

## n < p Boyutlu Biyolojik Verilerde Farklı Kümeleme Yöntemlerinin Karşılaştırılması Olarak İncelenmesi\*

Özkan ÖZTÜRK<sup>1</sup>, Necati YILDIZ<sup>2</sup>

<sup>1</sup> HÜ, Teknik Bilimler Meslek Yüksek Okulu, Van

<sup>2</sup> Bingöl Üniv. Ziraat Fak. Zootekni Bölümü, Bingöl

**Geliş Tarihi (Received): 22.08.2012**

**Kabul Tarihi (Accepted): 22.09.2012**

**Özet:** Araştırma, n < p boyutlu olan, 24 farklı Antepfıstığı (Pistacia vera L.) tipi ve bu tiplerden alınan 38 adet veri üzerinde yapıldı. Ancak bu tip bir veri matrisine kümeleme analizine ilişkin bazı çok değişkenli test istatistiklerinin uygulanabilmesi için değişken sayısının (n ≥ p) azaltılması gerekmektedir. Değişken sayısının azaltılmasında Temel Bileşenler (Principal Component) analizi, Ayırma (Diskriminant) analizi ve Korelasyon analizinden yararlanılmıştır. Söz konusu edilen yöntemlerle indirgenen değişkenler, farklı kümeleme yöntemleriyle karşılaştırılarak incelenmiştir. Sonuçta, kümelemede en uygun yöntemin Temel Bileşenler analizi ile birlikte kullanılan Ward metodunun olduğu saptandı. Küme sayısının belirlenmesinde ise en uygun ölçütün C<sub>max</sub>, Wilks Lambda ve Hotelling Lawley Trace istatistiklerinin olduğu belirlenmiştir.

**Anahtar Kelimeler:** Kümeleme Analizi, Sayısal Sınıflandırma, Temel Bileşenler Analizi, Ayırma Analizi, Korelasyon Analizi

### A Comparison of Different Clustering Methods on the Biological Data with n < p Dimensions

**Abstract :** This research was carried out on 38 variables from 24 types of pistachio (Pistacia vera L.) with n < p dimensions. However, in order to apply the some multivariate test statistics of clustering analysis to this type data matrix, the number of variables (n ≥ p) must be decreased. Principal Component, Discriminant and analysis of correlation were used to decrease the number of variables. The reduced number of variables by using these mentioned methods, was comparatively evaluated by using different clustering methods. Present results shown that the most suitable method for clustering is the Ward method used together with Principle Components variable reducing technique. It was also found out that the most suitable measurement index to determine the number of clusters were Wilk's Lambda, C<sub>max</sub> and Hotelling Lawley Trace statistics.

**Keywords:** Cluster Analysis, Numerical Taxonomy, Principal Component Analysis, Discriminant Analysis, Correlation Analysis

### GİRİŞ

Türlerin alt bölümleri arasındaki ekolojik ve genetik etkilerin tümünün olduğu türde varyasyon, türlere ait populasyon yapılarını ortaya koymaktadır (Jain, 1975).

Morfolojik varyasyonların ıslah çalışmalarında saptanması, değişik tür ve çeşitler arasındaki taksonomik ilişkileri aydınlatmak bakımından büyük bir önem arz etmektedir. Çünkü yetiştirilen türler içerisinde bulunan varyasyonların bilinmesi ve bu varyasyonun dağılımı durumu, ıslah programlarının uygulanması açısından çok önemlidir (Bliss, 1981).

Türler arasındaki varyasyonun saptanmasında sayısal sınıflandırma (Numerical Taxonomy) diğer bir ifade ile kümeleme analizi yöntemleri önemli rol oynamaktadır. Bu yöntemler ilk defa 1957'de Sneath tarafından organizmalardaki varyasyonu saptamak ve daha somut bir sınıflama yapmak amacıyla kullanılmıştır. Bitkiler ve hayvanların doğru ve nicel olarak sınıflandırılmaları için uygun istatistik analizlerin kullanılması gerekmektedir.

Özellikle gen bankalarında korumaya alınan materyalin benzer ve benzer olmayan özellikleri ile ilgili kayıtlarını analiz etmek kolay değildir. Çünkü benzerlik ve farklılık saptanmasında genetik, agronomik, morfolojik, kimyasal vb. pek çok özellik dikkate alınmaktadır. Çok fazla sayıda özelliğin incelendiği durumlarda, her bir özelliğin ölçülebilir verilerle ifadesi, bu ifadelerin benzerlik veya farklılıkını ortaya koyacak bir dizi hesaplamalarının yapılmasını gerektirmektedir (Tan, 1983).

İncelenen özellik sayısı arttıkça, benzerlik veya ayrılık ortaya koyan indekslerin hesaplamaları karmaşık ve zaman alıcı bir hal almaktadır. Çünkü klasik istatistiksel yöntemler, karakter adedi ve karşılaştırılan örnek sayısı arttıkça yetersiz kalmaktadır. Bu sorun, fazla sayıdaki karakteri sınıflandırma indeksleri kapsamında formüle edilen kümeleme analizinin ortaya konması ile büyük ölçüde çözüme kavuşturulmuştur.

\* Doktora tezinden alınmıştır.

\* Sorumlu yazar: Öztürk, Ö., ozirfan23@yahoo.com

Kümeleme analizi sadece hayvan ve bitki genetiği kaynaklı materyalin değerlendirilmesinde kullanılmakla kalmayıp, ıslah çalışmaları da genetik varyasyonu ortaya koyma, hat seçimi, ırk saptama, hibrit indeksleme, konukçu - patojen - vektör ilişkilerini belirleme, toprak sınıflamaları için saptanan özelliklerden yararlanarak toprak haritaları yapmak amacıyla da kullanılmaktadır (Tan, 1983).

Günümüzde Biyoloji, Ziraat ve Tıp alanında yapılan araştırmalardan elde edilen birçok özelliği kapsayan karmaşık veriler, çok değişkenli analiz yöntemleri ile değerlendirilmektedir (Johnson and Wichern, 1992).

Çok değişkenli analiz yöntemlerinden en yaygın olarak kullanılanları; Çok Değişkenli Varyans analizi (Multivariate Analysis of Variance), Faktör analizi (Factor Analysis), Temel Bileşenler analizi (Principal Component Analysis), Ayırım analizi (Discriminant Analysis), Uyum analizi (Correspondence Analysis), Çok boyutlu ölçekleme (Multidimensional Scaling) ve Kümeleme analizi (Cluster Analysis) olarak sıralanabilir (Anderberg, 1973; Özdamar, 1999).

Veri matrisindeki birim ve/veya değişkenler hakkında ayrıntılı şekilde bilgi edinilmek isteniyorsa, bunların en çok benzer olanları bir araya getirilerek sınıflandırılırlar. Sınıflandırma sonucunda aynı grup içindeki birimlerin (bireyler) benzerliği fazla iken diğer gruplardakilerle olan benzerlikleri daha azdır (Anderberg, 1973).

