

Johnson, R. A., & Wichern, D. W. (2007). Applied multivariate statistical analysis (6th Ed.). Upper Saddle River, NJ: Prentice Hall.

## Chapter

# 12

## CLUSTERING, DISTANCE METHODS, AND ORDINATION

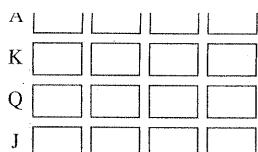
### 12.1 Introduction

Rudimentary, exploratory procedures are often quite helpful in understanding the complex nature of multivariate relationships. For example, throughout this book, we have emphasized the value of data plots. In this chapter, we shall discuss some additional displays based on certain measures of distance and suggested step-by-step rules (algorithms) for grouping objects (variables or items). Searching the data for a structure of "natural" groupings is an important exploratory technique. Groupings can provide an informal means for assessing dimensionality, identifying outliers, and suggesting interesting hypotheses concerning relationships.

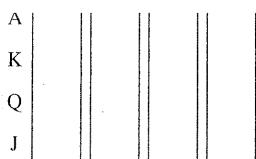
Grouping, or clustering, is distinct from the classification methods discussed in the previous chapter. Classification pertains to a *known* number of groups, and the operational objective is to assign new observations to one of these groups. Cluster analysis is a more primitive technique in that no assumptions are made concerning the number of groups or the group structure. Grouping is done on the basis of similarities or distances (dissimilarities). The inputs required are similarity measures or data from which similarities can be computed.

To illustrate the nature of the difficulty in defining a natural grouping, consider sorting the 16 face cards in an ordinary deck of playing cards into clusters of similar objects. Some groupings are illustrated in Figure 12.1. It is immediately clear that meaningful partitions depend on the definition of *similar*.

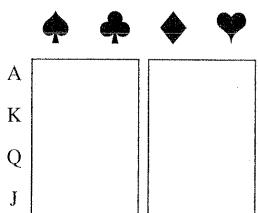
In most practical applications of cluster analysis, the investigator knows enough about the problem to distinguish "good" groupings from "bad" groupings. Why not enumerate all possible groupings and select the "best" ones for further study?



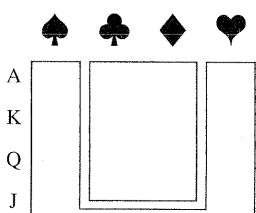
(a) Individual cards



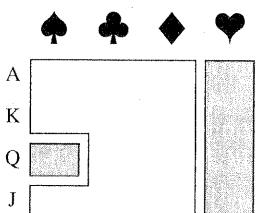
(b) Individual suits



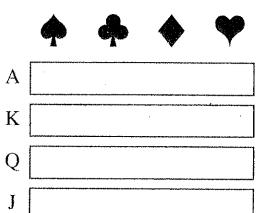
(c) Black and red suits



(d) Major and minor suits (bridge)



(e) Hearts plus queen of spades and other suits (hearts)



(f) Like face cards

**Figure 12.1** Grouping face cards.

For the playing-card example, there is one way to form a *single* group of 16 face cards, there are 32,767 ways to partition the face cards into *two* groups (of varying sizes), there are 7,141,686 ways to sort the face cards into *three* groups (of varying sizes), and so on.<sup>1</sup> Obviously, time constraints make it impossible to determine the best groupings of similar objects from a list of all possible structures. Even fast computers are easily overwhelmed by the typically large number of cases, so one must settle for *algorithms* that search for good, but not necessarily the best, groupings.

To summarize, the basic objective in cluster analysis is to discover natural groupings of the items (or variables). In turn, we must first develop a quantitative scale on which to measure the association (similarity) between objects. Section 12.2 is devoted to a discussion of similarity measures. After that section, we describe a few of the more common algorithms for sorting objects into groups.

<sup>1</sup>The number of ways of sorting  $n$  objects into  $k$  nonempty groups is a Stirling number of the second kind given by  $(1/k!)\sum_{j=0}^k (-1)^{k-j}\binom{k}{j}j^n$ . (See [1].) Adding these numbers for  $k = 1, 2, \dots, n$  groups, we obtain the total number of possible ways to sort  $n$  objects into groups.

Even without the precise notion of a natural grouping, we are often able to group objects in two- or three-dimensional plots by eye. Stars and Chernoff faces, discussed in Section 1.4, have been used for this purpose. (See Examples 1.11 and 1.12.) Additional procedures for depicting high-dimensional observations in two dimensions such that similar objects are, in some sense, close to one another are considered in Sections 12.5–12.7.

its

ts (bridge)

ds

## 12.2 Similarity Measures

Most efforts to produce a rather simple group structure from a complex data set require a measure of “closeness,” or “similarity.” There is often a great deal of subjectivity involved in the choice of a similarity measure. Important considerations include the nature of the variables (discrete, continuous, binary), scales of measurement (nominal, ordinal, interval, ratio), and subject matter knowledge.

When *items* (units or cases) are clustered, proximity is usually indicated by some sort of distance. By contrast, *variables* are usually grouped on the basis of correlation coefficients or like measures of association.

### Distances and Similarity Coefficients for Pairs of Items

We discussed the notion of distance in Chapter 1, Section 1.5. Recall that the Euclidean (straight-line) distance between two  $p$ -dimensional observations (items)  $\mathbf{x}' = [x_1, x_2, \dots, x_p]$  and  $\mathbf{y}' = [y_1, y_2, \dots, y_p]$  is, from (1-12),

$$\begin{aligned} d(\mathbf{x}, \mathbf{y}) &= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots + (x_p - y_p)^2} \\ &= \sqrt{(\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y})} \end{aligned} \quad (12-1)$$

The statistical distance between the same two observations is of the form [see (1-23)]

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})' \mathbf{A} (\mathbf{x} - \mathbf{y})} \quad (12-2)$$

Ordinarily,  $\mathbf{A} = \mathbf{S}^{-1}$ , where  $\mathbf{S}$  contains the sample variances and covariances. However, without prior knowledge of the distinct groups, these sample quantities cannot be computed. For this reason, Euclidean distance is often preferred for clustering.

Another distance measure is the Minkowski metric

$$d(\mathbf{x}, \mathbf{y}) = \left[ \sum_{i=1}^p |x_i - y_i|^m \right]^{1/m} \quad (12-3)$$

For  $m = 1$ ,  $d(\mathbf{x}, \mathbf{y})$  measures the “city-block” distance between two points in  $p$  dimensions. For  $m = 2$ ,  $d(\mathbf{x}, \mathbf{y})$  becomes the Euclidean distance. In general, varying  $m$  changes the weight given to larger and smaller differences.

single group of  
two groups (of  
three groups  
it impossible to  
ll possible struc-  
lly large number  
it not necessarily

discover natural  
a quantitative  
ects. Section 12.2  
n, we describe a  
s.

number of the second  
 $1, 2, \dots, n$  groups, we

Two additional popular measures of "distance" or dissimilarity are given by the Canberra metric and the Czekanowski coefficient. Both of these measures are defined for nonnegative variables only. We have

$$\text{Canberra metric: } d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{(x_i + y_i)} \quad (12-4)$$

$$\text{Czekanowski coefficient: } d(\mathbf{x}, \mathbf{y}) = 1 - \frac{2 \sum_{i=1}^p \min(x_i, y_i)}{\sum_{i=1}^p (x_i + y_i)} \quad (12-5)$$

Whenever possible, it is advisable to use "true" distances—that is, distances satisfying the distance properties of (1-25)—for clustering objects. On the other hand, most clustering algorithms will accept subjectively assigned distance numbers that may not satisfy, for example, the triangle inequality.

When items cannot be represented by meaningful  $p$ -dimensional measurements, pairs of items are often compared on the basis of the presence or absence of certain characteristics. Similar items have more characteristics in common than do dissimilar items. The presence or absence of a characteristic can be described mathematically by introducing a *binary variable*, which assumes the value 1 if the characteristic is present and the value 0 if the characteristic is absent. For  $p = 5$  binary variables, for instance, the "scores" for two items  $i$  and  $k$  might be arranged as follows:

	Variables				
	1	2	3	4	5
Item $i$	1	0	0	1	1
Item $k$	1	1	0	1	0

In this case, there are two 1–1 matches, one 0–0 match, and two mismatches.

Let  $x_{ij}$  be the score (1 or 0) of the  $j$ th binary variable on the  $i$ th item and  $x_{kj}$  be the score (again, 1 or 0) of the  $j$ th variable on the  $k$ th item,  $j = 1, 2, \dots, p$ . Consequently,

$$(x_{ij} - x_{kj})^2 = \begin{cases} 0 & \text{if } x_{ij} = x_{kj} = 1 \text{ or } x_{ij} = x_{kj} = 0 \\ 1 & \text{if } x_{ij} \neq x_{kj} \end{cases} \quad (12-6)$$

and the squared Euclidean distance,  $\sum_{j=1}^p (x_{ij} - x_{kj})^2$ , provides a count of the number of mismatches. A large distance corresponds to many mismatches—that is, dissimilar items. From the preceding display, the square of the distance between items  $i$  and  $k$  would be

$$\begin{aligned} \sum_{j=1}^5 (x_{ij} - x_{kj})^2 &= (1 - 1)^2 + (0 - 1)^2 + (0 - 0)^2 + (1 - 1)^2 + (1 - 0)^2 \\ &= 2 \end{aligned}$$

are given by the  
se measures are

(12-4)

(12-5)

distances satisfy-  
the other hand,  
nce numbers that

nisional measure-  
ence or absence of  
common than do  
can be described  
the value 1 if the  
bsent. For  $p = 5$   
ght be arranged as

ismatches.  
tem and  $x_{kj}$  be the  
,  $p$ . Consequently,

$i = 0$  (12-6)

ount of the number  
s—that is, dissimi-  
between items  $i$  and

$1)^2 + (1 - 0)^2$

Although a distance based on (12-6) might be used to measure similarity, it suffers from weighting the 1-1 and 0-0 matches equally. In some cases, a 1-1 match is a stronger indication of similarity than a 0-0 match. For instance, in grouping people, the evidence that two persons both read ancient Greek is stronger evidence of similarity than the absence of this ability. Thus, it might be reasonable to discount the 0-0 matches or even disregard them completely. To allow for differential treatment of the 1-1 matches and the 0-0 matches, several schemes for defining similarity coefficients have been suggested.

To introduce these schemes, let us arrange the frequencies of matches and mismatches for items  $i$  and  $k$  in the form of a contingency table:

		Item $k$		Totals
		1	0	
Item $i$	1	$a$	$b$	$a + b$
	0	$c$	$d$	$c + d$
Totals		$a + c$	$b + d$	$p = a + b + c + d$

(12-7)

In this table,  $a$  represents the frequency of 1-1 matches,  $b$  is the frequency of 1-0 matches, and so forth. Given the foregoing five pairs of binary outcomes,  $a = 2$  and  $b = c = d = 1$ .

Table 12.1 lists common similarity coefficients defined in terms of the frequencies in (12-7). A short rationale follows each definition.

**Table 12.1** Similarity Coefficients for Clustering Items\*

Coefficient	Rationale
1. $\frac{a + d}{p}$	Equal weights for 1-1 matches and 0-0 matches.
2. $\frac{2(a + d)}{2(a + d) + b + c}$	Double weight for 1-1 matches and 0-0 matches.
3. $\frac{a + d}{a + d + 2(b + c)}$	Double weight for unmatched pairs.
4. $\frac{a}{p}$	No 0-0 matches in numerator.
5. $\frac{a}{a + b + c}$	No 0-0 matches in numerator or denominator. (The 0-0 matches are treated as irrelevant.)
6. $\frac{2a}{2a + b + c}$	No 0-0 matches in numerator or denominator. Double weight for 1-1 matches.
7. $\frac{a}{a + 2(b + c)}$	No 0-0 matches in numerator or denominator. Double weight for unmatched pairs.
8. $\frac{a}{b + c}$	Ratio of matches to mismatches with 0-0 matches excluded.

\*[ $p$  binary variables; see (12-7).]

Coefficients 1, 2, and 3 in the table are monotonically related. Suppose coefficient 1 is calculated for two contingency tables, Table I and Table II. Then if  $(a_{\text{I}} + d_{\text{I}})/p \geq (a_{\text{II}} + d_{\text{II}})/p$ , we also have  $2(a_{\text{I}} + d_{\text{I}})/[2(a_{\text{I}} + d_{\text{I}}) + b_{\text{I}} + c_{\text{I}}] \geq 2(a_{\text{II}} + d_{\text{II}})/[2(a_{\text{II}} + d_{\text{II}}) + b_{\text{II}} + c_{\text{II}}]$ , and coefficient 3 will be at least as large for Table I as it is for Table II. (See Exercise 12.4.) Coefficients 5, 6, and 7 also retain their relative orders.

Monotonicity is important, because some clustering procedures are not affected if the definition of similarity is changed in a manner that leaves the relative orderings of similarities unchanged. The single linkage and complete linkage hierarchical procedures discussed in Section 12.3 are not affected. For these methods, any choice of the coefficients 1, 2, and 3 in Table 12.1 will produce the same groupings. Similarly, any choice of the coefficients 5, 6, and 7 will yield identical groupings.

**Example 12.1 (Calculating the values of a similarity coefficient)** Suppose five individuals possess the following characteristics:

	Height	Weight	Eye color	Hair color	Handedness	Gender
Individual 1	68 in	140 lb	green	blond	right	female
Individual 2	73 in	185 lb	brown	brown	right	male
Individual 3	67 in	165 lb	blue	blond	right	male
Individual 4	64 in	120 lb	brown	brown	right	female
Individual 5	76 in	210 lb	brown	brown	left	male

Define six binary variables  $X_1, X_2, X_3, X_4, X_5, X_6$  as

$$\begin{aligned} X_1 &= \begin{cases} 1 & \text{height} \geq 72 \text{ in.} \\ 0 & \text{height} < 72 \text{ in.} \end{cases} & X_4 &= \begin{cases} 1 & \text{blond hair} \\ 0 & \text{not blond hair} \end{cases} \\ X_2 &= \begin{cases} 1 & \text{weight} \geq 150 \text{ lb} \\ 0 & \text{weight} < 150 \text{ lb} \end{cases} & X_5 &= \begin{cases} 1 & \text{right handed} \\ 0 & \text{left handed} \end{cases} \\ X_3 &= \begin{cases} 1 & \text{brown eyes} \\ 0 & \text{otherwise} \end{cases} & X_6 &= \begin{cases} 1 & \text{female} \\ 0 & \text{male} \end{cases} \end{aligned}$$

The scores for individuals 1 and 2 on the  $p = 6$  binary variables are

		$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$
Individual	1	0	0	0	1	1	1
	2	1	1	1	0	1	0

and the number of matches and mismatches are indicated in the two-way array

		Individual 2		Total
		1	0	
Individual 1	1	1	2	3
	0	3	0	3
Totals		4	2	6

related. Suppose that  $a + b_I + c_I \geq d_I$ . Then  $a + b_I + c_I \geq d_I + b_I + c_I$  at least as large as  $d_I$ . Since 6 and 7 also relate,

are not affected by relative orderings of age or hierarchical methods, any choice of upings. Similarly,

suppose five indi-

ness	Gender
1	female
2	male
3	male
4	female
5	male

r  
I hair  
ded  
ed

	$X_5$	$X_6$
1	1	0
2	0	1

two-way array

Employing similarity coefficient 1, which gives equal weight to matches, we compute

$$\frac{a + d}{p} = \frac{1 + 0}{6} = \frac{1}{6}$$

Continuing with similarity coefficient 1, we calculate the remaining similarity numbers for pairs of individuals. These are displayed in the  $5 \times 5$  symmetric matrix

		Individual				
		1	2	3	4	5
Individual	1	1				
	2	$\frac{1}{6}$	1			
	3	$\frac{4}{6}$	$\frac{3}{6}$	1		
	4	$\frac{4}{6}$	$\frac{3}{6}$	$\frac{2}{6}$	1	
	5	0	( $\frac{5}{6}$ )	$\frac{2}{6}$	$\frac{2}{6}$	1

Based on the magnitudes of the similarity coefficient, we should conclude that individuals 2 and 5 are most similar and individuals 1 and 5 are least similar. Other pairs fall between these extremes. If we were to divide the individuals into two relatively homogeneous subgroups on the basis of the similarity numbers, we might form the subgroups (1 3 4) and (2 5).

