |x_{ik} - x_{jk}|^2 \right]^{1/2}$$

Burada w_k k'inci değişkenin standart sapma değerinin (s_k) ya da dağılım aralığı değerinin tersidir. w_k 'nin standart sapma değerinin tersi olması durumunda bu uzaklığa Karl-Pearson uzaklığı da denir.

e) Mahalanobis Uzaklığı:

$$d(x_i, x_j) = (x_i - x_j)^T S^{-1} (x_i - x_j)$$

Burada S örnek ya da küme içi kovaryans matrisidir.

f) Hotelling T² Uzaklığı:

İki grup ya da kümenin ortalama vektörlerinin karşılaştırılmasında kullanılır.

$$d(x_i, x_j) = T^2 = \frac{n_1 n_2}{n} (\bar{x}_i - \bar{x}_j)^T S^{-1} (\bar{x}_i - \bar{x}_j)$$

g) Pearson Korelasyon Katsayısı:

Gözlemler arasındaki uzaklık Pearson korelasyon katsayısı kullanılarak da hesaplanabilir. Pearson korelasyon katsayısı

$$r_{ij} = \frac{\sum_{k=1}^p (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^p (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^p (x_{jk} - \bar{x}_j)^2}}$$

biçiminde tanımlanır. Burada \bar{x}_i i'inci gözlem üzerinden ölçülen tüm p değişken değerlerinin ortalaması olup

$$\bar{x}_i = \frac{1}{p} \sum_{k=1}^p x_{ik}$$

biçiminde hesaplanır. Korelasyon katsayısı kullanılarak iki gözlem vektörü arasındaki uzaklık:

$$d(x_i, x_j) = (1 - r_{ij})$$

şeklinde bulunur (13, 10).

Karışık veriler için:

Değişkenlerin kesikli (isimsel veya sıralı ölçekli) olması durumunda sürekli veriler için verilen formüllere aşağıda verilen katsayı eklenmektedir (13):

$$w_k = \begin{cases} 1, & \text{nicel veriler için} \\ \frac{1}{k' \text{inci de\u0131işkenin da\u011fılım aralığı}}, & \text{nitel veriler için} \end{cases}$$

KÜMELEME YÖNTEMLERİ

Kümeleme yöntemleri iki ana grupta incelenebilir. Bunlar hiyerarşik kümeleme ve hiyerarşik olmayan kümelemedir. En çok kullanılan yöntemler hiyerarşik kümeleme yöntemleridir.

Hiyerarşik kümeleme

Hiyerarşik kümeleme yöntemleri birimleri birbirleri ile değişik aşamalarda bir araya getirerek ardışık biçimde kümeler belirlemeyi ve bu kümelere girecek elemanların hangi uzaklık (ya da benzerlik) düzeyinde küme elemanı olduğunu belirlemeye yarayan

yöntemdir. Hiyerarşik kümeleme iki grupta incelenebilir, bunlar yığılmalı hiyerarşik kümeleme ve bölünmeli hiyerarşik kümelemedir. Yığılmalı hiyerarşik kümeleme verideki her bir gözlemi bir küme olarak düşünür. Birleştirme işlemleri uygulanarak kümeler tek bir küme elde edilinceye kadar devam ettirilir. Bölünmeli hiyerarşik kümelemede, başlangıçta tüm birimlerin bir küme oluşturduğu kabul edilerek, birimleri aşamalı olarak kümelere ayırır. Yığılmalı hiyerarşik kümeleme yöntemleri aşağıdaki gibi verilebilir(6,8,12,13):

- a) Tek bağlantı yöntemi: En yakın komşuluk tekniği olarak da bilinir. En kısa mesafe esasına dayanır. Bu yöntemde uzaklıklar matrisi kullanılarak birbirine en yakın kümeler birleştirilmek suretiyle birleştirmeler art arda tekrarlanmaktadır.
- b) Tam bağlantı yöntemi: En uzak komşuluk tekniği olarak da bilinir. Tek bağlantı yöntemine çok benzer ancak burada iki küme arasında uzaklık olarak her kümedeki eleman çiftlerinin arasındaki uzaklığın en büyüğü alınır.
- c) Ortalama bağlantı yöntemi: Bu yöntemde ayrı gruplarda yer alan gözlem çiftleri arasındaki ortalama uzaklık iki küme arasındaki uzaklık olarak alınır.
- d) Ağırlıklı ortalama bağlantı yöntemi: Bu teknikte ortalama bağlantı tekniğinden farklı olarak, yeni oluşan küme ile diğer kümeler arasındaki uzaklık her bir kümedeki küme sayısı ile ağırlıklandırılır.
- e) Merkezi bağlantı yöntemi: İki küme arasındaki uzaklık kümelerin kendi merkezleri arasındaki uzaklık olarak alınır.
- f) Medyan bağlantı yöntemi: İki küme arasındaki uzaklık, iki kümenin merkezleri arasındaki uzaklığın eşit ağırlıklı olarak hesaplanmasıyla elde edilir.
- g) Ward's bağlantı yöntemi: Bir kümenin ortasına düşen gözlemin, aynı kümenin içinde bulunan gözlemlerden ortalama uzaklığını esas alır.