Kümeleme analizinden en yüksek yarar elde etmek için, aşağıdaki hususlara dikkat edilmesi gerektiğini Anderberg, (1973) belirtmiştir;

a) Verilerde küme yapısının olmaması veya tek bir küme yapısının olması gibi özel durumlarla karşılaşılabilir. Bu durumları unutmamak gerekir.

b) Aynı veri setine farklı kümeleme yöntemleri uygulandığında farklı sonuçlar verebilir. Bu sebepten tek bir kümeleme yöntemine bağlı kalmayıp birden fazla yöntemin denenmesi daha güvenilir sonuçlar verecektir.

c) Sınıflandırma sonucunda, karmaşık bir veri yapısı ile karşılaşılabilir.

d) Kümeleme yöntemleri hipotez testleri için amaçtır.

Çalışmanın amacı verileri en uygun kümelere ayırmaktır. Ancak veri setinde bireylere ait çok fazla değişken (p) bulunması; verileri kümelere ayırmada, yorumlamada, elde edilen kümelere çok değişkenli istatistik analizlerinin uygulanmasında, anlamsız değişkenlerin bireyleri yanlış kümelere sınıflandırması gibi sıkıntılar ortaya çıkarabilecektir. Bu nedenle öncelikle değişkenlerin boyutu azaltılarak birimlerin kümelenebilmesine çözüm arandı.

Ayrıca bireylerin çok boyutlu uzaydaki dağılımına bağlı olarak farklı kümeleme yöntemlerine göre farklı sonuçlar elde edilebileceğinden, sınıflandırma yapılırken birden fazla kümeleme yöntemi kullanılarak sonuçlar değerlendirilecektir.

## MATERYAL ve METOT

Bu çalışmanın verileri, anlırfa T GEM Ceylanpınar Tarım İstasyonunda Acar, (1997) tarafından yapılan araştırmadan elde edilmiştir. Ölçümler, 24 farklı erkek Antepfıstığı (Pistacia vera L.) tipleri üzerinde gerçekleştirilmiştir. Fıstık ağaçlarını türlerine ayırmak amacıyla 38 adet farklı karakterden ölçüm alınmıştır.

Ancak bilindiği üzere kümeleme analizi ile birlikte bazı Çok Değişkenli test istatistiklerinin uygulanabilmesi için veri matrisinin en azından ( $n \geq p$ ) boyutunda olması gerekmektedir. Değişken sayısının (p) birim (fıstık tipleri) sayısından (n) fazla olduğu durumlarda, ( $n < p$  durumunda) verileri analiz edebilmek için, değişken sayısının bazı boyut indirgeme yöntemleri ile azaltılması gerekmektedir. Bu nedenle araştırmada Korelasyon analizi, Temel Bileşenler analizi ve Stepwise Diskriminant analizi ile değişkenlerin boyutu ( $n \geq p$ )'ye indirgenerek gerçek değişkenlerle elde edilen kümelemeye en yakın boyut indirgeme yöntemi belirlenmeye çalışıldı. Bununla birlikte yukarıda sözü edilen kümeleme yöntemlerinden en yaygın olarak kullanılanları, değişken boyutu indirgenmiş 24 farklı fıstık ağacı tipi üzerine uygulanarak en iyi küme yapısını ortaya koyan yöntem tespit edilmeye çalışıldı.

Değişken boyutu indirgeme işleminden sonra farklı kümeleme yöntemlerine göre elde edilen kümeler, küme sayısını belirlemede kullanılan  $C_{max}$ ,  $k$ , ve  $M_{min}$  ölçütleri ile Wilk's Lambda ve Hotelling Lawley  $\lambda$  istatistikleri kullanılarak test edilip en iyi küme yapısı belirlenecektir.

Kümeleme analizinde kullanılan yöntemlerin teorik temelleri, çalışmada kolaylık sağlamak amacıyla aşağıda verildi.

Sayısal sınıflandırma analizinde nicel veriler için kullanılan uzaklık ölçütü, Minkowsky ölçütü olup aynı zamanda genel uzaklık ölçüsüdür (Anderberg, 1973).

$X_i$  ve  $X_j$  (p) tane karaktere göre ölçülmüş  $i$ . ve  $j$ . bireyler olarak tanımlandığında, Minkowsky ölçütü;

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

( $i, j = 1, 2, \dots, n$ ) olarak verilebilir.

Bu ölçüt Minkowsky uzaklık ölçütünün,  $\lambda=2$  olduğu durumdaki özel hali Öklid uzaklık ölçütünü verir (Anderberg, 1973; Tatlıdil, 1992).

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

Mariott (1974) ve Kendall (1975)'in Öklid uzaklığına gözlem değerleri standardize edildikten sonra kullanıldıktan, karakterlerin farklı ölçü birimleri ile ölçülmü olmasından etkilenmediğini bildirmektedir (Özdamar, 1988). Bu üstünlüğünden

dolayı, kümeleme analizinde sıklıkla kullanılan uzaklık ölçüsü Öklid oldu, ayrıca gruplar arası hata kareler ortalamasını minimize etti için genellikle tercih edildiğini bildirmiştir (Tatlıdil, 1992; Huang, 2001).

### Kümeleme kriterleri

Test istatistikleri	Wilk's Lambda	Hotelling Lawley z	$C_{max}$	$M_{min}$	k
Formüller	$\Lambda = \frac{ W }{ W+B }$	$T_o^2 = \sum_{r=1}^p r$	$C = \left[ \frac{iz(B)}{k-1} \right] / \left[ \frac{iz(W)}{N-k} \right]$	$M = k^2  W $	$k \cong \left( \frac{n}{2} \right)^{1/2}$
Kritik değer	$t^2_{p(k-1),r}$	$t^2_{p(k-1),r}$			

### Kümeleme Yöntemleri

**1-Tek Bala yöntemi (en yakın komu) kümeleme yöntemi**

$$d(u,v)_w = \min \{duw, dvw\}$$

**2-Tam Bala yöntemi (en uzak komu) kümeleme yöntemi**

$$d(u,v)_w = \max \{duw, dvw\}$$

**3- Gruplar için Ortalama Bala yöntemi kümeleme yöntemi**

$T_u$ ,  $(u)$  kümesi içindeki bireylerin tüm mümkün çiftleri arasındaki benzerliklerin toplamı,  $N_u$ 'da  $(u)$  kümesi içindeki bireylerin sayısı olarak kabul edilecek olursa; kümeler birlemeden önce  $T_u = 0$  ve  $N_u = 1$  olacaktır.  $(u)$  ve  $(v)$  birleştirildiğinde yeni oluşacak  $(uv)$  kümesi içindeki bireylerin tüm mümkün çiftleri arasındaki benzerliklerin toplamı,

$$T(uv) = T_u + T_v + duv \text{ ile}$$

$(uv)$  kümesindeki birimlerin sayısı ise

$$N(uv) = N_u + N_v$$

şeklinde hesaplanır.

$(u)$  ve  $(v)$  kümelerine ait çiftlerin benzerlik ölçülerinden ve birey sayılarından yararlanarak en benzer çiftin bulunmasında grup içi benzerlik ortalaması;

$$\bar{d}_{uv} = \frac{T(uv)}{N(uv) \cdot \frac{(N(uv)-1)}{2}} = \frac{T_u + T_v + duv}{(N_u + N_v) \cdot \frac{(N_u + N_v - 1)}{2}}$$

olarak hesaplanabilir.

Hesaplanan bu ortalamalar minimum oldu unda  $(u)$  ve  $(v)$  kümeleri birleştirilir.  $(u)$  ve  $(v)$ 'nin satır ve sütunları D matrisinden çıkarılır. Bu işlemler tekrarlanarak suretiyle tüm bireylerin kümeleme işlemi devam eder (Anderberg, 1973; Tatlıdil, 1992; Mac Queen, 1967).