Note that  $X_3 = 0$  implies an absence of brown eyes, so that two people, one with blue eyes and one with green eyes, will yield a 0–0 match. Consequently, it may be inappropriate to use similarity coefficient 1, 2, or 3 because these coefficients give the same weights to 1–1 and 0–0 matches.

We have described the construction of distances and similarities. It is always possible to construct similarities from distances. For example, we might set

$$\tilde{s}_{ik} = \frac{1}{1 + d_{ik}} \quad (12-8)$$

where  $0 < \tilde{s}_{ik} \leq 1$  is the similarity between items  $i$  and  $k$  and  $d_{ik}$  is the corresponding distance.

However, distances that must satisfy (1-25) cannot always be constructed from similarities. As Gower [11, 12] has shown, this can be done only if the matrix of similarities is nonnegative definite. With the nonnegative definite condition, and with the maximum similarity scaled so that  $\tilde{s}_{ii} = 1$ ,

$$d_{ik} = \sqrt{2(1 - \tilde{s}_{ik})} \quad (12-9)$$

has the properties of a distance.

## Similarities and Association Measures for Pairs of Variables

Thus far, we have discussed similarity measures for items. In some applications, it is the variables, rather than the items, that must be grouped. Similarity measures for variables often take the form of sample correlation coefficients. Moreover, in some clustering applications, negative correlations are replaced by their absolute values.

When the variables are binary, the data can again be arranged in the form of a contingency table. This time, however, the variables, rather than the items, delineate the categories. For each pair of variables, there are  $n$  items categorized in the table. With the usual 0 and 1 coding, the table becomes as follows:

		Variable $k$		Totals	(12-10)
		1	0		
Variable $i$	1	$a$	$b$	$a + b$	
	0	$c$	$d$	$c + d$	
Totals		$a + c$	$b + d$	$n = a + b + c + d$	

For instance, variable  $i$  equals 1 and variable  $k$  equals 0 for  $b$  of the  $n$  items.

The usual product moment correlation formula applied to the binary variables in the contingency table of (12-10) gives (see Exercise 12.3)

$$r = \frac{ad - bc}{[(a+b)(c+d)(a+c)(b+d)]^{1/2}} \quad (12-11)$$

This number can be taken as a measure of the similarity between the two variables.

The correlation coefficient in (12-11) is related to the chi-square statistic ( $r^2 = \chi^2/n$ ) for testing the independence of two categorical variables. For  $n$  fixed, a large similarity (or correlation) is consistent with the presence of dependence.

Given the table in (12-10), measures of association (or similarity) exactly analogous to the ones listed in Table 12.1 can be developed. The only change required is the substitution of  $n$  (the number of items) for  $p$  (the number of variables).

### Concluding Comments on Similarity

To summarize this section, we note that there are many ways to measure the similarity between pairs of objects. It appears that most practitioners use distances [see (12-1) through (12-5)] or the coefficients in Table 12.1 to cluster *items* and correlations to cluster *variables*. However, at times, inputs to clustering algorithms may be simple frequencies.

---

**Example 12.2 (Measuring the similarities of 11 languages)** The meanings of words change with the course of history. However, the meaning of the numbers 1, 2, 3, ... represents one conspicuous exception. Thus, a first comparison of languages might be based on the numerals alone. Table 12.2 gives the first 10 numbers in English, Polish, Hungarian, and eight other modern European languages. (Only languages that use the Roman alphabet are considered, and accent marks, cedillas, dieresis, etc., are omitted.) A cursory examination of the spelling of the numerals in the table suggests that the first five languages (English, Norwegian, Danish, Dutch, and German) are very much alike. French, Spanish, and Italian are in even closer agreement. Hungarian and Finnish seem to stand by themselves, and Polish has some of the characteristics of the languages in each of the larger subgroups.

in the form of a  
e items, delineate  
ized in the table.

(12-10)

$c + d$

$\geq n$  items.  
binary variables

(12-11)

the two variables.  
square statistic  
oles. For  $n$  fixed, a  
dependence.  
ty) exactly analo-  
hange required is  
variables).

measure the simi-  
use distances [see  
*items* and correla-  
algorithms may be

meanings of words  
umbers 1, 2, 3, ...  
f languages might  
mbers in English,  
(Only languages  
cedillas, dieresis,  
nerals in the table  
, Dutch, and Ger-  
closer agreement.  
has some of the

**Table 12.2** Numerals in 11 Languages

English (E)	Norwegian (N)	Danish (Da)	Dutch (Du)	German (G)	French (Fr)	Spanish (Sp)	Italian (I)	Polish (P)	Hungarian (H)	Finnish (Fi)
one	en	en	een	eins	un	uno	uno	jeden	egy	yksi
two	to	to	twee	zwei	deux	dos	due	dwa	ketto	kaksi
three	tre	tre	drie	drei	trois	tres	tre	trzy	harom	kolme
four	fire	fire	vier	vier	quatre	cuatro	quattro	cztery	nagy	neljä
five	fem	fem	vijf	funf	cinq	cinco	cinque	pięc	öt	viisi
six	seks	seks	zes	sechs	six	seis	sei	szesc	hat	kuusi
seven	sju	sju	zeven	sieben	sept	siete	sette	siedem	het	seitsemän
eight	atte	atte	otte	acht	huit	ocho	otto	osiem	nyolc	kahdeksan
nine	ni	ni	negen	neun	neuf	nueve	nove	dziewiec	kilenc	yhdeksan
ten	ti	ti	tien	zehn	dix	diez	dieci	dziesiec	tiz	kymmenen

**Table 12.3** Concordant First Letters for Numbers in 11 Languages

	E	N	Da	Du	G	Fr	Sp	I	P	H	Fi
E	10										
N	8	10									
Da	8	9	10								
Du	3	5	4	10							
G	4	6	5	5	10						
Fr	4	4	4	1	3	10					
Sp	4	4	5	1	3	8	10				
I	4	4	5	1	3	9	9	10			
P	3	3	4	0	2	5	7	6	10		
H	1	2	2	2	1	0	0	0	0	10	
Fi	1	1	1	1	1	1	1	1	2	10	

The words for 1 in French, Spanish, and Italian all begin with *u*. For illustrative purposes, we might compare languages by looking at the *first letters* of the numbers. We call the words for the same number in two different languages *concordant* if they have the same first letter and *discordant* if they do not. From Table 12.2, the table of concordances (frequencies of matching first initials) for the numbers 1–10 is given in Table 12.3. We see that English and Norwegian have the same first letter for 8 of the 10 word pairs. The remaining frequencies were calculated in the same manner.

The results in Table 12.3 confirm our initial visual impression of Table 12.2. That is, English, Norwegian, Danish, Dutch, and German seem to form a group. French, Spanish, Italian, and Polish might be grouped together, whereas Hungarian and Finnish appear to stand alone. ■

In our examples so far, we have used our visual impression of similarity or distance measures to form groups. We now discuss less subjective schemes for creating clusters.

## 12.3 Hierarchical Clustering Methods

We can rarely examine all grouping possibilities, even with the largest and fastest computers. Because of this problem, a wide variety of clustering algorithms have emerged that find “reasonable” clusters without having to look at all configurations.

Hierarchical clustering techniques proceed by either a series of successive mergers or a series of successive divisions. *Agglomerative hierarchical methods* start with the individual objects. Thus, there are initially as many clusters as objects. The most similar objects are first grouped, and these initial groups are merged according to their similarities. Eventually, as the similarity decreases, all subgroups are fused into a single cluster.

*Divisive hierarchical methods* work in the opposite direction. An initial single group of objects is divided into two subgroups such that the objects in one subgroup are “far from” the objects in the other. These subgroups are then further divided into dissimilar subgroups; the process continues until there are as many subgroups as objects—that is, until each object forms a group.

Languages	
H	Fi
10	
2	10

For illustrative purposes of the numbers, we say two numbers are *concordant* if they have the same sign. In Table 12.2, the table of concordance for pairs of numbers 1–10 is given in the last column. The letter for 8 of the numbers is the same as that for 10 of the numbers. That is, 10 is a group. French, German, and Hungarian are concordant.

of similarity or dissimilarity schemes for creating

largest and fastest algorithms have all configurations. series of successive hierarchical methods start with objects. The merged according subgroups are fused

1. An initial single step is to put all entities in one subgroup. This can further divide the entities into as many subgroups as

The results of both agglomerative and divisive methods may be displayed in the form of a two-dimensional diagram known as a *dendrogram*. As we shall see, the dendrogram illustrates the mergers or divisions that have been made at successive levels.

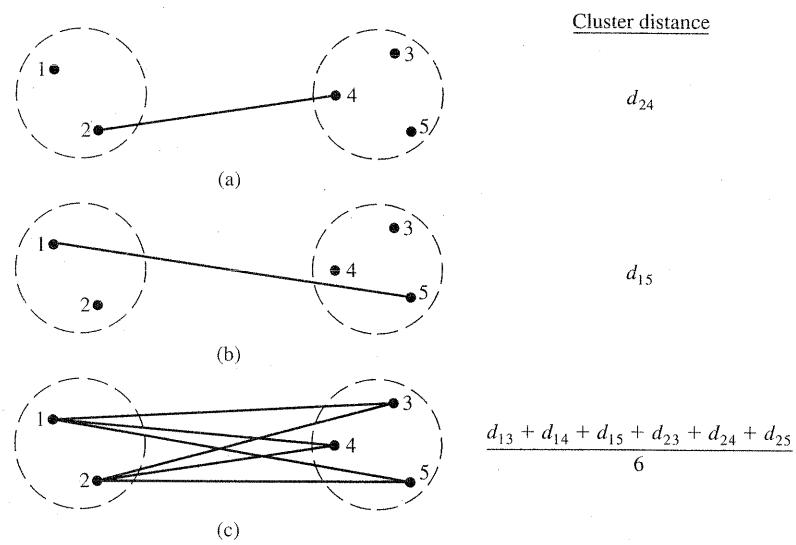
In this section we shall concentrate on agglomerative hierarchical procedures and, in particular, *linkage methods*. Excellent elementary discussions of divisive hierarchical procedures and other agglomerative techniques are available in [3] and [8].

Linkage methods are suitable for clustering items, as well as variables. This is not true for all hierarchical agglomerative procedures. We shall discuss, in turn, *single linkage* (minimum distance or nearest neighbor), *complete linkage* (maximum distance or farthest neighbor), and *average linkage* (average distance). The merging of clusters under the three linkage criteria is illustrated schematically in Figure 12.2.

From the figure, we see that single linkage results when groups are fused according to the distance between their nearest members. Complete linkage occurs when groups are fused according to the distance between their farthest members. For average linkage, groups are fused according to the average distance between pairs of members in the respective sets.

The following are the steps in the agglomerative hierarchical clustering algorithm for grouping  $N$  objects (items or variables):

1. Start with  $N$  clusters, each containing a single entity and an  $N \times N$  symmetric matrix of distances (or similarities)  $\mathbf{D} = \{d_{ik}\}$ .
2. Search the distance matrix for the nearest (most similar) pair of clusters. Let the distance between “most similar” clusters  $U$  and  $V$  be  $d_{UV}$ .



**Figure 12.2** Intercluster distance (dissimilarity) for (a) single linkage, (b) complete linkage, and (c) average linkage.

3. Merge clusters  $U$  and  $V$ . Label the newly formed cluster ( $UV$ ). Update the entries in the distance matrix by (a) deleting the rows and columns corresponding to clusters  $U$  and  $V$  and (b) adding a row and column giving the distances between cluster ( $UV$ ) and the remaining clusters.
4. Repeat Steps 2 and 3 a total of  $N - 1$  times. (All objects will be in a *single* cluster after the algorithm terminates.) Record the identity of clusters that are merged and the levels (distances or similarities) at which the mergers take place. (12-12)

The ideas behind any clustering procedure are probably best conveyed through examples, which we shall present after brief discussions of the input and algorithmic components of the linkage methods.

### Single Linkage

The inputs to a single linkage algorithm can be distances or similarities between pairs of objects. Groups are formed from the individual entities by merging nearest neighbors, where the term *nearest neighbor* connotes the smallest distance or largest similarity.

Initially, we must find the smallest distance in  $\mathbf{D} = \{d_{ik}\}$  and merge the corresponding objects, say,  $U$  and  $V$ , to get the cluster ( $UV$ ). For Step 3 of the general algorithm of (12-12), the distances between ( $UV$ ) and any other cluster  $W$  are computed by

$$d_{(UV)W} = \min \{d_{UW}, d_{VW}\} \quad (12-13)$$

Here the quantities  $d_{UW}$  and  $d_{VW}$  are the distances between the nearest neighbors of clusters  $U$  and  $W$  and clusters  $V$  and  $W$ , respectively.

The results of single linkage clustering can be graphically displayed in the form of a *dendrogram*, or tree diagram. The branches in the tree represent clusters. The branches come together (merge) at nodes whose positions along a distance (or similarity) axis indicate the level at which the fusions occur. Dendograms for some specific cases are considered in the following examples.