Hiyerarşik kümeleme tekniklerinde kümeler art arda birleştirilir ve bir grup diğeri ile bir kez birleştirildikten sonra, devam eden adımlarda bir daha ayrılmaz. Bu teknikler ele alınan değişkenler için hiyerarşik bir yapı oluştururlar. Hiyerarşik kümeleme tekniklerinde küme sayısına görsel olarak karar verilir. Bu durumda genellikle *dendogram* olarak bilinen *ağaç diyagramı* kullanılır (4).

Hiyerarşik olmayan kümeleme

Küme sayısı konusunda ön bilgi var ise ya da araştırmacı anlamlı olacak küme sayısına karar vermiş ise bu durumda hiyerarşik olmayan kümeleme yöntemi kullanılabilir. Bu kümeleme yönteminde birimlerin kümelere parçalanması rastgele yapılabilir. Birimlerin ayrılacakları küme sayısı belirlendikten sonra, küme belirleme kriterine göre birimlerin hangi kümelere gireceklerine karar verilir ve atama işlemleri yapılır. Hiyerarşik olmayan kümeleme yöntemleri aşağıdaki gibi verilebilir (4,13):

- a) k-ortalama tekniği: Mac Queen tarafından geliştirilmiştir. Bu yöntemde önce araştırmacının ön bilgisine ve tecrübesine dayanarak küme sayısı belirlenir. Sonra her kümenin tipik bir gözlemi seçilir, benzer gözlemler tipik gözlemin etrafında birer birer kümelendirilir. Burada bazı istatistiksel testler kullanılarak her kümeyi oluşturan gözlemlerin değişkenlere göre ortalamalarına bakılır. Güvenilir olması en belirgin üstünlüğüdür. Buna karşılık yorumlaması zordur.

- b) En çok olabilirlik tekniđi: Her bir gözlem (birey) en büyük olabilirlik değeri verecek biçimde daha önceden belirlenen kümelere atanmaktadır. Bu yöntem kuramsal dayanađı güçlü bir yöntemdir.

KÜMELEME ANALİZİNİN TIP ALANINDA KULLANIMINA ÖRNEKLER

Her alanda olduđu gibi kümeleme analizi tıp alanında da yoğun bir şekilde kullanılmaktadır. Özellikle son yıllarda klinik olayların zenginleşmesi ve karmaşık vakaların incelenmesine bađlı olarak, daha etkin ve yorumlanması kolay bir yöntem olan kümeleme analizine sıklıkla başvurulmaktadır. Tıp alanında bu analizin kullanıldıđı çalışmalara ilişkin örnekler aşıđıdaki gibi verilebilir.

Kanser hastalıđına yakalanan bireylerin bu hastalıkla nasıl başa çıktıkları araştırılmıştır. Kanserle mücadele çok boyutlu olduđu için, konuyu daha net araştırmak amacıyla kümeleme analizinden yararlanılmıştır. Kanserle başa çıkma yolları yapılan kümeleme analizi sonucunda 4 gruba ayrılmıştır. Bu araştırmada kümeleme analizi, kanser tedavisi sırasında ve sonrasında hastaların yaşam kalitesinin nasıl etkilediđini daha kolay incelemek için gruplar elde edilmesini sađlamıştır (11).

Kızılderili toplumundan alınan örnekte bulunan çeşitli HIV risk gruplarının belirlenmesinde kümeleme analizinin yararlılıđı gösterilmiştir. Belirli sayıda Kızılderili'den oluşan bir grup içindeki HIV riski/ korunmasını 4 kümeye ayrılmasında bu yöntem kullanılmıştır. Bu analizde kümeleme analizi yöntemlerinden hiyerarşik yöntem ve k-ortalama yöntemi kullanılmıştır (9).