**4- Gruplar Arası Ortalama Bala yöntemi kümeleme yöntemi**

$$\bar{d}_{uv} = \frac{duv}{N_u \cdot N_v}$$

olarak hesaplanır.

**5- Merkezi (Centroid) Bala yöntemi kümeleme yöntemi**

$$d(uv)_w = \frac{N_u \cdot duw + N_v \cdot dvw}{N(uv)} - \frac{N_u \cdot N_v \cdot duv}{N(uv)^2}$$

olarak hesaplanır.

**6. Ortanca (Median) Bala yöntemi kümeleme yöntemi**

$$d(uv)_w = \frac{1}{2} \cdot duw + \frac{1}{2} \cdot dvw - \frac{1}{4} \cdot duv$$

olarak hesaplanır.

**7. Ward kümeleme tekniği**

$$d(uv)_w = \frac{\{(N_w + N_u) \cdot duw + (N_w + N_v) \cdot dvw - N_w \cdot duv\}}{N_w + N(uv)}$$

olarak hesaplanır.

**8. McQuitty Bala yöntemi kümeleme yöntemi**

$$d(uv)_w = \frac{d(uw) + d(vw)}{2}$$

olarak hesaplanır.

**9. k-Ortalama Yöntemi:**

$X_1, X_2, X_3, \dots, X_n$  bireyin her birinden elde edilen  $p$  karakterli gözlem vektörleri, çok boyutlu  $X$  uzayında birer nokta olarak düşünülecek olursa ve aynı uzayda, her grup birey için küme merkezleri olarak,  $a_{1n}, a_{2n}, \dots, a_{kn}$ ;  $j = 1, 2, \dots, k$  seçildiğinde;

$$W_n = \frac{1}{n} \sum_{i=1}^n \min_{1 \leq j \leq k} \left\| X_i - a_{jn} \right\|^2$$

kuralına göre bireyler en yakın kümeye sınıflanmaktadır (Hawkins et al., 1982; Minitab, 1995; Blashfield, 1976).

### Kümeleme analizinden önce de i ken sayısının azaltılması

Çalışmamızda de i ken sayısı kümelenecek birim sayısından fazla oldu undan ( $n=24 < p=38$ ), bazı çok de i kenli analizlerin yapılabilmesi için de i ken sayısının azaltılması ( $n \geq p$ ) gerekmektedir.

Orijinal de i kenlerin ölçüm de i kenlerinin, de i ken aralıklarının ve ölçü birimlerinin fazla oldu u durumlarda, de i ken ( $p$ ) sayısının birim ( $n$ ) sayısından çok fazla oldu u ( $n < p$ ) vb. durumlarda korelasyon yada kovaryans matrislerini nonsingular (tekil olmayan matris) hale getirmek için veri indirgemesi yapmak ve kümeleme analizi uygulamak için temel bile en skorları hesaplayarak Temel Bile enler analizinden yararlanılmaktadır (Özdamar, 1999).

Bu durumları göz önüne alarak aşağıda sıralanacak boyut indirgeme yöntemleri kullanıldı:

#### 1. Temel Bile enler analizi (Principal Component Analysis)

Aralarında korelasyon bulunan  $p$  sayıda orijinal de i kenin açıkladığı varyans yapısını daha az sayıda ( $p > k$ ) orijinal de i kenlerin doğrusal bile enleri olan de i kenlerle ifade etme yöntemine Temel Bile enler analizi denilir. Temel bile enler, bizzat kendileri bir sonuç olmaktan ziyade sonuç almayı sağlayıcı özelliğe sahiptir. Çünkü Temel bile enler daha geniş incelemeler için bir ara adım özelliği taşımaktadırlar (Özdamar, 1999; Sharma, 1996).

statistik analizlerde de i kenler arasında önemli düzeyde yüksek birlikte de i ken ya da korelasyonların bulunması arzu edilmediğinden, veri setinin bir ekinde bu de i kenlerden arındırılarak kullanılması daha uygun olacaktır (Özdamar, 1999).

Temel Bile enler analizi ile en etkili yani varyansın büyük bir bölümünü açıklayan Temel Bile enler ait skorlar belirlenerek, kümeleme analizinde boyutu indirgenmiş de i kenler olarak kullanıldı.

#### 2. Korelasyon analizi

De i kenler arasındaki ikili korelasyon katsayıları hesaplanarak, aralarında yüksek derecede ilişkili bulunan de i kenler veri setinden çıkarıldı ve kalan de i kenlerin yardımcı birimler (Fıstık tipleri) sözü edilen kümeleme yöntemleri ile sınıflandırıldı.

#### 3. Adımlı Ayırım (Stepwise Diskriminant) analizi

Veri setlerine diskriminant analizinin uygulanabilmesi için, veri setlerinin MANOVA (çok de i kenli

varyans analizi) uygulaması için gerekli varsayımları taşıması gerekmektedir (Özdamar, 1999). Bu varsayımlar;

a- Veri matrisi çok de i kenli normal dağılımı göstermeli

b- De i kenlerin varyans ve kovaryansları homojen olmalı

c- De i kenlerin ortalamaları ve varyansları arasında korelasyon bulunmamalı

d- Veri seti grupların birbirinden ayrılmasında rol oynamayacak gereksiz de i kenleri içermemelidir.

Söz konusu varsayımlar sağlandı nda; birimlerin belirli kümelere atanmasında en etkili olan de i kenler amaçlı olarak *Adımlı Ayırım* (Stepwise Diskriminant) analizi ile veri setine dâhil edilir. Böylece kümelemede en az etkisi olan de i kenler veri setinden elemine edilmiş olur.

### ARA TIRMA BULGULARI ve TARTI MA

Bu çalışmada, Acar, (1997) tarafından anlırfa Ceylanpınar Tarım İstasyonunda (T GEM) yapılan bir ara tırma sonucunda elde edilen veriler kullanıldı. Ölçümler 24 farklı fıstık tipine ait erkek fıstık açları (*Pistacia vera* L.) üzerinde gerçekleştirildi. Fıstık açlarını türlerine ayırmak amacıyla 38 adet farklı karakter için ölçüm alınmıştır. Bu karakterler (de i kenler) morfolojik ve fenolojik gözlemlere dayalı olarak Acar, (1997) tarafından geleneksel metotlara dayanarak tespit edilmiştir. Fıstık tipleri erkenci ve geççi olmak üzere temel iki gruba ayrılmaktadır.

Bununla birlikte erkenci ve geççi tipler arasında yer alan ve bazen erkenci tiplerle, bazen de geççi tiplerle kümelenen bir kısım melez tipler orta geççi tipler olmak üzere üçüncü bir grup oluşturmuştur. Sözü edilen fıstık tipleri Çizelge 1'de verilmiştir.

Fıstık tiplerine ilişkin de i kenlerin bir kısmı aralık ölçümlerinde (sürekli ve kesikli) de i ken özelliği gösterirken, oranlı özelliklerde vardır. Bunun yanı sıra bazı karakterler adlandırma (Nominal), bazıları da sıralama (ordinal) ölçümlerinde de i ken özellikleri olup, bu tipler de i kenlere puanlar verildi. Kümeleme analizinde kullanılmak üzere her bir karakter (C1,C2,...,C38 olarak) yeniden tanımlanıp, bu karakterlere ilişkin de i ken tiplerini Çizelge 2'de görmek mümkündür.