---

**Example 12.3 (Clustering using single linkage)** To illustrate the single linkage algorithm, we consider the hypothetical distances between pairs of five objects as follows:

$$\mathbf{D} = \{d_{ik}\} = \begin{bmatrix} & 1 & 2 & 3 & 4 & 5 \\ 1 & 0 & & & & \\ 2 & 9 & 0 & & & \\ 3 & 3 & 7 & 0 & & \\ 4 & 6 & 5 & 9 & 0 & \\ 5 & 11 & 10 & ② & 8 & 0 \end{bmatrix}$$

Treating each object as a cluster, we commence clustering by merging the two closest items. Since

$$\min_{i,k} (d_{ik}) = d_{53} = 2$$

). Update the entries corresponding to the distances between

vill be in a *single* set of clusters that the mergers take  
(12-12)

conveyed through it and algorithmic

similarities between merging nearest distance or largest

and merge the step 3 of the general cluster  $W$  are

(12-13)

nearest neighbors

played in the form present clusters. The a distance (or programs for some

the single linkage of five objects as

merging the two

objects 5 and 3 are merged to form the cluster (35). To implement the next level of clustering, we need the distances between the cluster (35) and the remaining objects, 1, 2, and 4. The nearest neighbor distances are

$$\begin{aligned}d_{(35)1} &= \min \{d_{31}, d_{51}\} = \min \{3, 11\} = 3 \\d_{(35)2} &= \min \{d_{32}, d_{52}\} = \min \{7, 10\} = 7 \\d_{(35)4} &= \min \{d_{34}, d_{54}\} = \min \{9, 8\} = 8\end{aligned}$$

Deleting the rows and columns of  $\mathbf{D}$  corresponding to objects 3 and 5, and adding a row and column for the cluster (35), we obtain the new distance matrix

$$(35) \quad \begin{matrix} (35) & 1 & 2 & 4 \\ 0 & & & \\ \textcircled{3} & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{matrix}$$

The smallest distance between pairs of clusters is now  $d_{(35)1} = 3$ , and we merge cluster (1) with cluster (35) to get the next cluster, (135). Calculating

$$\begin{aligned}d_{(135)2} &= \min \{d_{(35)2}, d_{12}\} = \min \{7, 9\} = 7 \\d_{(135)4} &= \min \{d_{(35)4}, d_{14}\} = \min \{8, 6\} = 6\end{aligned}$$

we find that the distance matrix for the next level of clustering is

$$(135) \quad \begin{matrix} (135) & 2 & 4 \\ 0 & & \\ 7 & 0 & \\ 6 & \textcircled{5} & 0 \end{matrix}$$

The minimum nearest neighbor distance between pairs of clusters is  $d_{42} = 5$ , and we merge objects 4 and 2 to get the cluster (24).

At this point we have two distinct clusters, (135) and (24). Their nearest neighbor distance is

$$d_{(135)(24)} = \min \{d_{(135)2}, d_{(135)4}\} = \min \{7, 6\} = 6$$

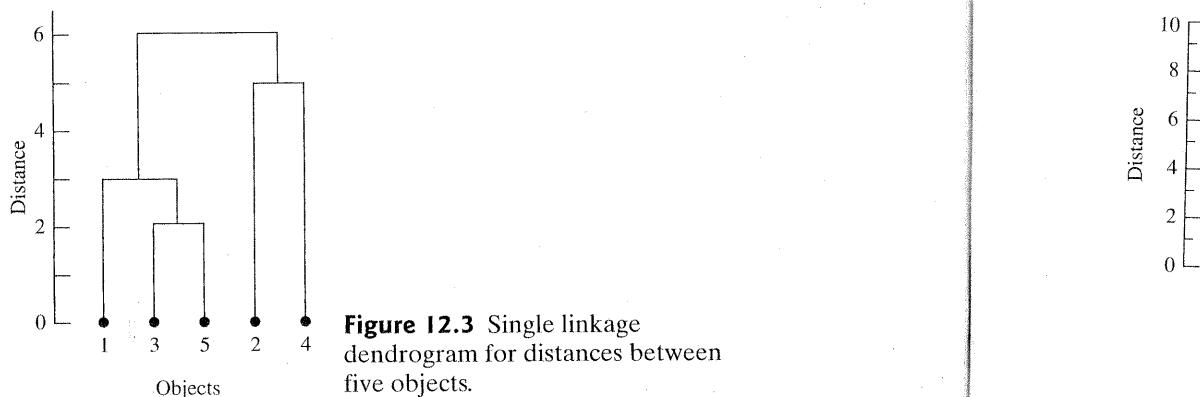
The final distance matrix becomes

$$(135) \quad \begin{matrix} (135) & (24) \\ 0 & \\ \textcircled{6} & 0 \end{matrix}$$

Consequently, clusters (135) and (24) are merged to form a single cluster of all five objects, (12345), when the nearest neighbor distance reaches 6.

The dendrogram picturing the hierarchical clustering just concluded is shown in Figure 12.3. The groupings and the distance levels at which they occur are clearly illustrated by the dendrogram. ■

In typical applications of hierarchical clustering, the intermediate results—where the objects are sorted into a moderate number of clusters—are of chief interest.



**Figure 12.3** Single linkage dendrogram for distances between five objects.

**Example 12.4 (Single linkage clustering of 11 languages)** Consider the array of concordances in Table 12.3 representing the closeness between the numbers 1–10 in 11 languages. To develop a matrix of distances, we subtract the concordance figure of 10 that each language has with itself. The subsequent assignments of distances are

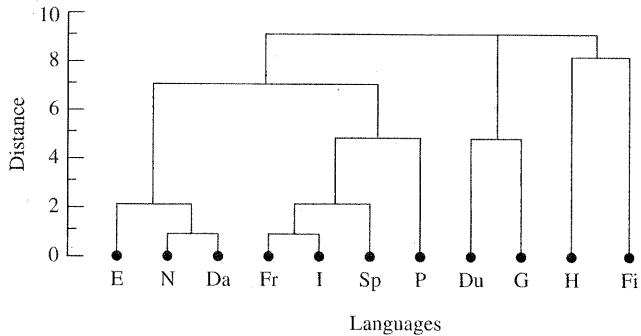
	E	N	Da	Du	G	Fr	Sp	I	P	H	Fi
E	0										
N	2	0									
Da	2	(1)	0								
Du	7	5	6	0							
G	6	4	5	5	0						
Fr	6	6	6	9	7	0					
Sp	6	6	5	9	7	2	0				
I	6	6	5	9	7	(1)	(1)	0			
P	7	7	6	10	8	5	3	4	0		
H	9	8	8	8	9	10	10	10	10	0	
Fi	9	9	9	9	9	9	9	9	8	0	

We first search for the minimum distance between pairs of languages (clusters). The minimum distance, 1, occurs between Danish and Norwegian, Italian and French, and Italian and Spanish. Numbering the languages in the order in which they appear across the top of the array, we have

$$d_{32} = 1; \quad d_{86} = 1; \quad \text{and } d_{87} = 1$$

Since  $d_{76} = 2$ , we can merge only clusters 8 and 6 or clusters 8 and 7. We cannot merge clusters 6, 7, and 8 at level 1. We choose first to merge 6 and 8, and then to update the distance matrix and merge 2 and 3 to obtain the clusters (68) and (23). Subsequent computer calculations produce the dendrogram in Figure 12.4.

From the dendrogram, we see that Norwegian and Danish, and also French and Italian, cluster at the minimum distance (maximum similarity) level. When the allowable distance is increased, English is added to the Norwegian–Danish group,



**Figure 12.4** Single linkage dendrograms for distances between numbers in 11 languages.

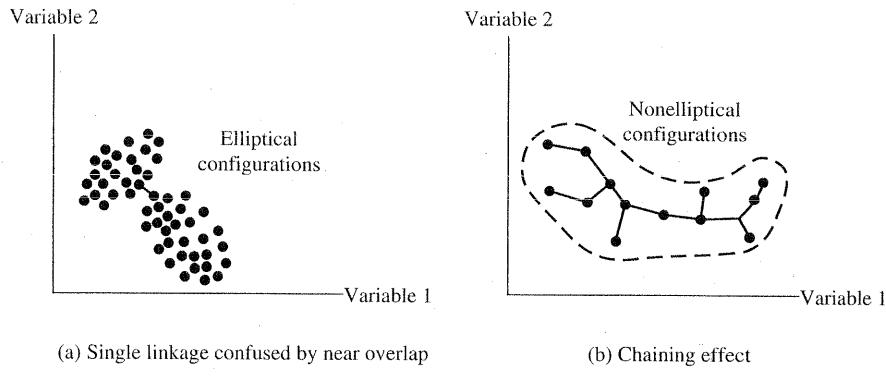
between

the array of numbers 1–10 in 11 ordinances from the f. The subsequent

Fi

and Spanish merges with the French–Italian group. Notice that Hungarian and Finnish are more similar to each other than to the other clusters of languages. However, these two clusters (languages) do not merge until the distance between nearest neighbors has increased substantially. Finally, all the clusters of languages are merged into a single cluster at the largest nearest neighbor distance, 9. ■

Since single linkage joins clusters by the shortest link between them, the technique cannot discern poorly separated clusters. [See Figure 12.5(a).] On the other hand, single linkage is one of the few clustering methods that can delineate nonelliptical clusters. The tendency of single linkage to pick out long stringlike clusters is known as *chaining*. [See Figure 12.5(b).] Chaining can be misleading if items at opposite ends of the chain are, in fact, quite dissimilar.



**Figure 12.5** Single linkage clusters.

The clusters formed by the single linkage method will be unchanged by any assignment of distance (similarity) that gives the same relative orderings as the initial distances (similarities). In particular, any one of a set of similarity coefficients from Table 12.1 that are monotonic to one another will produce the same clustering.

### Complete Linkage

Complete linkage clustering proceeds in much the same manner as single linkage clusterings, with one important exception: At each stage, the distance (similarity) between clusters is determined by the distance (similarity) between the two

and 7. We cannot and 8, and then to ters (68) and (23). gure 12.4. id also French and ) level. When the an–Danish group,

elements, one from each cluster, that are *most distant*. Thus, complete linkage ensures that all items in a cluster are within some maximum distance (or minimum similarity) of each other.

The general agglomerative algorithm again starts by finding the minimum entry in  $\mathbf{D} = \{d_{ik}\}$  and merging the corresponding objects, such as  $U$  and  $V$ , to get cluster  $(UV)$ . For Step 3 of the general algorithm in (12-12), the distances between  $(UV)$  and any other cluster  $W$  are computed by

$$d_{(UV)W} = \max \{d_{UW}, d_{VW}\} \quad (12-14)$$

Here  $d_{UW}$  and  $d_{VW}$  are the distances between the most distant members of clusters  $U$  and  $W$  and clusters  $V$  and  $W$ , respectively.

**Example 12.5 (Clustering using complete linkage)** Let us return to the distance matrix introduced in Example 12.3:

$$\mathbf{D} = \{d_{ik}\} = \begin{array}{c} \begin{matrix} & 1 & 2 & 3 & 4 & 5 \\ 1 & 0 & & & & \\ 2 & 9 & 0 & & & \\ 3 & 3 & 7 & 0 & & \\ 4 & 6 & 5 & 9 & 0 & \\ 5 & 11 & 10 & ② & 8 & 0 \end{matrix} \end{array}$$

At the first stage, objects 3 and 5 are merged, since they are most similar. This gives the cluster  $(35)$ . At stage 2, we compute

$$\begin{aligned} d_{(35)1} &= \max \{d_{31}, d_{51}\} = \max \{3, 11\} = 11 \\ d_{(35)2} &= \max \{d_{32}, d_{52}\} = 10 \\ d_{(35)4} &= \max \{d_{34}, d_{54}\} = 9 \end{aligned}$$

and the modified distance matrix becomes

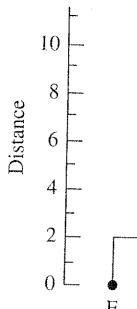
$$(35) \quad \begin{array}{c} (35) \quad 1 \quad 2 \quad 4 \\ \begin{bmatrix} 0 & & & \\ 11 & 0 & & \\ 10 & 9 & 0 & \\ 9 & 6 & ⑤ & 0 \end{bmatrix} \end{array}$$

The next merger occurs between the most similar groups, 2 and 4, to give the cluster  $(24)$ . At stage 3, we have

$$\begin{aligned} d_{(24)(35)} &= \max \{d_{2(35)}, d_{4(35)}\} = \max \{10, 9\} = 10 \\ d_{(24)1} &= \max \{d_{21}, d_{41}\} = 9 \end{aligned}$$

and the distance matrix

$$(35) \quad (24) \quad \begin{array}{c} (35) \quad (24) \quad 1 \\ \begin{bmatrix} 0 & & \\ 10 & 0 & \\ 11 & ⑨ & 0 \end{bmatrix} \end{array}$$



complete linkage  
ance (or minimum

the minimum entry  
and  $V$ , to get cluster  
ances between  $(UV)$

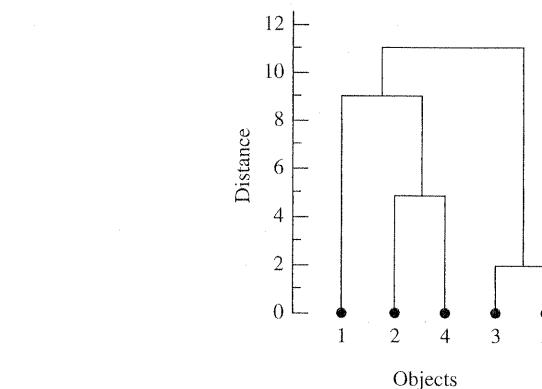
(12-14)

members of clusters

rn to the distance

ance (or minimum

t similar. This gives



**Figure 12.6** Complete linkage dendrogram for distances between five objects.

The next merger produces the cluster (124). At the final stage, the groups (35) and (124) are merged as the single cluster (12345) at level

$$d_{(124)(35)} = \max \{d_{1(35)}, d_{(24)(35)}\} = \max \{11, 10\} = 11$$

The dendrogram is given in Figure 12.6. ■

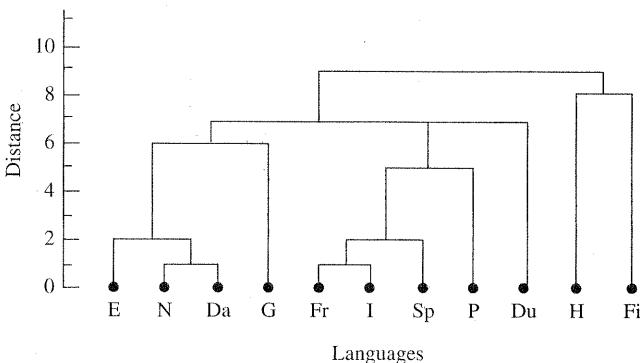
Comparing Figures 12.3 and 12.6, we see that the dendograms for single linkage and complete linkage differ in the allocation of object 1 to previous groups.

---

**Example 12.6 (Complete linkage clustering of 11 languages)** In Example 12.4, we presented a distance matrix for numbers in 11 languages. The complete linkage clustering algorithm applied to this distance matrix produces the dendrogram shown in Figure 12.7.

Comparing Figures 12.7 and 12.4, we see that both hierarchical methods yield the English–Norwegian–Danish and the French–Italian–Spanish language groups. Polish is merged with French–Italian–Spanish at an intermediate level. In addition, both methods merge Hungarian and Finnish only at the penultimate stage.

However, the two methods handle German and Dutch differently. Single linkage merges German and Dutch at an intermediate distance, and these two languages remain a cluster until the final merger. Complete linkage merges German



**Figure 12.7** Complete linkage dendrogram for distances between numbers in 11 languages.

with the English–Norwegian–Danish group at an intermediate level. Dutch remains a cluster by itself until it is merged with the English–Norwegian–Danish–German and French–Italian–Spanish–Polish groups at a higher distance level. The final complete linkage merger involves two clusters. The final merger in single linkage involves three clusters.