Kalp seslerinin farklı morfolojilerini belirlemede ve bunların fizyolojik durumlarının sınıflandırılmasında kümeleme analizi yöntemlerinden hiyerarşik yöntem kullanılmıştır (1).

Hastalıklar belli bir bölgede veya belli bir zaman diliminde yığılma gösterebilir. Radyasyon Onkolojisi Anabilim Dalında tedavi gören nazofarinks kanseri tanısı ile tedavi gören hastalara, hastalıđın yere göre sınıflandırılması amacıyla kümeleme analizi uygulanmıştır (7).

SONUÇ

Bilimsel araştırmalarda, araştırmaya konu olan olaylar veya nesnelere her birey için aynı anda ölçülebilen bir veya birden çok deđişken tarafından etkilenebilir. Birden çok deđişkene sahip araştırmalarda deđişkenleri ayrı ayrı incelemek ve analiz etmek gerçekte olan durumu yansıtmayabilir. Deđişkenleri ayrı ayrı incelemek aralarındaki ilişkiyi göz ardı etmek demektir. İlişki terimi bađımlılıđı beraberinde getirir. Bu nedenle çok deđişkenli istatistiksel analiz yöntemleri geliştirilmiştir. Bu analiz yöntemlerinden biri olan kümeleme analizi karşılıklı bađıntıları göz önüne alır.

Kümeleme analizi ile bireylerin sınıflandırılması ayrıntılı bir şekilde açıklanır. Gruplanan kümeler sayesinde, kalabalık ve karmaşık olan veri topluluđundan uzaklaşarak daha rahat incelenebilir ve analiz edilebilir bir problem elde edilir. Karmaşık bir veri topluluđundan kurtulmak tıp alanındaki çalışmalarda da çok gereklidir. Kümeleme analizi yardımıyla yapılan gruplamalarla bilginin düzenlenmesi ve sonuçların elde edilmesi daha

etkin sađlanır. Bu alıřmada veriyi basitleřtirerek daha etkin deđiřkenler üzerinde alıřma imkanı sađlayan kmeleme analizi incelenmiř ve tıp alanındaki uygulamalarına deđinilmiřtir.

KAYNAKLAR

1. Amit G., Gavriely, N., Intrator, Cluster Analysis and Classification of Heart Sounds. Biomedical Signal Processing and Control. 2009; 4:26-36
2. Blbl, S., Gler, M.F., Kandemir, A.ř., Propensity Skor Uygulamalarında Kmeleme Analizinin Test Amalı Kullanımı.
3. Dođan, i., Kmeleme Analizi ile Seleksiyon. Turk J Vet Anim Sci. 2002;26: 47-53
4. Everitt, B., Dunn, G.: Applied Multivariate Data Analysis, Oxford Uni. Press, New York, 1992:101.
5. Hair, Jr. F.J., Anderson, E. R., Tatham, L. R., et al.: Multivariate Data Analysis With Readings, 5. Ed., Prentice-Hall, USA, 1998.
6. Kalaycı, ř.: SPSS Uygulamalı ok Deđiřkenli İstatistik Teknikleri, 4. Baskı, Asil Yayın Dađıtım, 2009.
7. Karabulut, E., Alpar, R., zyar, E., Hastalıkların Yere Gre Kmelenmesinde Kullanılan Yntemler. İnn niversitesi Tıp Fakltesi Dergisi. 2006; 13(1): 37-43
8. Lorr, M.: Cluster Analysis for Social Scientists, Jossey-Bass Publishers, London, 1983.
9. Mitchell, C.H., Kaufman C.E., Beals, J., et al. Identifying Diverse HIV Groups Among American Indian Young Adults: The Utility of Cluster Analysis. AIDS and Behavior. 2004; 8(3):263-275
10. Servi, T., ok Deđiřkenli Karma Dađılım Modeline Dayalı Kmeleme Analizi, ukurova niversitesi Fen Bilimleri Enstits, Doktora Tezi, Adana, 2009.
11. Shapiro, D.E., Rodricue, J.R., Boggs, S.R., et al., Cluster Analysis of The Medical Coping Modes Questionnaire: Evidence for Coping with Cancer Styles? Journal of Psvchosomeric Research. 1994; 38(2):151-159
12. zdamar, K.: Paket Programlarla İstatistiksel Veri Analizi 2, Kaan Kitabevi, Eskiřehir, 1999.
13. Tatlıdil, H.: Uygulamalı ok Deđiřkenli İstatistiksel Analiz, Cem Web Ofset, Ankara, 1996.