Daha öncede belirtildiği gibi veriler, farklı ölçüm ölçütüne sahip ise en uygun uzaklık ölçüsü Öklid (Euclid) veya Öklid uzaklığının karesi (kareli Öklid) dir.

Çizelge 1. Kümeleme Analizinde Kullanılan Fıstık Tipleri ve Kod Numaraları

Erkenci Fıstık Tipleri	Orta Geççi Fıstık Tipleri	Geççi Fıstık Tipleri
3	1	7
4	2	13
5	8	14
6	9	17
10	12	18
11	16	19
15	22	20
24	23	21

Çizelge 2. Karakterlere İlişkin Tanımlama ve Değişken Tipleri

Tanımlama	Karakterler	Değişken Tipi
C1	Büyüme Biçimi	Adlandırma
C2	Sürgün Uzunluğu	Sürekli (Oran ölçeri)
C3	Dal Sayısı	Kesikli
C4	Sürgün Rengi	Nominal (Sıralama)
C5	Dal Kalınlığı	Sürekli (Oran ölçeri)
C6	Yaprak Uzunluğu	Sürekli (Oran ölçeri)
C7	Yaprak Geniliği	Sürekli (Oran ölçeri)
C8	Yaprakçık Sayısı	Kesikli
C9	Uç Yaprak Uzunluğu	Sürekli (Oran ölçeri)
C10	Uç Yaprak Geniliği	Sürekli (Oran ölçeri)
C11	Uç Yaprak Rengi	Sıralama
C12	Uç Yaprak Şekli	Adlandırma
C13	Uç Yaprak Tepesi	Adlandırma
C14	Uç Yaprak Tabanı	Adlandırma
C15	Uç Yaprak Kenarı	Adlandırma
C16	Yaprak Rengi	Adlandırma
C17	Karagöz Ağırlığı	Sürekli (Oran ölçeri)
C18	Karagöz Şekli	Adlandırma
C19	Karagöz Rengi	Adlandırma
C20	Sürgündeki Karagöz Sayısı	Kesikli
C21	Karagöz Uzunluğu	Sürekli (Oran ölçeri)
C22	Karagöz Geniliği	Sürekli (Oran ölçeri)
C23	Gözlerin Kabarması	(Aralık ölçeri.)
C24	Gözlerin Patlaması	(Aralık ölçeri)
C25	Çiçeklenme Başlangıcı	(Aralık ölçeri)
C26	Tam Çiçeklenme	(Aralık ölçeri)
C27	Çiçeklenme Sonu	(Aralık ölçeri)
C28	Çiçeklenme Süresi	Kesikli
C29	Salkım Sayısı	Kesikli
C30	Salkım Uzunluğu	Sürekli (Oran ölçeri)
C31	Bir Salkımdaki Çiçek Sayısı	Kesikli
C32	Verim Potansiyeli	Sürekli (Oran ölçeri)
C33	Çiçek Tozu Üretimi	Sürekli (Oran ölçeri)
C34	Salkım Ağırlığı	Sürekli (Oran ölçeri)
C35	Bir Çiçekteki Anter Sayısı	Kesikli
C36	Bir Anterdeki Çiçek Tozu Sayısı	Kesikli
C37	Çiçek Tozu Canlılık Oranı	Orantı
C38	% 20 Sakkarozda Çiçek Tozu Çimlenme Oranı	Orantı

Veri setinde değişken sayısı birim (denek) sayısından fazla olduğundan ( $n < p$ ), bu tip bir veri matrisine kümeleme analizine ilişkin bazı çok değişkenli test istatistiklerinin uygulanabilmesi için değişken sayısı Temel Bileşenler, Adımlı Ayırım (Stepwise Diskriminant) ve Korelasyon analizleri ile ( $n \geq p$ ) azaltıldı.

Değişkenler ( $p$ ) arası korelasyon analizi sonuçlarından yararlanarak, aralarında yüksek derecede (\*\*:  $P < 0.01$ ) ilişkisi bulunan değişkenler tespit edilerek veri setinden çıkarıldı. Bu suretle hem fıstık ağaçlarının sınıflandırılmasında birbirinden bağımsız olan karakterler (değişkenler) tespit edilmiş, hem de veri

matrisinin boyutu ( $n \geq p$ ) indirgenmiştir (Ancak yapılan analizler de korelasyonla elemine edilen değişkenlerle elde edilen kümeleme dağınıkları olumlu sonuçlar vermediğinden, detayının verilmesine gerek görülmedi. Bununla birlikte Stepwise Diskriminant analizi ile yapılan boyut indirgemede de olumlu sonuçlar alınmadığından detayı verilmedi).

Üçüncü yöntem olarak, Temel Bileşenler analizinden yararlanarak  $n < p$  boyutunda olan veri matrisinin boyutu indirgenerek değişken sayısı azaltıldı. Bireylerin kümelenebilmesinde en etkili ve en fazla varyasyonu açıklayan Temel Bileşenler ise Çizelge 3'deki gibi elde edildi.

Çizelge 3. Temel Bileşenler (Principal Component Analysis) Analizi Sonuçları

Bileşenler	Özdeğerler	Varyans	Eklemeli Varyans (Cumulative)
1	9.7703	0.257	0.257
2	7.0284	0.185	0.442
3	4.6561	0.123	0.565
4	2.9303	0.077	0.642
5	2.1037	0.055	0.697
6	1.8353	0.048	0.745
7	0.0440	0.044	0.789
8	0.0380	0.038	0.828
9	0.0370	0.037	0.864
10	0.0270	0.027	0.892
11	0.0230	0.023	0.914
12	0.0200	0.020	0.935
13	0.6365	0.017	0.952
14	0.5900	0.016	0.967
15	0.4076	0.011	0.978
16	0.2923	0.008	0.986
17	0.1746	0.005	0.990
18	0.1576	0.004	0.994
19	0.0890	0.002	0.997
20	0.0595	0.002	0.998
21	0.0428	0.001	0.999
22	0.0272	0.001	1.000
23	0.0000	0.000	1.000
24	0.0000	0.000	1.000
.	.	.	.
.	.	.	.
38	0.0000	0.000	1.000

Analizde kullanılacak en etkili Temel Bileşenlerin sayısını belirlemek amacıyla özdeğerlerin, varyans açıklama oranlarından yararlanarak Temel Bileşenlere (Ana Bileşenlere) ilişkin grafik (Scree Graph) çizdirildi. Böylece, Temel Bileşenlerin sayısı ekil 1'de verilen grafik vasıtasıyla belirleme imkânı doğmuştur.