**Example 12.7 (Clustering variables using complete linkage)** Data collected on 22 U.S. public utility companies for the year 1975 are listed in Table 12.4. Although it is more interesting to group companies, we shall see here how the complete linkage algorithm can be used to cluster variables. We measure the similarity between pairs of

**Table 12.4** Public Utility Data (1975)

Company	Variables							
	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$
1. Arizona Public Service	1.06	9.2	151	54.4	1.6	9077	0.	.628
2. Boston Edison Co.	.89	10.3	202	57.9	2.2	5088	25.3	1.555
3. Central Louisiana Electric Co.	1.43	15.4	113	53.0	3.4	9212	0.	1.058
4. Commonwealth Edison Co.	1.02	11.2	168	56.0	.3	6423	34.3	.700
5. Consolidated Edison Co. (N.Y.)	1.49	8.8	192	51.2	1.0	3300	15.6	2.044
6. Florida Power & Light Co.	1.32	13.5	111	60.0	-2.2	11127	22.5	1.241
7. Hawaiian Electric Co.	1.22	12.2	175	67.6	2.2	7642	0.	1.652
8. Idaho Power Co.	1.10	9.2	245	57.0	3.3	13082	0.	.309
9. Kentucky Utilities Co.	1.34	13.0	168	60.4	7.2	8406	0.	.862
10. Madison Gas & Electric Co.	1.12	12.4	197	53.0	2.7	6455	39.2	.623
11. Nevada Power Co.	.75	7.5	173	51.5	6.5	17441	0.	.768
12. New England Electric Co.	1.13	10.9	178	62.0	3.7	6154	0.	1.897
13. Northern States Power Co.	1.15	12.7	199	53.7	6.4	7179	50.2	.527
14. Oklahoma Gas & Electric Co.	1.09	12.0	96	49.8	1.4	9673	0.	.588
15. Pacific Gas & Electric Co.	.96	7.6	164	62.2	-0.1	6468	.9	1.400
16. Puget Sound Power & Light Co.	1.16	9.9	252	56.0	9.2	15991	0.	.620
17. San Diego Gas & Electric Co.	.76	6.4	136	61.9	9.0	5714	8.3	1.920
18. The Southern Co.	1.05	12.6	150	56.7	2.7	10140	0.	1.108
19. Texas Utilities Co.	1.16	11.7	104	54.0	-2.1	13507	0.	.636
20. Wisconsin Electric Power Co.	1.20	11.8	148	59.9	3.5	7287	41.1	.702
21. United Illuminating Co.	1.04	8.6	204	61.0	3.5	6650	0.	2.116
22. Virginia Electric & Power Co.	1.07	9.3	174	54.3	5.9	10093	26.6	1.306

KEY:  $X_1$ : Fixed-charge coverage ratio (income/debt).

$X_2$ : Rate of return on capital.

$X_3$ : Cost per KW capacity in place.

$X_4$ : Annual load factor.

$X_5$ : Peak kWh demand growth from 1974 to 1975.

$X_6$ : Sales (kWh use per year).

$X_7$ : Percent nuclear.

$X_8$ : Total fuel costs (cents per kWh).

Source: Data courtesy of H. E. Thompson.

el. Dutch remains  
-Danish-German  
el. The final com-  
single linkage in-

a collected on 22  
2.4. Although it is  
mplete linkage al-  
✓ between pairs of

6	$X_7$	$X_8$
77	0.	.628
88	25.3	1.555
12	0.	1.058
23	34.3	.700
00	15.6	2.044
27	22.5	1.241
42	0.	1.652
82	0.	.309
06	0.	.862
55	39.2	.623
41	0.	.768
54	0.	1.897
79	50.2	.527
73	0.	.588
68	.9	1.400
91	0.	.620
14	8.3	1.920
40	0.	1.108
07	0.	.636
87	41.1	.702
50	0.	2.116
93	26.6	1.306

**Table 12.5** Correlations Between Pairs of Variables (Public Utility Data)

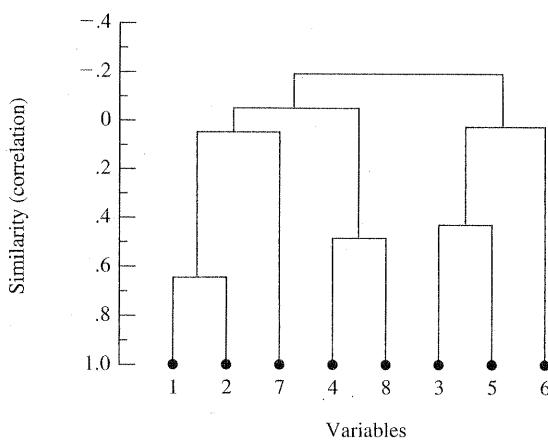
$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$
1.000							
.643	1.000						
-.103	-.348	1.000					
-.082	-.086	.100	1.000				
-.259	-.260	.435	.034	1.000			
-.152	-.010	.028	-.288	.176	1.000		
.045	.211	.115	-.164	-.019	-.374	1.000	
-.013	-.328	.005	.486	-.007	-.561	-.185	1.000

variables by the product-moment correlation coefficient. The correlation matrix is given in Table 12.5.

When the sample correlations are used as similarity measures, variables with large negative correlations are regarded as very dissimilar; variables with large positive correlations are regarded as very similar. In this case, the "distance" between clusters is measured as the *smallest* similarity between members of the corresponding clusters. The complete linkage algorithm, applied to the foregoing similarity matrix, yields the dendrogram in Figure 12.8.

We see that variables 1 and 2 (fixed-charge coverage ratio and rate of return on capital), variables 4 and 8 (annual load factor and total fuel costs), and variables 3 and 5 (cost per kilowatt capacity in place and peak kilowatthour demand growth) cluster at intermediate "similarity" levels. Variables 7 (percent nuclear) and 6 (sales) remain by themselves until the final stages. The final merger brings together the (12478) group and the (356) group. ■

As in single linkage, a "new" assignment of distances (similarities) that have the same relative orderings as the initial distances will not change the configuration of the complete linkage clusters.



**Figure 12.8** Complete linkage dendrogram for similarities among eight utility company variables.

## Average Linkage

Average linkage treats the distance between two clusters as the average distance between all pairs of items where one member of a pair belongs to each cluster.

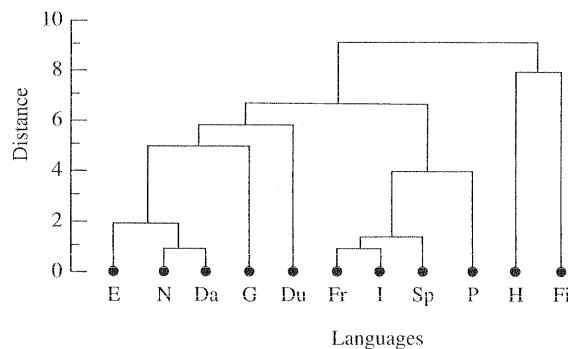
Again, the input to the average linkage algorithm may be distances or similarities, and the method can be used to group objects or variables. The average linkage algorithm proceeds in the manner of the general algorithm of (12-12). We begin by searching the distance matrix  $\mathbf{D} = \{d_{ik}\}$  to find the nearest (most similar) objects—for example,  $U$  and  $V$ . These objects are merged to form the cluster  $(UV)$ . For Step 3 of the general agglomerative algorithm, the distances between  $(UV)$  and the other cluster  $W$  are determined by

$$d_{(UV)W} = \frac{\sum_i \sum_k d_{ik}}{N_{(UV)} N_W} \quad (12-15)$$

where  $d_{ik}$  is the distance between object  $i$  in the cluster  $(UV)$  and object  $k$  in the cluster  $W$ , and  $N_{(UV)}$  and  $N_W$  are the number of items in clusters  $(UV)$  and  $W$ , respectively.

---

**Example 12.8 (Average linkage clustering of 11 languages)** The average linkage algorithm was applied to the “distances” between 11 languages given in Example 12.4. The resulting dendrogram is displayed in Figure 12.9.



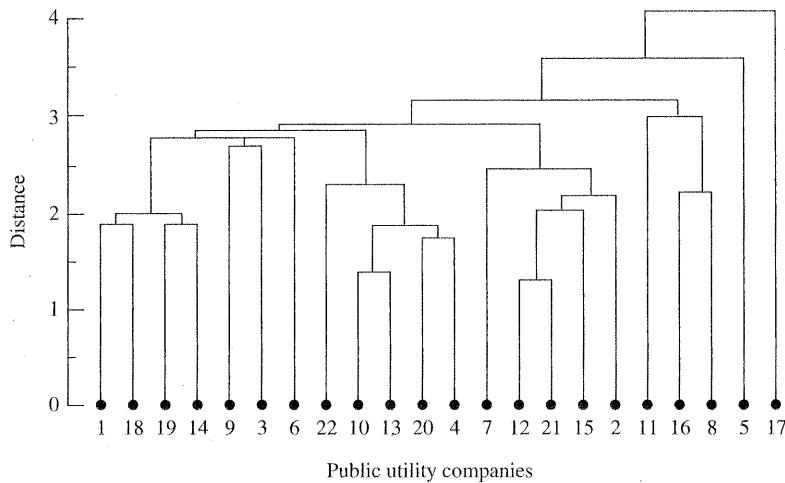
**Figure 12.9** Average linkage dendrogram for distances between numbers in 11 languages.

A comparison of the dendrogram in Figure 12.9 with the corresponding single linkage dendrogram (Figure 12.4) and complete linkage dendrogram (Figure 12.7) indicates that average linkage yields a configuration very much like the complete linkage configuration. However, because distance is defined differently for each case, it is not surprising that mergers take place at different levels. ■

---

**Example 12.9 (Average linkage clustering of public utilities)** An average linkage algorithm applied to the Euclidean distances between 22 public utilities (see Table 12.6) produced the dendrogram in Figure 12.10 on page 692.





**Figure 12.10** Average linkage dendrogram for distances between 22 public utility companies.

Concentrating on the intermediate clusters, we see that the utility companies tend to group according to geographical location. For example, one intermediate cluster contains the firms 1 (Arizona Public Service), 18 (The Southern Company—primarily Georgia and Alabama), 19 (Texas Utilities Company), and 14 (Oklahoma Gas and Electric Company). There are some exceptions. The cluster (7, 12, 21, 15, 2) contains firms on the eastern seaboard and in the far west. On the other hand, all these firms are located near the coasts. Notice that Consolidated Edison Company of New York and San Diego Gas and Electric Company stand by themselves until the final amalgamation stages.

It is, perhaps, not surprising that utility firms with similar locations (or types of locations) cluster. One would expect regulated firms in the same area to use, basically, the same type of fuel(s) for power plants and face common markets. Consequently, types of generation, costs, growth rates, and so forth should be relatively homogeneous among these firms. This is apparently reflected in the hierarchical clustering.

For average linkage clustering, changes in the assignment of distances (similarities) can affect the arrangement of the final configuration of clusters, even though the changes preserve relative orderings.

### Ward's Hierarchical Clustering Method

Ward [32] considered hierarchical clustering procedures based on minimizing the 'loss of information' from joining two groups. This method is usually implemented with loss of information taken to be an increase in an error sum of squares criterion,

ESS. First, for a given cluster  $k$ , let  $\text{ESS}_k$  be the sum of the squared deviations of every item in the cluster from the cluster mean (centroid). If there are currently  $K$  clusters, define ESS as the sum of the  $\text{ESS}_k$  or  $\text{ESS} = \text{ESS}_1 + \text{ESS}_2 + \dots + \text{ESS}_K$ . At each step in the analysis, the union of every possible pair of clusters is considered, and the two clusters whose combination results in the smallest increase in ESS (minimum loss of information) are joined. Initially, each cluster consists of a single item, and, if there are  $N$  items,  $\text{ESS}_k = 0$ ,  $k = 1, 2, \dots, N$ , so  $\text{ESS} = 0$ . At the other extreme, when all the clusters are combined in a single group of  $N$  items, the value of ESS is given by

$$\text{ESS} = \sum_{j=1}^N (\mathbf{x}_j - \bar{\mathbf{x}})'(\mathbf{x}_j - \bar{\mathbf{x}})$$

where  $\mathbf{x}_j$  is the multivariate measurement associated with the  $j$ th item and  $\bar{\mathbf{x}}$  is the mean of all the items.

The results of Ward's method can be displayed as a dendrogram. The vertical axis gives the values of ESS at which the mergers occur.

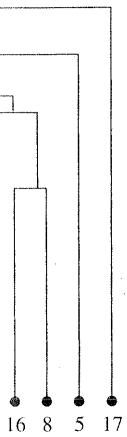
Ward's method is based on the notion that the clusters of multivariate observations are expected to be roughly elliptically shaped. It is a hierarchical precursor to nonhierarchical clustering methods that optimize some criterion for dividing data into a *given* number of elliptical groups. We discuss nonhierarchical clustering procedures in the next section. Additional discussion of optimization methods of cluster analysis is contained in [8].

---

**Example 12.10 (Clustering pure malt scotch whiskies)** Virtually all the world's pure malt Scotch whiskies are produced in Scotland. In one study (see [22]), 68 binary variables were created measuring characteristics of Scotch whiskey that can be broadly classified as color, nose, body, palate, and finish. For example, there were 14 color characteristics (descriptions), including white wine, yellow, very pale, pale, bronze, full amber, red, and so forth. LaPointe and Legendre clustered 109 pure malt Scotch whiskies, each from a different distillery. The investigators were interested in determining the major types of single-malt whiskies, their chief characteristics, and the best representative. In addition, they wanted to know whether the groups produced by the hierarchical clustering procedure corresponded to different geographical regions, since it is known that whiskies are affected by local soil, temperature, and water conditions.

Weighted similarity coefficients  $\{s_{ik}\}$  were created from binary variables representing the presence or absence of characteristics. The resulting "distances," defined as  $\{d_{ik} = 1 - s_{ik}\}$ , were used with Ward's method to group the 109 pure (single-) malt Scotch whiskies. The resulting dendrogram is shown in Figure 12.11. (An average linkage procedure applied to a similarity matrix produced almost exactly the same classification.)

The groups labelled A-L in the figure are the 12 groups of similar Scotches identified by the investigators. A follow-up analysis suggested that these 12 groups have a large geographic component in the sense that Scotches with similar characteristics tend to be produced by distilleries that are located reasonably



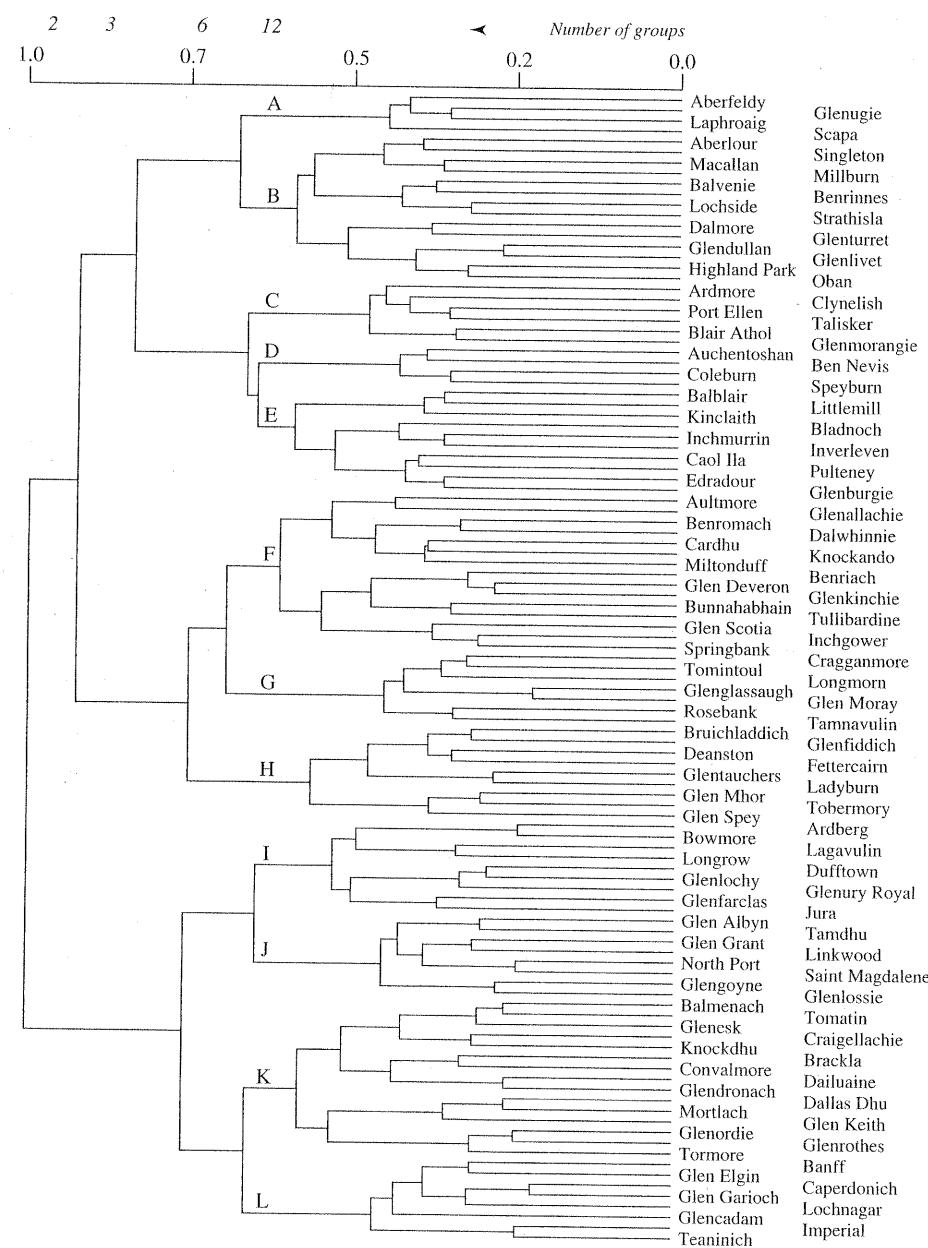
public utility

utility companies  
one intermediate  
thern Company—  
and 14 (Oklahoma  
ter (7, 12, 21, 15, 2)  
the other hand, all  
Edison Company  
y themselves until

ations (or types of  
e area to use, basi-  
n markets. Conse-  
ould be relatively  
in the hierarchical

distances (simili-  
sters, even though

on minimizing the  
ually implemented  
f squares criterion,



**Figure 12.11** A dendrogram for similarities between 109 pure malt Scotch whiskies.

close to one another. Consequently, the investigators concluded, "The relationship with geographic features was demonstrated, supporting the hypothesis that whiskies are affected not only by distillery secrets and traditions but also by factors dependent on region such as water, soil, microclimate, temperature and even air quality."

### Final Comments—Hierarchical Procedures

There are many agglomerative hierarchical clustering procedures besides single linkage, complete linkage, and average linkage. However, all the agglomerative procedures follow the basic algorithm of (12-12).

As with most clustering methods, sources of error and variation are not formally considered in hierarchical procedures. This means that a clustering method will be sensitive to outliers, or "noise points."

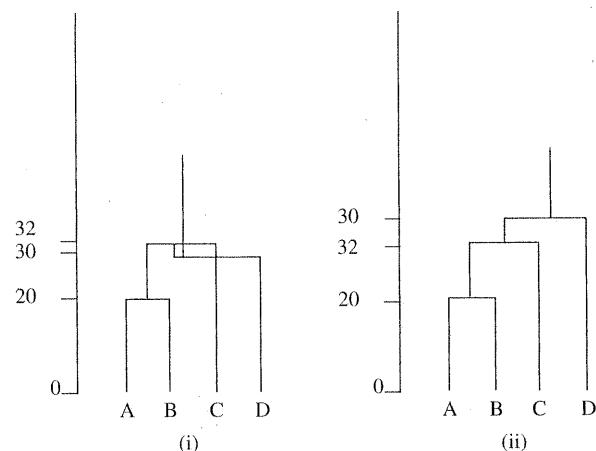
In hierarchical clustering, there is no provision for a reallocation of objects that may have been "incorrectly" grouped at an early stage. Consequently, the final configuration of clusters should always be carefully examined to see whether it is sensible.

For a particular problem, it is a good idea to try several clustering methods and, within a given method, a couple different ways of assigning distances (similarities). If the outcomes from the several methods are (roughly) consistent with one another, perhaps a case for "natural" groupings can be advanced.

The *stability* of a hierarchical solution can sometimes be checked by applying the clustering algorithm before and after *small* errors (perturbations) have been added to the data units. If the groups are fairly well distinguished, the clusterings before perturbation and after perturbation should agree.

Common values (ties) in the similarity or distance matrix can produce multiple solutions to a hierarchical clustering problem. That is, the dendograms corresponding to different treatments of the tied similarities (distances) can be different, particularly at the lower levels. This is not an inherent problem of any method; rather, multiple solutions occur for certain kinds of data. Multiple solutions are not necessarily bad, but the user needs to know of their existence so that the groupings (dendograms) can be properly interpreted and different groupings (dendograms) compared to assess their overlap. A further discussion of this issue appears in [27].

Some data sets and hierarchical clustering methods can produce *inversions*. (See [27].) An inversion occurs when an object joins an existing cluster at a smaller distance (greater similarity) than that of a previous consolidation. An inversion is represented two different ways in the following diagram:



Glenugie  
 Scapa  
 Singleton  
 Millburn  
 Benrinnes  
 Strathisla  
 Glenturret  
 Glenlivet  
 Oban  
 Clynelish  
 Talisker  
 Glenmorangie  
 Ben Nevis  
 Speyburn  
 Littlemill  
 Bladnoch  
 Inverleven  
 Pulteney  
 Glenburgie  
 Glenlachie  
 Dalwhinnie  
 Knockando  
 Benriach  
 Glenkinchie  
 Tullibardine  
 Inchgower  
 Cragganmore  
 Longmorn  
 Glen Moray  
 Tamnavulin  
 Glenfiddich  
 Fettercairn  
 Ladyburn  
 Tobermory  
 Ardbeg  
 Lagavulin  
 Dufftown  
 Glenury Royal  
 Jura  
 Tamduh  
 Linkwood  
 Saint Magdalene  
 Glenlossie  
 Tomatin  
 Craigellachie  
 Brackla  
 Dailhuaine  
 Dallas Dhu  
 Glen Keith  
 Glenrothes  
 Banff  
 Caperdonich  
 Lochnagar  
 Imperial  
 Scotch

"The relationship is hypothesis that is but also by factorperature and even

In this example, the clustering method joins A and B at distance 20. At the next step, C is added to the group (AB) at distance 32. Because of the nature of the clustering algorithm, D is added to group (ABC) at distance 30, a smaller distance than the distance at which C joined (AB). In (i) the inversion is indicated by a dendrogram with crossover. In (ii), the inversion is indicated by a dendrogram with a non-monotonic scale.

Inversions can occur when there is no clear cluster structure and are generally associated with two hierarchical clustering algorithms known as the centroid method and the median method. The hierarchical procedures discussed in this book are not prone to inversions.

## 12.4 Nonhierarchical Clustering Methods

Nonhierarchical clustering techniques are designed to group *items*, rather than *variables*, into a collection of  $K$  clusters. The number of clusters,  $K$ , may either be specified in advance or determined as part of the clustering procedure. Because a matrix of distances (similarities) does not have to be determined, and the basic data do not have to be stored during the computer run, nonhierarchical methods can be applied to much larger data sets than can hierarchical techniques.

Nonhierarchical methods start from either (1) an initial partition of items into groups or (2) an initial set of seed points, which will form the nuclei of clusters. Good choices for starting configurations should be free of overt biases. One way to start is to randomly select seed points from among the items or to randomly partition the items into initial groups.

In this section, we discuss one of the more popular nonhierarchical procedures, the  $K$ -means method.

### *K*-means Method

MacQueen [25] suggests the term *K-means* for describing an algorithm of his that assigns each item to the cluster having the nearest centroid (mean). In its simplest version, the process is composed of these three steps:

1. Partition the items into  $K$  initial clusters.
2. Proceed through the list of items, assigning an item to the cluster whose centroid (mean) is nearest. (Distance is usually computed using Euclidean distance with either standardized or unstandardized observations.) Recalculate the centroid for the cluster receiving the new item and for the cluster losing the item.
3. Repeat Step 2 until no more reassignments take place. (12-16)

Rather than starting with a partition of all items into  $K$  preliminary groups in Step 1, we could specify  $K$  initial centroids (seed points) and then proceed to Step 2.

The final assignment of items to clusters will be, to some extent, dependent upon the initial partition or the initial selection of seed points. Experience suggests that most major changes in assignment occur with the first reallocation step.

ice 20. At the next nature of the cluster-aller distance than ated by a dendrogram with a non-

and are generally used as the centroid discussed in this book

s, rather than *variables*, may either be specified. Because a matrix of basic data do not have centroids can be applied

ition of items into nuclei of clusters. biases. One way to to randomly parti-

chical procedures,

gorithm of his that an). In its simplest

ter whose centroid mean distance with calculate the centroid g the item.

(12-16)

reliminary groups d then proceed to

extent, dependent xperience suggests cation step.

**Example 12.11 (Clustering using the K-means method)** Suppose we measure two variables  $X_1$  and  $X_2$  for each of four items  $A, B, C$ , and  $D$ . The data are given in the following table:

Item	Observations	
	$x_1$	$x_2$
A	5	3
B	-1	1
C	1	-2
D	-3	-2

The objective is to divide these items into  $K = 2$  clusters such that the items within a cluster are closer to one another than they are to the items in different clusters. To implement the  $K = 2$ -means method, we arbitrarily partition the items into two clusters, such as  $(AB)$  and  $(CD)$ , and compute the coordinates  $(\bar{x}_1, \bar{x}_2)$  of the cluster centroid (mean). Thus, at Step 1, we have

Cluster	Coordinates of centroid	
	$\bar{x}_1$	$\bar{x}_2$
$(AB)$	$\frac{5 + (-1)}{2} = 2$	$\frac{3 + 1}{2} = 2$
$(CD)$	$\frac{1 + (-3)}{2} = -1$	$\frac{-2 + (-2)}{2} = -2$

At Step 2, we compute the Euclidean distance of each item from the group centroids and reassign each item to the nearest group. If an item is moved from the initial configuration, the cluster centroids (means) must be updated before proceeding. The  $i$ th coordinate,  $i = 1, 2, \dots, p$ , of the centroid is easily updated using the formulas:

$$\bar{x}_{i,new} = \frac{n\bar{x}_i + x_{ji}}{n+1} \quad \text{if the } j\text{th item is added to a group}$$

$$\bar{x}_{i,new} = \frac{n\bar{x}_i - x_{ji}}{n-1} \quad \text{if the } j\text{th item is removed from a group}$$

Here  $n$  is the number of items in the “old” group with centroid  $\bar{x}' = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)$ .

Consider the initial clusters  $(AB)$  and  $(CD)$ . The coordinates of the centroids are  $(2, 2)$  and  $(-1, -2)$  respectively. Suppose item  $A$  with coordinates  $(5, 3)$  is moved to the  $(CD)$  group. The new groups are  $(B)$  and  $(ACD)$  with updated centroids:

$$\text{Group } (B) \quad \bar{x}_{1,new} = \frac{2(2) - 5}{2-1} = -1 \quad \bar{x}_{2,new} = \frac{2(2) - 3}{2-1} = 1, \text{ the coordinates of } B$$

$$\text{Group } (ACD) \quad \bar{x}_{1,new} = \frac{2(-1) + 5}{2+1} = 1 \quad \bar{x}_{2,new} = \frac{2(-2) + 3}{2+1} = -.33$$

Returning to the initial groupings in Step 1, we compute the squared distances

$$\begin{aligned} d^2(A, (AB)) &= (5 - 2)^2 + (3 - 2)^2 = 10 && \text{if } A \text{ is not moved} \\ d^2(A, (CD)) &= (5 + 1)^2 + (3 + 2)^2 = 61 \\ d^2(A, (B)) &= (5 + 1)^2 + (3 - 1)^2 = 40 && \text{if } A \text{ is moved to the } (CD) \text{ group} \\ d^2(A, (ACD)) &= (5 - 1)^2 + (3 + .33)^2 = 27.09 \end{aligned}$$

Since  $A$  is closer to the center of  $(AB)$  than it is to the center of  $(ACD)$ , it is not reassigned.

Continuing, we consider reassigning  $B$ . We get

$$\begin{aligned} d^2(B, (AB)) &= (-1 - 2)^2 + (1 - 2)^2 = 10 && \text{if } B \text{ is not moved} \\ d^2(B, (CD)) &= (-1 + 1)^2 + (1 + 2)^2 = 9 \\ d^2(B, (A)) &= (-1 - 5)^2 + (1 - 3)^2 = 40 && \text{if } B \text{ is moved to the } (CD) \text{ group} \\ d^2(B, (BCD)) &= (-1 + 1)^2 + (1 + 1)^2 = 4 \end{aligned}$$

Since  $B$  is closer to the center of  $(BCD)$  than it is to the center of  $(AB)$ ,  $B$  is reassigned to the  $(CD)$  group. We now have the clusters  $(A)$  and  $(BCD)$  with centroid coordinates  $(5, 3)$  and  $(-1, -1)$  respectively.

We check  $C$  for reassignment.

$$\begin{aligned} d^2(C, (A)) &= (1 - 5)^2 + (-2 - 3)^2 = 41 && \text{if } C \text{ is not moved} \\ d^2(C, (BCD)) &= (1 + 1)^2 + (-2 + 1)^2 = 5 \\ d^2(C, (AC)) &= (1 - 3)^2 + (-2 - .5)^2 = 10.25 && \text{if } C \text{ is moved to the } (A) \text{ group} \\ d^2(C, (BD)) &= (1 + 2)^2 + (-2 + .5)^2 = 11.25 \end{aligned}$$

Since  $C$  is closer to the center of the  $BCD$  group than it is to the center of the  $AC$  group,  $C$  is not moved. Continuing in this way, we find that no more reassessments take place and the final  $K = 2$  clusters are  $(A)$  and  $(BCD)$ .