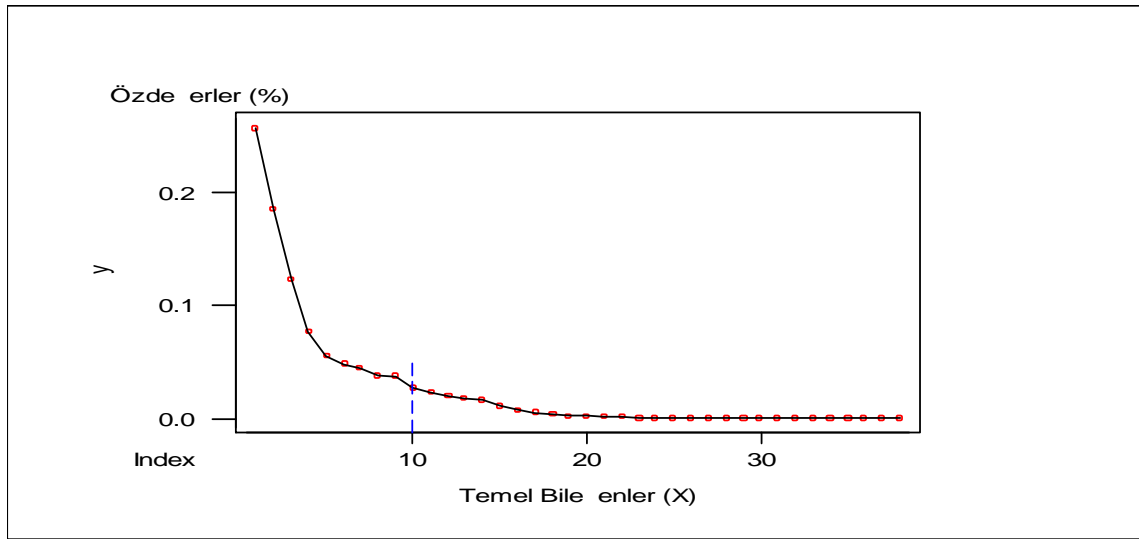
Grafik (Scree graph) vasıtasıyla belirlenen ve varyansın % 89'unu açıklayan ilk on temel bileşene ait skorlar kümeleme analizinde boyutu indirgenmiş de i kenler olarak kullanıldı (Çizelge 4).

Küme yapılarını ve küme sayılarını belirlemeden önce ilk iki temel bileşen vasıtasıyla da birimlerin dağılımını ve kaç küme ayrılacaklarını ekil 2'de çizilmiş olan grafik vasıtasıyla tahmin etmek mümkündür.

Temel Bileşenler, Adımlı Ayırım (Stepwise Diskriminant) ve Korelasyon analizi ile boyutu

indirgenmiş olan de i kenlere yukarıda sözü edilen kümeleme analizi metodları ayrı ayrı uygulandı. Elde edilen küme yapıları incelendi inde bu dağılıma en çok uygunluk gösteren yöntemin, Temel Bileşenler analizi ile boyutu indirgenmiş de i kenlerin Ward metodu ile birlikte kullanımı sonucu olduğu saptandı ( ekil 3). Nedenine gelince, uygulanan tüm kümeleme yöntemlerine ait dendogramlar incelendi inde erkenci ve geçici fıstık tiplerini birbirinden en iyi ayıran yöntemin Ward yöntemi olduğu bunu Tam Balantı ve k-ortalama yöntemi izlemiştir.

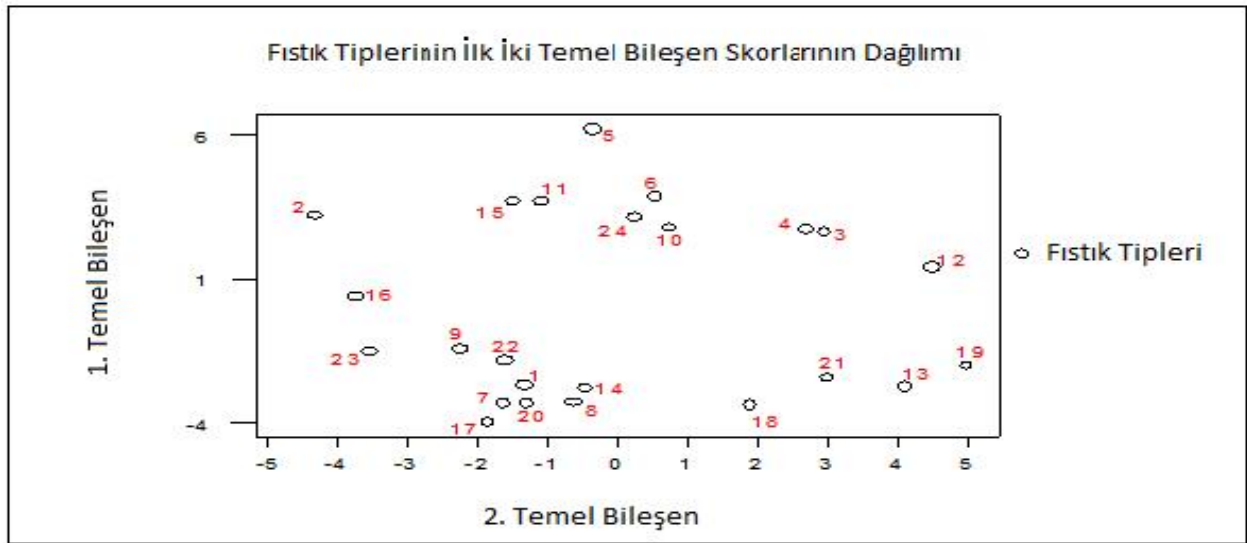
Temel Bileşenler analizi ile boyutu indirgenen Ward ve Tam balantı yöntemi ekil 3 ve ekil 4'de karışıklık incelendi inde erkenci ve geçici çeşitlerin yaklaşık benzer bir şekilde sınıflandırıldığı gözlenmektedir.



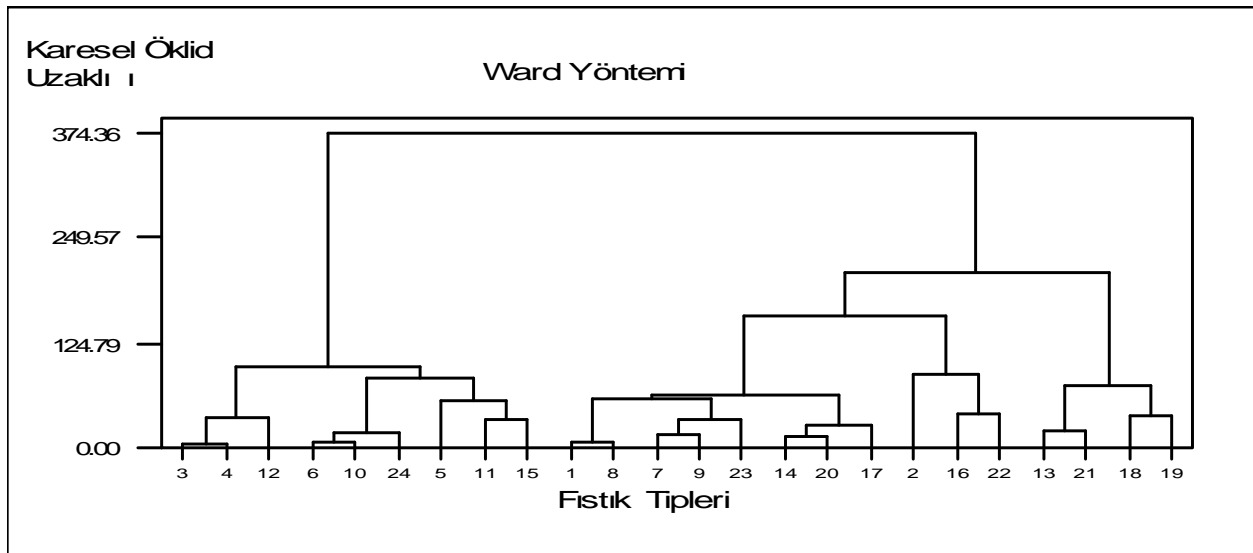
ekil 1. En Etkili Olan Temel Bileşen Sayısının (Scree Graph) Grafik Metodu ile Belirlenmesi