For the final clusters, we have

		Squared distances to group centroids			
		Item			
Cluster		$A$	$B$	$C$	$D$
$A$		0	40	41	89
$(BCD)$		52	4	5	5

The within cluster sum of squares (sum of squared distances to centroid) are

Cluster  $A$ : 0

Cluster  $(BCD)$ :  $4 + 5 + 5 = 14$

Equivalently, we can determine the  $K = 2$  clusters by using the criterion

$$\min E = \sum d_{i,c(i)}^2$$

ired distances

ved

to the  $(CD)$  group

of  $(ACD)$ , it is not

ved

to the  $(CD)$  group

of  $(AB)$ ,  $B$  is reas-  
 $CD$ ) with centroid

ved

to the  $(A)$  group

e center of the  $AC$   
ore reassessments

D  
39  
5

centroid) are

riterion

where the minimum is over the number of  $K = 2$  clusters and  $d_{i,c(i)}^2$  is the squared distance of case  $i$  from the centroid (mean) of the assigned cluster.

In this example, there are seven possibilities for  $K = 2$  clusters:

- $A, (BCD)$
- $B, (ACD)$
- $C, (ABD)$
- $D, (ABC)$
- $(AB), (CD)$
- $(AC), (BD)$
- $(AD), (BC)$

For the  $A, (BCD)$  pair:

$$\begin{array}{ll} A & d_{A,c(A)}^2 = 0 \\ (BCD) & d_{B,c(B)}^2 + d_{C,c(C)}^2 + d_{D,c(D)}^2 = 4 + 5 + 5 = 14 \end{array}$$

Consequently,  $\sum d_{i,c(i)}^2 = 0 + 14 = 14$

For the remaining pairs, you may verify that

$$\begin{array}{ll} B, (ACD) & \sum d_{i,c(i)}^2 = 48.7 \\ C, (ABD) & \sum d_{i,c(i)}^2 = 27.7 \\ D, (ABC) & \sum d_{i,c(i)}^2 = 31.3 \\ (AB), (CD) & \sum d_{i,c(i)}^2 = 28 \\ (AC), (BD) & \sum d_{i,c(i)}^2 = 27 \\ (AD), (BC) & \sum d_{i,c(i)}^2 = 51.3 \end{array}$$

Since the smallest  $\sum d_{i,c(i)}^2$  occurs for the pair of clusters  $(A)$  and  $(BCD)$ , this is the final partition. ■

To check the stability of the clustering, it is desirable to rerun the algorithm with a new initial partition. Once clusters are determined, intuitions concerning their interpretations are aided by rearranging the list of items so that those in the first cluster appear first, those in the second cluster appear next, and so forth. A table of the cluster centroids (means) and within-cluster variances also helps to delineate group differences.

---

**Example 12.12 ( $K$ -means clustering of public utilities)** Let us return to the problem of clustering public utilities using the data in Table 12.4. The  $K$ -means algorithm for several choices of  $K$  was run. We present a summary of the results for  $K = 4$  and  $K = 5$ . In general, the choice of a particular  $K$  is not clear cut and depends upon subject-matter knowledge, as well as data-based appraisals. (Data-based appraisals might include choosing  $K$  so as to maximize the between-cluster variability relative

to the within-cluster variability. Relevant measures might include  $|\mathbf{W}|/|\mathbf{B} + \mathbf{W}|$  [see (6-38)] and  $\text{tr}(\mathbf{W}^{-1}\mathbf{B})$ .) The summary is as follows:

$K = 4$

Cluster	Number of firms	Firms
1	5	{ Idaho Power Co. (8), Nevada Power Co. (11), Puget Sound Power & Light Co. (16), Virginia Electric & Power Co. (22), Kentucky Utilities Co. (9).
2	6	{ Central Louisiana Electric Co. (3), Oklahoma Gas & Electric Co. (14), The Southern Co. (18), Texas Utilities Co. (19), Arizona Public Service (1), Florida Power & Light Co. (6).
3	5	{ New England Electric Co. (12), Pacific Gas & Electric Co. (15), San Diego Gas & Electric Co. (17), United Illuminating Co. (21), Hawaiian Electric Co. (7).
4	6	{ Consolidated Edison Co. (N.Y.) (5), Boston Edison Co. (2), Madison Gas & Electric Co. (10), Northern States Power Co. (13), Wisconsin Electric Power Co. (20), Commonwealth Edison Co. (4).

Distances between Cluster Centers

$$\begin{array}{ccccc} & & 1 & 2 & 3 & 4 \\ & 1 & \left[ \begin{array}{ccccc} 0 & & & & \\ 3.08 & 0 & & & \\ 3.29 & 3.56 & 0 & & \\ 3.05 & 2.84 & 3.18 & 0 & \end{array} \right] \\ & 2 & & & & \\ & 3 & & & & \\ & 4 & & & & \end{array}$$

$K = 5$

Cluster	Number of firms	Firms
1	5	{ Nevada Power Co. (11), Puget Sound Power & Light Co. (16), Idaho Power Co. (8), Virginia Electric & Power Co. (22), Kentucky Utilities Co. (9).
2	6	{ Central Louisiana Electric Co. (3), Texas Utilities Co. (19), Oklahoma Gas & Electric Co. (14), The Southern Co. (18), Arizona Public Service (1), Florida Power & Light Co. (6).
3	5	{ New England Electric Co. (12), Pacific Gas & Electric Co. (15), San Diego Gas & Electric Co. (17), United Illuminating Co. (21), Hawaiian Electric Co. (7).
4	2	{ Consolidated Edison Co. (N.Y.) (5), Boston Edison Co. (2).
5	4	{ Commonwealth Edison Co. (4), Madison Gas & Electric Co. (10), Northern States Power Co. (13), Wisconsin Electric Power Co. (20).

ie  $|\mathbf{W}|/|\mathbf{B} + \mathbf{W}|$

---



---

Puget  
tric &

a Gas & Electric  
es Co. (19),  
Light Co. (6).  
z Electric

tric Co. (7).  
Edison Co.  
ern States

---



---

Light  
& Power Co.

es Co. (19),  
rn Co.  
& Light Co. (6).

Electric  
ited  
( ).

Electric Co. (10),  
ric Power Co. (20).

---

Distances between Cluster Centers

	1	2	3	4	5
1	0				
2	3.08	0			
3	3.29	3.56	0		
4	3.63	3.46	2.63	0	
5	3.18	2.99	3.81	2.89	0

The cluster profiles ( $K = 5$ ) shown in Figure 12.12 order the eight variables according to the ratios of their between-cluster variability to their within-cluster variability. [For univariate  $F$ -ratios, see Section 6.4.] We have

$$F_{\text{nuc}} = \frac{\text{mean square percent nuclear between clusters}}{\text{mean square percent nuclear within clusters}} = \frac{3.335}{.255} = 13.1$$

so firms within different clusters are widely separated with respect to percent nuclear, but firms within the same cluster show little percent nuclear variation. Fuel costs (FUEL) and annual sales (SALES) also seem to be of some importance in distinguishing the clusters.

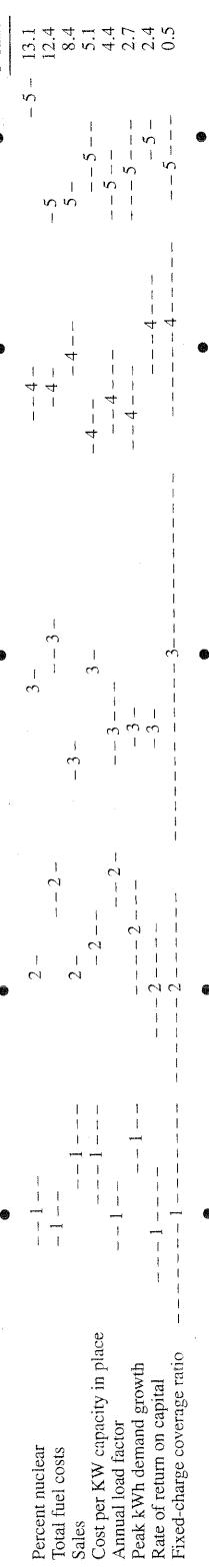
Reviewing the firms in the five clusters, it is apparent that the  $K$ -means method gives results generally consistent with the average linkage hierarchical method. (See Example 12.9.) Firms with common or compatible geographical locations cluster. Also, the firms in a given cluster seem to be roughly the same in terms of percent nuclear.

We must caution, as we have throughout the book, that the importance of *individual* variables in clustering must be judged from a multivariate perspective. All of the variables (multivariate observations) determine the cluster means and the reassignment of items. In addition, the values of the descriptive statistics measuring the importance of individual variables are functions of the number of clusters and the final configuration of the clusters. On the other hand, descriptive measures can be helpful, after the fact, in assessing the "success" of the clustering procedure.

### Final Comments—Nonhierarchical Procedures

There are strong arguments for not fixing the number of clusters,  $K$ , in advance, including the following:

1. If two or more seed points inadvertently lie within a single cluster, their resulting clusters will be poorly differentiated.

Cluster profiles—variables are ordered by *F*-ratio size

Each column describes a cluster.

The cluster number is printed at the mean of each variable.

Dashes indicate one standard deviation above and below mean.

Figure 12.12 Cluster profiles ( $K = 5$ ) for public utility data.

2. The existence of an outlier might produce at least one group with very disperse items.
3. Even if the population is known to consist of  $K$  groups, the sampling method may be such that data from the rarest group do not appear in the sample. Forcing the data into  $K$  groups would lead to nonsensical clusters.

In cases in which a single run of the algorithm requires the user to specify  $K$ , it is always a good idea to rerun the algorithm for several choices.

Discussions of other nonhierarchical clustering procedures are available in [3], [8], and [16].

## 12.5 Clustering Based on Statistical Models

The popular clustering methods discussed earlier in this chapter, including single linkage, complete linkage, average linkage, Ward's method and  $K$ -means clustering, are intuitively reasonable procedures but that is as much as we can say without having a model to explain how the observations were produced. Major advances in clustering methods have been made through the introduction of statistical models that indicate how the collection of  $(p \times 1)$  measurements  $\mathbf{x}_j$ , from the  $N$  objects, was generated. The most common model is one where cluster  $k$  has expected proportion  $p_k$  of the objects and the corresponding measurements are generated by a probability density function  $f_k(\mathbf{x})$ . Then, if there are  $K$  clusters, the observation vector for a single object is modeled as arising from the *mixing distribution*

$$f_{Mix}(\mathbf{x}) = \sum_{k=1}^K p_k f_k(\mathbf{x})$$

where each  $p_k \geq 0$  and  $\sum_{k=1}^K p_k = 1$ . This distribution  $f_{Mix}(\mathbf{x})$  is called a mixture of the  $K$  distributions  $f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_K(\mathbf{x})$  because the observation is generated from the component distribution  $f_k(\mathbf{x})$  with probability  $p_k$ . The collection of  $N$  observation vectors generated from this distribution will be a mixture of observations from the component distributions.

The most common mixture model is a mixture of multivariate normal distributions where the  $k$ -th component  $f_k(\mathbf{x})$  is the  $N_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  density function.

The normal mixture model for one observation  $\mathbf{x}$  is

$$\begin{aligned} f_{Mix}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) \\ = \sum_{k=1}^K p_k \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right) \end{aligned} \quad (12-17)$$

Clusters generated by this model are ellipsoidal in shape with the heaviest concentration of observations near the center.

Inferences are based on the likelihood, which for  $N$  objects and a fixed number of clusters  $K$ , is

$$\begin{aligned} L(p_1, \dots, p_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K) &= \prod_{j=1}^N f_{Mix}(\mathbf{x}_j | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K) \\ &= \prod_{j=1}^N \left( \sum_{k=1}^K p_k \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x}_j - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_k) \right) \right) \end{aligned} \quad (12-18)$$

where the proportions  $p_1, \dots, p_k$ , the mean vectors  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ , and the covariance matrices  $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k$  are unknown. The measurements for different objects are treated as independent and identically distributed observations from the mixture distribution.

There are typically far too many unknown parameters for parameters for making inferences when the number of objects to be clustered is at least moderate. However, certain conclusions can be made regarding situations where a heuristic clustering method should work well. In particular, the likelihood based procedure under the normal mixture model with all  $\boldsymbol{\Sigma}_k$  the same multiple of the identity matrix,  $\eta \mathbf{I}$ , is approximately the same as  $K$ -means clustering and Ward's method. To date, no statistical models have been advanced for which the cluster formation procedure is approximately the same as single linkage, complete linkage or average linkage.

Most importantly, under the sequence of mixture models (12-17) for different  $K$ , the problems of choosing the number of clusters and choosing an appropriate clustering method has been reduced to the problem of selecting an appropriate statistical model. This is a major advance.

A good approach to selecting a model is to first obtain the maximum likelihood estimates  $\hat{p}_1, \dots, \hat{p}_K, \hat{\boldsymbol{\mu}}_1, \dots, \hat{\boldsymbol{\mu}}_K, \hat{\boldsymbol{\Sigma}}_1, \dots, \hat{\boldsymbol{\Sigma}}_K$  for a fixed number of clusters  $K$ . These estimates must be obtained numerically using special purpose software. The resulting value of the maximum of the likelihood

$$L_{\max} = L(\hat{p}_1, \dots, \hat{p}_K, \hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\Sigma}}_1, \dots, \hat{\boldsymbol{\mu}}_K, \hat{\boldsymbol{\Sigma}}_K)$$

provides the basis for model selection. How do we decide on a reasonable value for the number of clusters  $K$ ? In order to compare models with different numbers of parameters, a penalty is subtracted from twice the maximized value of the log-likelihood to give

$$-2 \ln L_{\max} - \text{Penalty}$$

where the penalty depends on the number of parameters estimated and the number of observations  $N$ . Since the probabilities  $p_k$  sum to 1, there are only  $K - 1$  probabilities that must be estimated,  $K \times p$  means and  $K \times p(p + 1)/2$  variances and covariances. For the Akaike information criterion (AIC), the penalty is  $2N \times (\text{number of parameters})$  so

$$\text{AIC} = 2 \ln L_{\max} - 2N \left( K \frac{1}{2} (p + 1)(p + 2) - 1 \right) \quad (12-19)$$

Assumed for  
for  $\boldsymbol{\Sigma}_k$   
 $\boldsymbol{\Sigma}_k = \eta \mathbf{I}$   
 $\boldsymbol{\Sigma}_k = \eta_k \mathbf{I}$   
 $\boldsymbol{\Sigma}_k = \eta_k \text{ Dia}$

nd a fixed number

,  $\Sigma_K$ )

$j - \mu_k)\right)$  (12-18)

nd the covariance  
erent objects are  
from the mixture

rameters for making at least moderate  
where a heuristic  
based procedure  
le of the identity  
d Ward's method.  
cluster formation  
inkage or average

2-17) for different  
ng an appropriate  
n appropriate sta-

iximum likelihood  
usters  $K$ . These es-  
vare. The resulting

asonable value for  
different numbers  
ized value of the

ed and the number  
only  $K - 1$  proba-  
)/2 variances and  
, the penalty is

1) (12-19)

The Bayesian information criterion (BIC) is similar but uses the logarithm of the number of parameters in the penalty function

$$\text{BIC} = 2 \ln L_{\max} - 2 \ln(N) \left( K \frac{1}{2} (p + 1)(p + 2) - 1 \right) \quad (12-20)$$

There is still occasional difficulty with too many parameters in the mixture model so simple structures are assumed for the  $\Sigma_k$ . In particular, progressively more complicated structures are allowed as indicated in the following table.