Çizelge 4. Temel Bileşenler Analizi ile Elde Edilen Yeni Değişkenlere İlişkin Scorlar

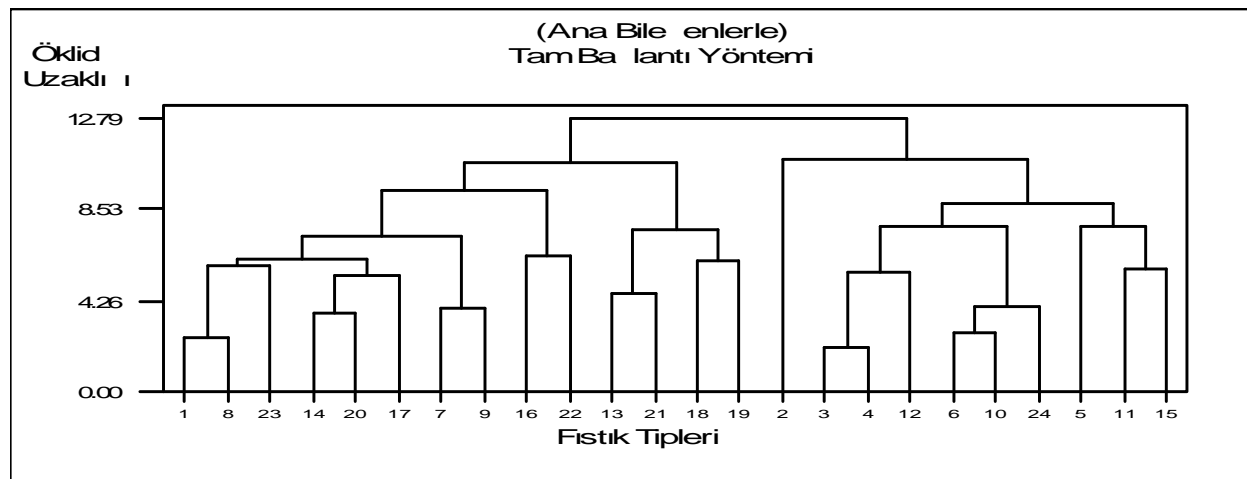
Fıstık Tipleri	Değişken Scorları									
	1	D2	D3	D4	D5	D6	D7	D8	D9	D10
1	-2,644	-1,323	-3,851	2,313	0,182	-0,801	0,019	0,627	0,308	0,408
2	3,267	-4,355	3,286	3,647	-0,416	1,094	-0,028	2,418	-0,314	1,265
3	2,679	2,944	-0,218	-1,818	-1,070	-1,115	0,158	0,268	1,584	0,970
4	2,764	2,684	0,311	-1,545	-0,020	-1,646	-0,908	0,727	1,156	0,011
5	6,254	-0,356	-1,688	-0,543	2,606	1,377	0,735	-1,921	0,325	2,310
6	3,898	0,548	-3,331	-0,093	-1,253	0,013	-0,420	0,511	-1,166	-1,058
7	-3,288	-1,644	-1,489	-1,327	0,042	1,620	-0,026	2,295	1,730	-0,290
8	-3,229	-0,641	-3,921	1,029	0,259	-0,108	-0,307	-0,456	-0,918	1,135
9	-1,410	-2,259	-1,226	-2,717	0,371	-1,011	0,760	1,816	0,572	-0,011
10	2,855	0,729	-2,309	1,665	-2,190	0,875	-1,212	-0,034	-1,152	-0,947
11	3,751	-1,105	1,301	-0,560	-0,697	-2,126	-0,105	0,182	-0,603	-1,640
12	1,450	4,513	0,652	1,617	1,358	-0,294	-2,294	-0,877	1,339	-0,118
13	-2,719	4,122	-0,872	1,634	-0,272	0,953	2,986	0,294	0,568	-0,313
14	-2,761	-0,471	1,485	0,952	0,854	-0,519	-0,505	0,480	0,745	-0,238
15	3,739	-1,500	1,275	-0,187	3,957	0,956	0,777	0,353	0,462	-1,977
16	0,405	-3,719	3,862	-0,561	-1,546	-0,072	-0,234	-0,960	-0,559	0,765
17	-3,967	-1,859	-0,031	1,297	0,835	0,830	-2,446	-1,842	0,388	-1,255
18	-3,343	1,897	1,751	-1,269	-0,133	-0,526	-2,333	0,772	-0,592	1,638
19	-1,953	5,001	1,874	-1,771	1,304	1,825	0,174	1,008	-3,511	-0,082
20	-3,267	-1,293	1,014	-0,161	0,380	-2,554	0,988	-1,184	-0,598	-0,261
21	-2,365	3,013	2,968	2,340	-0,756	-0,536	2,141	-0,949	0,730	-0,057
22	-1,804	-1,611	0,808	-2,812	-2,438	3,127	0,096	-1,667	1,247	-0,482
23	-1,509	-3,552	-1,207	-1,322	1,273	-1,218	0,995	-1,199	-1,368	0,201
24	3,201	0,237	-0,444	0,191	-2,631	-0,142	0,987	-0,662	-0,372	0,025



ekil 2. İlk iki Temel Bileşen Vasıtasıyla Elde Edilen Grupların (Birimlerin) Dağılımı



ekil 3. Temel Bileşenler Analizi ile Boyutu İndirgenen Değişkenlere İlişkin Ward Yöntemi ile Elde Edilen Açık Grafiği



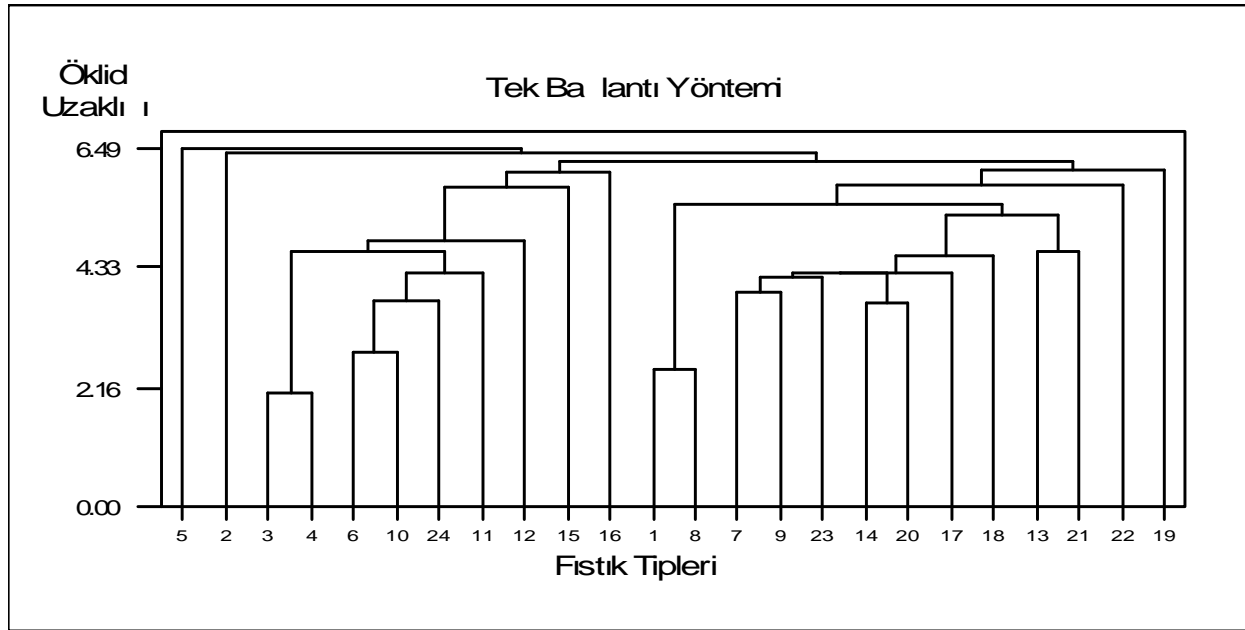
ekil 4. Temel Bileşenler Analizi ile Boyutu İndirgenmiş Olan Değişkenlere İlişkin Fıstık Tipleri ile Tam Bağılantı Yöntemine Göre Elde Edilen Açık Grafiği



Ancak ekil 5'de Tek Ba lantı yöntemine göre elde edilen dendogram Çizelge 1'de verilen fıstık tipleri ile kar ıla tırıldı nda erkenci ve geççi tiplerin birbiri ile kar ıla tırıldı ı görülmektedir. Di er kümeleme yöntemleri de benzer bir durum sergiledi inden ayrıntılı dendogramları burada vermeye gerek görülmedi.

Temel Bile enler analizi ile elemine edilen de i kenler, Çizelge 5'te verilen Hotelling Lawley z ölçütüne göre de erlendirildi inde; fıstık tiplerinin

( $p < 0.001$ ) önem düzeyinde 2, 3, 4, 5 kümeye ayrılabilce ini, ancak 2 kümeye ayrılması durumunda F de erinin daha büyük oldu unu görmekteyiz. ekil 3 incelendi inde fıstık tiplerinin erkenci ve geççi olmak üzere temel 2 kümeye ayrıldıkları görülmektedir. Bir alt kümede (3, 4 ve 5' inci kümelerde) ise fıstık tiplerinin birbirinden farklı özelliklerde alt türlere ayrıldıkları görülmektedir.



ekil 5. Temel Bile enler Analizi le Boyutu ndirgenmi Olan De i kenlere li kin Tek Ba lantı Yöntemi le Elde Edilen A aç Grafi i

Çizelge 5. Temel Bile enler Analizi le Boyutu ndirgenmi De i kenlere li kin Farklı Küme Sayılarının Hotelling Lawley Trace ( z) Ölçütü le Önemlilik Düzeylerinin Belirlenmesi

Kümeleme Yöntemi	Küme Sayısı	Hotelling Lawley z Ölçütü	Test statisti i F	Önemlilik Düzeyi P
k-Ortalama	2	15.56241	<b>20.231</b>	0.000<0.001
	3	15.59810	8.579	0.000<0.001
	4	23.41867	7.546	0.000<0.001
	5	73.53192	15.626	0.000<0.001
Ward Yöntemi	2	33.56928	<b>43.640</b>	0.000<0.001
	3	37.27117	20.499	0.000<0.001
	4	55.35712	17.837	0.000<0.001
	5	68.05547	14.462	0.000<0.001
Tam Ba lantı	2	22.56820	<b>29.339</b>	0.000<0.001
	3	40.64567	22.355	0.000<0.001
	4	44.54142	14.352	0.000<0.001
	5	61.49244	13.067	0.000<0.001

Temel Bile enleri esas olarak kümeleme yöntemleri kendi aralarında önemlilik seviyelerine göre de erlendirecek olursa; fıstık tipleri 2 ve 4 kümeye ayrıldı nda en iyi küme yapısını Ward yönteminin verdi ini Çizelge 6'da görmemiz mümkündür. Ancak, Fıstık tipleri 3 kümeye ayrıldı nda Tam Ba lantı yönteminin en anlamlı küme yapısını verdi ini yine

aynı cetvelden görebiliriz. Fıstık tipleri 5 kümeye ayrıldı nda ise k-Ortalama yöntemine ait F de erinin daha büyük oldu u ve di er yöntemlerden daha anlamlı küme yapısını verdi ini görmekteyiz.

Farklı boyut indirgeme metotlarına göre elemine edilen de i kenlerle Fıstık tipleri kümelendikten sonra, elde edilen kümelere bazı çok de i kenli test

istatistikleri uygulanarak Çizelge 7'deki sonuçlar elde edildi. Daha öncede sözü edildiği gibi Temel Bilemler analizi ile boyutu indirgenen de i kenler en iyi sonucu verdi i tespit edildi. Çizelge 7 incelendi inde  $C_{max}$ ,

Wilk's Lambda ve Hotelling Lawley z istatistikleri Fıstık tiplerinin iki kümeye ayrılması gerektiği görülmektedir. Bu durumu ekil 2 ve ekil 3'ün incelenmesinde görülebilmektedir.

Çizelge 6. Temel Bilemler Analizi ile Boyutu indirgenmiş de i kenlere li kin Farklı Küme Sayılarının Wilk's Lambda Test istatistiği ile Önemlilik Düzeylerinin Belirlenmesi

Kümeleme Yöntemi	Küme Sayısı	Wilk's Lambda Ölçütü	Test istatistiği F	Önemlilik Düzeyi P
K-Ortalama	2	0.06038	<b>20.231</b>	0.000<0.001
	3	0.01975	7.339	0.000<0.001
	4	0.00272	7.121	0.000<0.001
	5	0.00009	10.695	0.000<0.001
Ward Yöntemi	2	0.02893	<b>43.640</b>	0.000<0.001
	3	0.00627	13.956	0.000<0.001
	4	0.00083	11.209	0.000<0.001
	5	0.00013	9.629	0.000<0.001
Tam Ba lantı	2	0.04243	<b>29.339</b>	0.000<0.001
	3	0.00576	14.612	0.000<0.001
	4	0.00120	9.777	0.000<0.001
	5	0.00022	8.220	0.000<0.001

Çizelge 7. Kümeleme Kriterleri ile Küme Sayısının Ward Yöntemine Göre Belirlenmesi

De i ken Tipleri	Küme Sayıları	$C_{max}$	$M_{min}$	Wilk's Lambda		Hotelling Lawley z	
				F	p	F	P
Korelasyon analizi ile elemine edilen de i kenler	2	5.84	20.76E+30	<b>14.035</b>	0.000	<b>14.035</b>	0.000<0.001
	<b>3</b>	<b>17.39</b>	28.41E+29	6.338	0.000	6.887	0.001<0.01
	4	12.54	21.37E+28	7.845	0.000	7.286	0.000<0.001
	5	9.13	<b>12.73E+27</b>	6.918	0.000	5.697	0.000<0.001
Temel Bilemler analizi ile elemine edilen de i kenler	2	<b>6.98</b>	20.76E+30	<b>43.640</b>	0.000	<b>43.640</b>	0.000<0.001
	3	6.31	28.41E+29	13.956	0.000	20.499	0.000<0.001
	4	6.08	21.37E+28	11.209	0.000	17.837	0.000<0.001
	5	5.56	<b>12.73E+27</b>	9.629	0.000	14.462	0.000<0.001
Stepwise Diskriminant analizi ile elemine edilen de i kenler	2	1.99	34416674657	<b>57.40</b>	0.000	<b>57.40</b>	0.000<0.001
	3	3.34	5,70502E+12	12.85	0.000	12.54	0.000<0.001
	<b>4</b>	<b>19.70</b>	<b>2715327701</b>	13.00	0.000	17.61	0.000<0.001
	5	14.46	9846561289	20.53	0.000	20.24	0.000<0.001