Assumed form for $\Sigma_k$	Total number of parameters	BIC
$\Sigma_k = \eta \mathbf{I}$	$K(p + 1)$	$\ln L_{\max} - 2 \ln(N)K(p + 1)$
$\Sigma_k = \eta_k \mathbf{I}$	$K(p + 2) - 1$	$\ln L_{\max} - 2 \ln(N)(K(p + 2) - 1)$
$\Sigma_k = \eta_k \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$	$K(p + 2) + p - 1$	$\ln L_{\max} - 2 \ln(N)(K(p + 2) + p - 1)$

Additional structures for the covariance matrices are considered in [6] and [9].

Even for a fixed number of clusters, the estimation of a mixture model is complicated. One current software package, *MCLUST*, available in the R software library, combines hierarchical clustering, the EM algorithm and the BIC criterion to develop an appropriate model for clustering. In the 'E'-step of the EM algorithm, a  $(N \times K)$  matrix is created whose  $j$ th row contains estimates of the conditional (on the current parameter estimates) probabilities that observation  $\mathbf{x}_j$  belongs to cluster 1, 2, ...,  $K$ . So, at convergence, the  $j$ th observation (object) is assigned to the cluster  $k$  for which the conditional probability

$$p(k | \mathbf{x}_j) = \hat{p}_j f(\mathbf{x}_j | k) / \sum_{i=1}^K \hat{p}_i f(\mathbf{x}_i | k)$$

of membership is the largest. (See [6] and [9] and the references therein.)

**Example 12.13 (A model based clustering of the iris data)** Consider the Iris data in Table 11.5. Using *MCLUST* and specifically the *me* function, we first fit the  $p = 4$  dimensional normal mixture model restricting the covariance matrices to satisfy  $\Sigma_k = \eta_k \mathbf{I}$ ,  $k = 1, 2, 3$ .

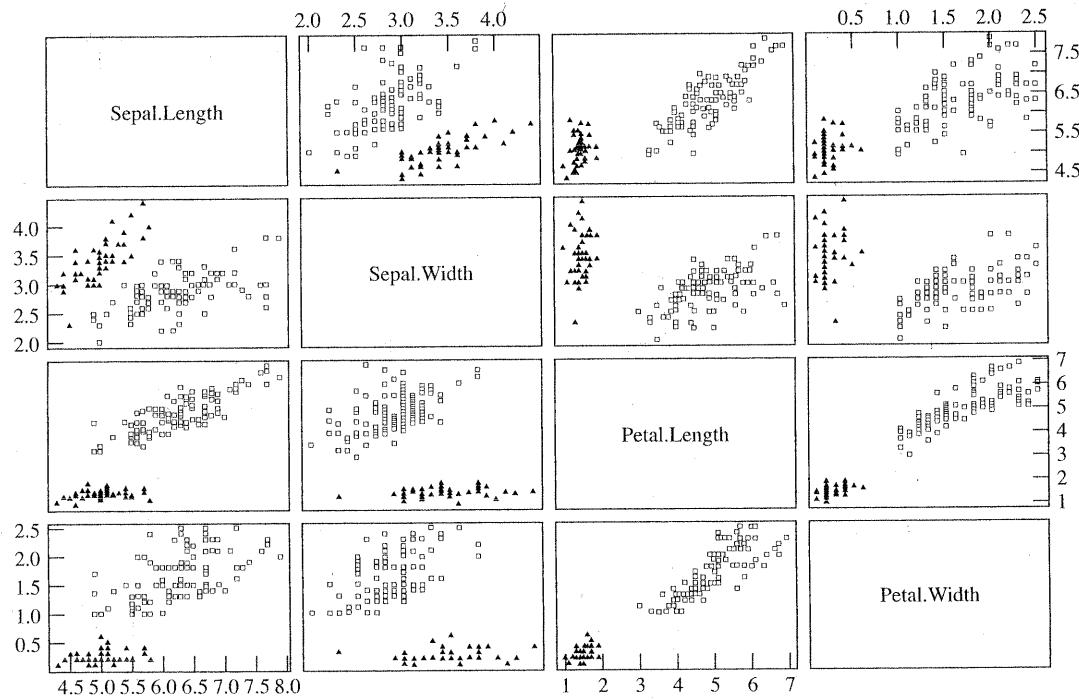
Using the BIC criterion, the software chooses  $K = 3$  clusters with estimated centers

$$\boldsymbol{\mu}_1 = \begin{bmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{bmatrix}, \quad \boldsymbol{\mu}_2 = \begin{bmatrix} 5.90 \\ 2.75 \\ 4.40 \\ 1.43 \end{bmatrix}, \quad \boldsymbol{\mu}_3 = \begin{bmatrix} 6.85 \\ 3.07 \\ 5.73 \\ 2.07 \end{bmatrix},$$

and estimated variance-covariance scale factors  $\hat{\eta}_1 = .076$ ,  $\hat{\eta}_2 = .163$  and  $\hat{\eta}_3 = .163$ . The estimated mixing proportions are  $\hat{p}_1 = .3333$ ,  $\hat{p}_2 = .4133$  and  $\hat{p}_3 = .2534$ . For this solution,  $\text{BIC} = -853.8$ . A matrix plot of the clusters for pairs of variables is shown in Figure 12.13.

Once we have an estimated mixture model, a new object  $\mathbf{x}_j$  will be assigned to the cluster for which the conditional probability of membership is the largest (see [9]).

Assuming the  $\Sigma_k = \eta_k \mathbf{I}$  covariance structure and allowing up to  $K = 7$  clusters, the BIC can be increased to  $\text{BIC} = -705.1$ .

**Figure 12.13** Multiple scatter plots of  $K = 3$  clusters for Iris data

Finally, using the BIC criterion with up to  $K = 9$  groups and several different covariance structures, the best choice is a two group mixture model with unconstrained covariances. The estimated mixing probabilities are  $\hat{p}_1 = .3333$  and  $\hat{p}_2 = .6667$ . The estimated group centers are

$$\boldsymbol{\mu}_1 = \begin{bmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{bmatrix}, \quad \boldsymbol{\mu}_2 = \begin{bmatrix} 6.26 \\ 2.87 \\ 4.91 \\ 1.68 \end{bmatrix}$$

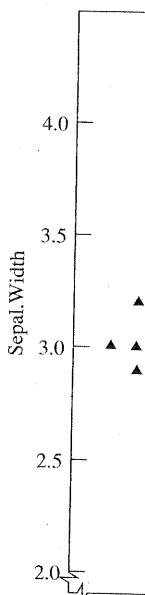
and the two estimated covariance matrices are

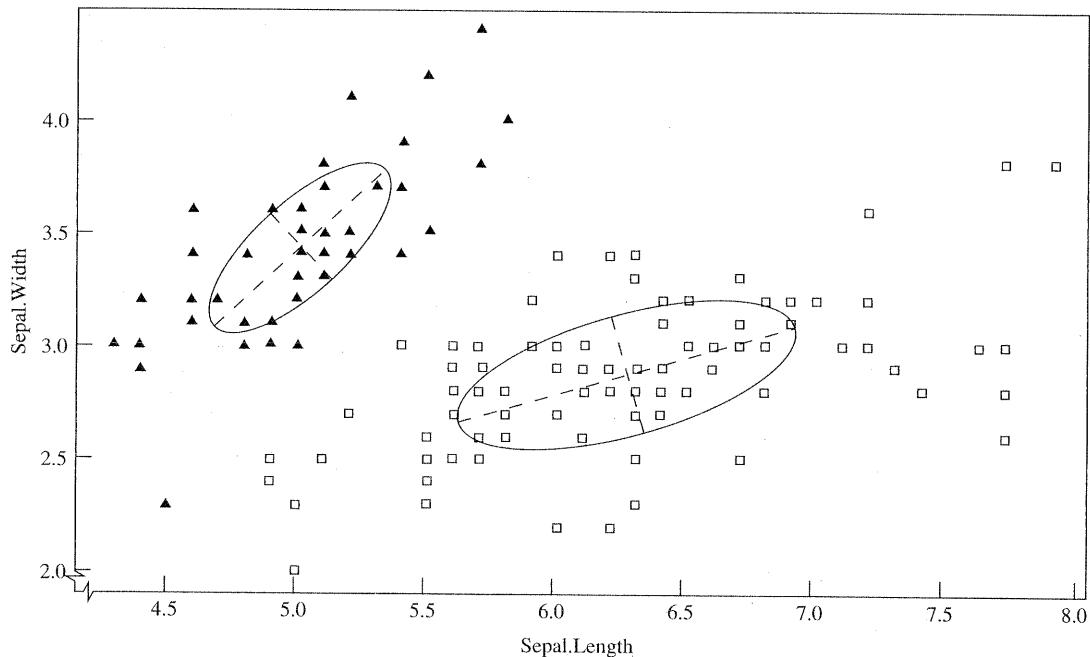
$$\hat{\Sigma}_1 = \begin{bmatrix} .1218 & .0972 & .0160 & .0101 \\ .0972 & .1408 & .0115 & .0091 \\ .0160 & .0115 & .0296 & .0059 \\ .0101 & .0091 & .0059 & .0109 \end{bmatrix} \quad \hat{\Sigma}_2 = \begin{bmatrix} .4530 & .1209 & .4489 & .1655 \\ .1209 & .1096 & .1414 & .0792 \\ .4489 & .1414 & .6748 & .2858 \\ .1655 & .0792 & .2858 & .1786 \end{bmatrix}$$

Essentially, two species of Iris have been put in the same cluster as the projected view of the scatter plot of the sepal measurements in Figure 12.14 shows. ■

## 12.6 Multidimensional Scaling

This section begins a discussion of methods for displaying (transformed) multivariate data in low-dimensional space. We have already considered this issue when we

**Figure 12.14**



**Figure 12.14** Scatter plot of sepal measurements for best model.

discussed plotting scores on, say, the first two principal components or the scores on the first two linear discriminants. The methods we are about to discuss differ from these procedures in the sense that their *primary* objective is to “fit” the original data into a low-dimensional coordinate system such that any distortion caused by a reduction in dimensionality is minimized. Distortion generally refers to the similarities or dissimilarities (distances) among the original data points. Although Euclidean distance may be used to measure the closeness of points in the final low-dimensional configuration, the notion of similarity or dissimilarity depends upon the underlying technique for its definition. A low-dimensional plot of the kind we are alluding to is called an *ordination* of the data.

Multidimensional scaling techniques deal with the following problem: For a set of observed similarities (or distances) between every pair of  $N$  items, find a representation of the items in few dimensions such that the interitem proximities “nearly match” the original similarities (or distances).

It may not be possible to match exactly the ordering of the original similarities (distances). Consequently, scaling techniques attempt to find configurations in  $q \leq N - 1$  dimensions such that the match is as close as possible. The numerical measure of closeness is called the *stress*.

It is possible to arrange the  $N$  items in a low-dimensional coordinate system using only the *rank orders* of the  $N(N - 1)/2$  original similarities (distances), and not their magnitudes. When only this ordinal information is used to obtain a geometric representation, the process is called *nonmetric multidimensional scaling*. If the actual magnitudes of the original similarities (distances) are used to obtain a geometric representation in  $q$  dimensions, the process is called *metric multidimensional scaling*. Metric multidimensional scaling is also known as *principal coordinate analysis*.

Scaling techniques were developed by Shepard (see [29] for a review of early work), Kruskal [19, 20, 21], and others. A good summary of the history, theory, and applications of multidimensional scaling is contained in [35]. Multidimensional scaling invariably requires the use of a computer, and several good computer programs are now available for the purpose.

### The Basic Algorithm

For  $N$  items, there are  $M = N(N - 1)/2$  similarities (distances) between pairs of different items. These similarities constitute the basic data. (In cases where the similarities cannot be easily quantified as, for example, the similarity between two colors, the rank orders of the similarities are the basic data.)

Assuming no ties, the similarities can be arranged in a strictly ascending order as

$$s_{i_1 k_1} < s_{i_2 k_2} < \cdots < s_{i_M k_M} \quad (12-21)$$

Here  $s_{i_1 k_1}$  is the smallest of the  $M$  similarities. The subscript  $i_1 k_1$  indicates the pair of items that are least similar—that is, the items with rank 1 in the similarity ordering. Other subscripts are interpreted in the same manner. We want to find a  $q$ -dimensional configuration of the  $N$  items such that the distances,  $d_{i k}^{(q)}$ , between pairs of items match the ordering in (12-21). If the distances are laid out in a manner corresponding to that ordering, a perfect match occurs when

$$d_{i_1 k_1}^{(q)} > d_{i_2 k_2}^{(q)} > \cdots > d_{i_M k_M}^{(q)} \quad (12-22)$$

That is, the descending ordering of the distances in  $q$  dimensions is exactly analogous to the ascending ordering of the initial similarities. As long as the order in (12-22) is preserved, the magnitudes of the distances are unimportant.

For a given value of  $q$ , it may not be possible to find a configuration of points whose pairwise distances are monotonically related to the original similarities. Kruskal [19] proposed a measure of the extent to which a geometrical representation falls short of a perfect match. This measure, the stress, is defined as

$$\text{Stress } (q) = \left\{ \frac{\sum_{i < k} (d_{i k}^{(q)} - \hat{d}_{i k}^{(q)})^2}{\sum_{i < k} [d_{i k}^{(q)}]^2} \right\}^{1/2} \quad (12-23)$$

The  $\hat{d}_{i k}^{(q)}$ 's in the stress formula are numbers known to satisfy (12-22); that is, they are monotonically related to the similarities. The  $\hat{d}_{i k}^{(q)}$ 's are *not* distances in the sense that they satisfy the usual distance properties of (1-25). They are merely reference numbers used to judge the nonmonotonicity of the observed  $d_{i k}^{(q)}$ 's.