## SONUÇ ve ÖNER LER

Sonuç olarak fıstık tiplerinin kümelenebilmesinde en önemli rol oynayan karakterlerin çiçeklenme zamanı ile ilgili de i kenler olduğu söylenebilir. Çünkü fıstık tipleri erkenci ve geççi olarak iki guruba ayrıldı nı görmekteyiz. Orta geççi fıstık tipleri ise büyük bir ihtimalle erkenci tiplerle geççi tiplerin birer melezi olduklarından, fenolojik olarak daha çok benzerlik gösterdikleri fıstık tipleri ile kümelendikleri ekil 2.'den görülebilmektedir.

De i ken sayısının kümelenecek birey sayısından fazla olması ( $n < p$ ) durumunda, de i ken sayısını azaltmak için en iyi yöntemin Temel Bilemler olduğu tespit edilmiştir. Özellikle küme sayısı hakkında önceden herhangi bir ön bilgi olmaması durumunda temel bilemlerle skorları ile tatmin edici sonuçlar elde edilebilmektedir.

Küme yapısının önceden bilinmesi durumunda, diskriminant analizi ile bireylerin doğru kümelere atanıp

atanmadıkları konusunda önemli bilgilerin elde edilmesi mümkündür. Ancak söz konusu veri setimize ait küme yapısının ve de i kenlerin etkinliği hakkında herhangi bir bilgi olmaması durumunda, Stepwise Diskriminant analizi ile etkili de i kenleri belirleyerek boyut indirgemek pek mümkün olamamaktadır.

Blashfield, (1976); Monte Carlo teknikleri ile karışık da ılıma sahip veriler türeterek, oluşturulan örneklerle farklı kümeleme yöntemlerini uygulamıştır. Sözü edilen karışık da ılımlı modeli en iyi belirleyen yöntemin Ward yöntemi olduğu belirlenmiştir olup, bu durum çalı mamızda ki sonuçları desteklemektedir.

Ancak en iyi kümeleme sonuçlarını veren yöntemin, her zaman Ward yöntemi olduğu söylemek pek doğrudur. Çünkü birimlerin kümelenebilmesinde, kümeleme yöntemlerinin gösterdikleri başarı verisi yapısına bağlı olarak değişmektedir. Yani Verilerin dağılımı ekli de i tikçe, farklı bir kümeleme yöntemi daha iyi sonuçlar verebilmektedir. Bununla ilgili bir kısım

çalı malar yapılmı olup, Fisher and Van Ness, (1971) ve Mainly (1994)'in de açıklamaları bu doğrultudadır.

Bu çalışmada sürecinde verileri tavsiye edilen herhangi bir kümeleme yöntemi ile de erlendirmenin çok zor olduğu tespit edilmiştir. Veriler hakkında geniş bir bilgi birikimi olmadan veya konunun uzmanı olmadan sadece bazı kümeleme yöntemlerine ve çok değişkenli istatistik kriterlerine bağlı olarak sonuçları de erlendirip yorumlamanın doğru olmadığı kanısına varmış bulunmaktayız. Çünkü kümeleme sonuçları veri setindeki değişkenlerin etkinliğine göre farklılık arz etmektedir. Bu durum verilerin çok boyutlu uzaydaki dağılımının ve değişken yapısı ile değişken seçiminin önemini ortaya koymaktadır.

#### TE EKKÜR

Doktora çalışmam sürecinde katkılarını ve yardımlarını esirgemeyen Prof. Dr. Yüksel BEK'e en içten teşekkürlerimi sunarım.

#### KAYNAKLAR

- Acar, .., 1997. Ceylanpınar Tarım İletmesinde Seçilmiş Bazı Erkek Antep Fıstığı Tiplerinin Morfolojik ve Biyolojik Özellikleri Üzerinde Bir Araştırma. Harran Üniv. Fen Bil. Enst. Y. Lisans Tezi (Basılmamış), 92 s.
- Anderberg, M.R., 1973. Cluster Analysis For Applications. Academic Press, New York 359 S.
- Blashfield, R.K., 1976. Mixture Model Test of Cluster Analysis: Accuracy of Four Agglomerative Hierarchical Methods. Psychological Bull., 83: 377-388.
- Bliss F.A., 1981. Utilization of Vegetable Germplasm. Hortscience 16(2): 129-132.
- Fisher, L., Van Ness, J.W., 1971. Admissible Clustering Procedures. Biometrika, 58: (1), 91-104.
- Hawkins, D. M., Müller, M.W., Kroonen, J.A., 1982. Topics in Applied Multivariate Analysis. Cambridge University Press, USA.
- Huang Tung C., 2001. The Effects of Linkage Between Business and Human Resource Management Strategies, Personal Review, 30(2):132-151.
- Jain, S.K., 1975. Population Structure and the Effects of Breeding System. in Crop Genetic Resources for Today and Tomorrow, Ed. Frankel, O.H. and Hawkes, J.G., Cambridge Uni. Press.
- Johnson, A.R., Wichern, D.W., 1992. Applied Multivariate Statistical Analysis. Prentice Hall International Inc. New Jersey.
- Mac Queen, J., 1967. Some Methods for Classifications and Stopping Rules: An Evaluation. Comp. J., 20: 359-363.
- Mainly, B.F.J., 1994. Multivariate Statistical Methods, Second Edition, Londra: Chapman- Hall. New York,
- Minitab Inc, 1995. Minitab Release, V: 10,51.
- Özdamar, K., 1988. Hastalık Olgularının İncelenmesinde Kümeleme Çözümlemesinin Kullanılması. T.C. Anadolu Üniversitesi Yayınları No:295, Tıp Fakültesi Yayınları No:25, Eskişehir.
- Özdamar, K., 1999. Paket Programlar ile İstatistiksel Veri Analizi. Kaan Kitabevi, Eskişehir, 214s.
- Sharma S., 1996. Applied Multivariate Techniques. John Wiley & Sons.Inc. New York.
- Sneath, P. 1957. The Application of Computers to Taxonomy. Journal of General Microbiology, 17, 201-226.
- Sokal R.R., Michener C.D. 1958. A Statistical Method for Evaluating Systematic Relationships. The University of Kansas Scientific Bulletin 38: 1409-1438.
- Tan, A., 1983. Sayısal Taksonomik Yöntemlerle Varyasyonun Saptanması. Ege Bölge Ziraat Enst. Yay. No: 30, Menemen- İzmir.
- Tatlıdil, H., 1992. Uygulamalı Çok Değişkenli İstatistik Analiz. Hacettepe Üniversitesi, Ankara.