The idea is to find a representation of the items as points in  $q$ -dimensions such that the stress is as small as possible. Kruskal [19] suggests the stress be informally interpreted according to the following guidelines:

Stress	Goodness of fit	
20%	Poor	
10%	Fair	
5%	Good	
2.5%	Excellent	
0%	Perfect	

(12-24)

*Goodness of fit* refers to the monotonic relationship between the similarities and the final distances.

a review of early history, theory, and Multidimensional scaling computer

) between pairs of es where the similarity between two col-

ascending order as  
(12-21)

indicates the pair in the similarity. We want to find a ces,  $d_{ik}^{(q)}$ , between id out in a manner

(12-22)

s is exactly analo-  
ig as the order in  
tant.  
guration of points  
ignal similarities.  
etrical representa-  
ied as

(12-23)

2-22); that is, they  
tances in the sense  
merely reference  
s.  
 $q$ -dimensions such  
ress be informally

(12-24)

A second measure of discrepancy, introduced by Takane et al. [31], is becoming the preferred criterion. For a given dimension  $q$ , this measure, denoted by SSStress, replaces the  $d_{ik}$ 's and  $\hat{d}_{ik}$ 's in (12-23) by their squares and is given by

$$\text{SSStress} = \left[ \frac{\sum \sum (d_{ik}^2 - \hat{d}_{ik}^2)^2}{\sum \sum d_{ik}^4} \right]^{1/2} \quad (12-25)$$

The value of SSStress is always between 0 and 1. Any value less than .1 is typically taken to mean that there is a good representation of the objects by the points in the given configuration.

Once items are located in  $q$  dimensions, their  $q \times 1$  vectors of coordinates can be treated as multivariate observations. For display purposes, it is convenient to represent this  $q$ -dimensional scatter plot in terms of its principal component axes. (See Chapter 8.)

We have written the stress measure as a function of  $q$ , the number of dimensions for the geometrical representation. For each  $q$ , the configuration leading to the minimum stress can be obtained. As  $q$  increases, minimum stress will, within rounding error, decrease and will be zero for  $q = N - 1$ . Beginning with  $q = 1$ , a plot of these stress ( $q$ ) numbers versus  $q$  can be constructed. The value of  $q$  for which this plot begins to level off may be selected as the "best" choice of the dimensionality. That is, we look for an "elbow" in the stress-dimensionality plot.

The entire multidimensional scaling algorithm is summarized in these steps:

1. For  $N$  items, obtain the  $M = N(N - 1)/2$  similarities (distances) between distinct pairs of items. Order the similarities as in (12-21). (Distances are ordered from largest to smallest.) If similarities (distances) cannot be computed, the rank orders must be specified.
2. Using a trial configuration in  $q$  dimensions, determine the interitem distances  $d_{ik}^{(q)}$  and numbers  $\hat{d}_{ik}^{(q)}$ , where the latter satisfy (12-22) and minimize the stress (12-23) or SSStress (12-25). (The  $\hat{d}_{ik}^{(q)}$  are frequently determined within scaling computer programs using regression methods designed to produce monotonic "fitted" distances.)
3. Using the  $\hat{d}_{ik}^{(q)}$ 's, move the points around to obtain an improved configuration. (For  $q$  fixed, an improved configuration is determined by a general function minimization procedure applied to the stress. In this context, the stress is regarded as a function of the  $N \times q$  coordinates of the  $N$  items.) A new configuration will have new  $d_{ik}^{(q)}$ 's, new  $\hat{d}_{ik}^{(q)}$ 's, and smaller stress. The process is repeated until the best (minimum stress) representation is obtained.
4. Plot minimum stress ( $q$ ) versus  $q$  and choose the best number of dimensions,  $q^*$ , from an examination of this plot.

We have assumed that the initial similarity values are symmetric ( $s_{ik} = s_{ki}$ ), that there are no ties, and that there are no missing observations. Kruskal [19, 20] has suggested methods for handling asymmetries, ties, and missing observations. In addition, there are now multidimensional scaling computer programs that will handle not only Euclidean distance, but any distance of the Minkowski type. [See (12-3).]

The next three examples illustrate multidimensional scaling with distances as the initial (dis)similarity measures.

---

**Example 12.14 (Multidimensional scaling of U.S. cities)** Table 12.7 displays the airline distances between pairs of selected U.S. cities.

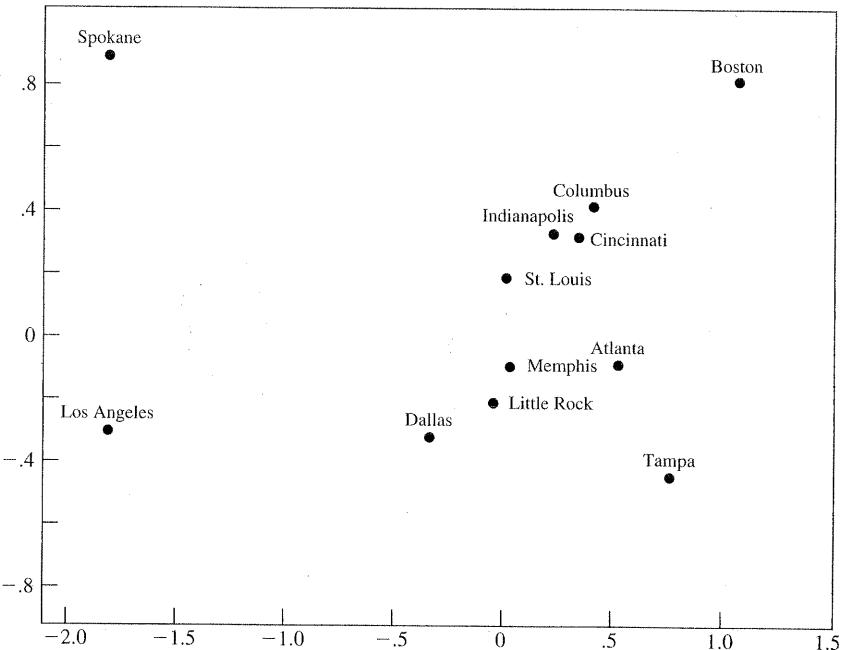
cur  
 $q =$   
 in t  
 sim  
 The  
 stre  
 sing  
 He

sior  
 Thi  
 with

**Exa**  
 the  
 spac  
 dist  
 proc

**Table 12.7** Airline-Distance Data

	Atlanta	Boston	Cincinnati	Columbus	Dallas	Indianapolis	Little Rock	Los Angeles	Memphis	St. Louis	Spokane	Tampa
(1)	0											
(2)	1068	0										
(3)	461	867	0									
(4)	549	769	107	0								
(5)	805	1819	943	1050	0							
(6)	508	941	108	172	882	0						
(7)	505	1494	618	725	325	562	0					
(8)	2197	3052	2186	2245	1403	2080	1701	0				
(9)	366	1355	502	586	464	436	137	1831	0			
(10)	558	1178	338	409	645	234	353	1848	294	0		
(11)	2467	2747	2067	2131	1891	1959	1988	1227	2042	1820	0	
(12)	467	1379	928	985	1077	975	912	2480	779	1016	2821	0



**Figure 12.15** A geometrical representation of cities produced by multidimensional scaling.

Since the cities naturally lie in a two-dimensional space (a nearly level part of the curved surface of the earth), it is not surprising that multidimensional scaling with  $q = 2$  will locate these items about as they occur on a map. Note that if the distances in the table are ordered from largest to smallest—that is, from a least similar to most similar—the first position is occupied by  $d_{\text{Boston}, \text{L.A.}} = 3052$ .

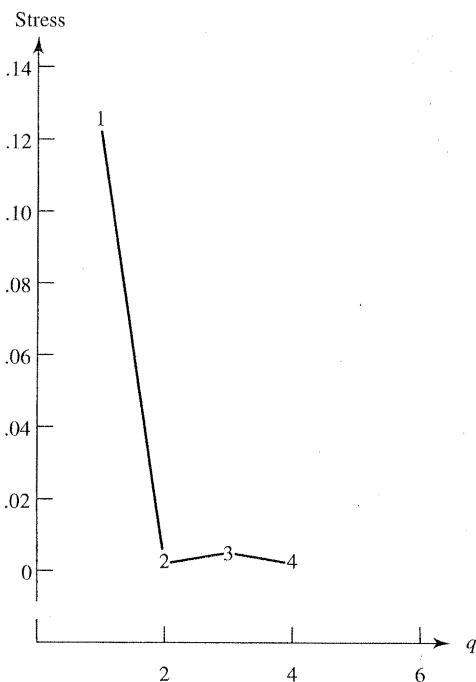
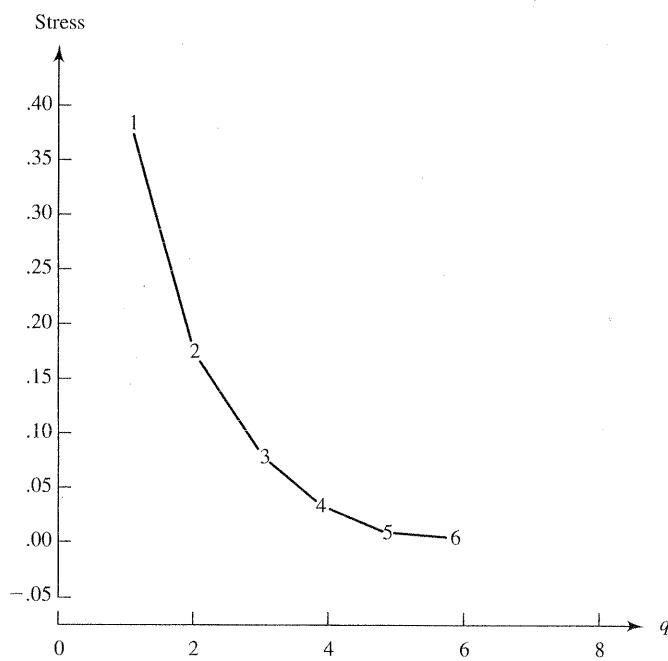
A multidimensional scaling plot for  $q = 2$  dimensions is shown in Figure 12.15. The axes lie along the sample principal components of the scatter plot.

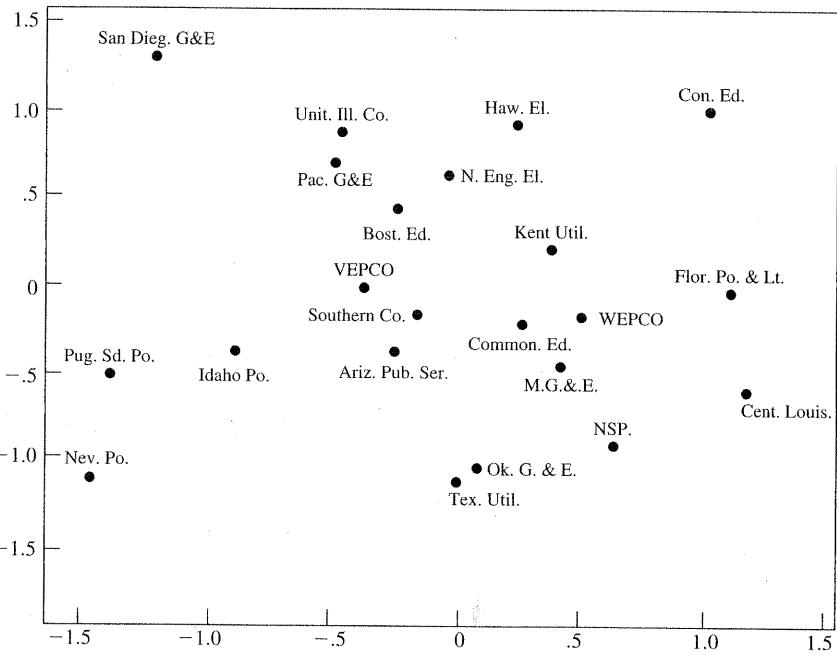
A plot of stress ( $q$ ) versus  $q$  is shown in Figure 12.16 on page 712. Since stress (1)  $\times 100\% = 12\%$ , a representation of the cities in one dimension (along a single axis) is not unreasonable. The “elbow” of the stress function occurs at  $q = 2$ . Here stress (2)  $\times 100\% = 0.8\%$ , and the “fit” is almost perfect.

The plot in Figure 12.16 indicates that  $q = 2$  is the best choice for the dimension of the final configuration. Note that the stress actually increases for  $q = 3$ . This anomaly can occur for extremely small values of stress because of difficulties with the numerical search procedure used to locate the minimum stress. ■

---

**Example 12.15 (Multidimensional scaling of public utilities)** Let us try to represent the 22 public utility firms discussed in Example 12.7 as points in a low-dimensional space. The measures of (dis)similarities between pairs of firms are the Euclidean distances listed in Table 12.6. Multidimensional scaling in  $q = 1, 2, \dots, 6$  dimensions produced the stress function shown in Figure 12.17.

**Figure 12.16** Stress function for airline distances between cities.**Figure 12.17** Stress function for distances between utilities.**Fig**out  
q =  
fici  
me;  
fin19%  
the  
util  
Gas  
(sin  
clos  
locdist:  
flex:  
trod**Exa**  
univ  
give



**Figure 12.18** A geometrical representation of utilities produced by multidimensional scaling.

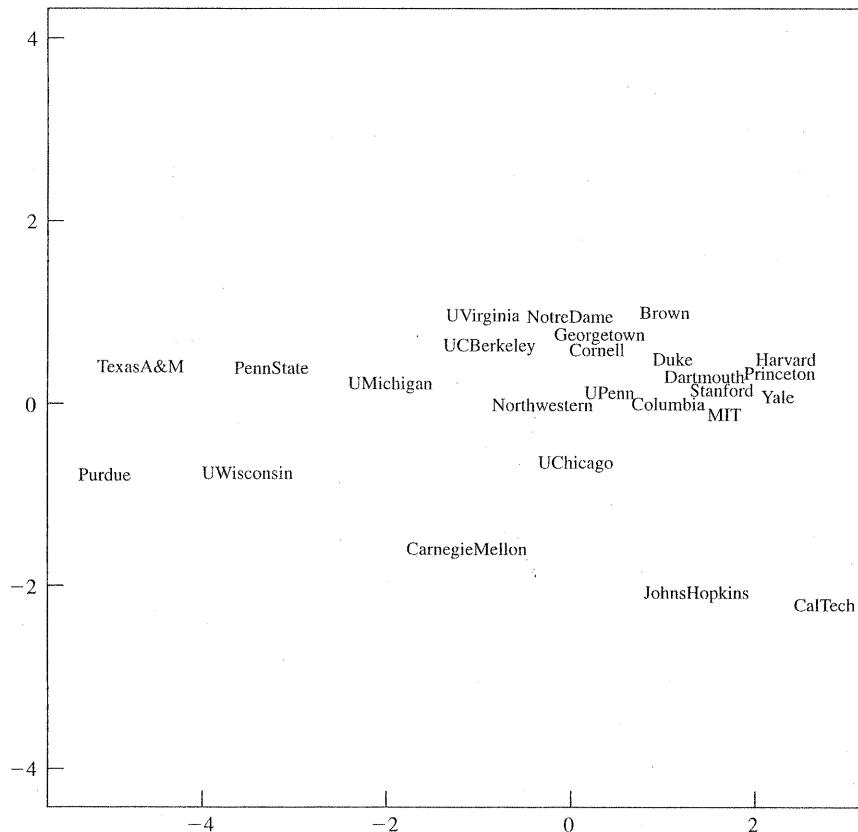
The stress function in Figure 12.17 has no sharp elbow. The plot appears to level out at “good” values of stress (less than or equal to 5%) in the neighborhood of  $q = 4$ . A good four-dimensional representation of the utilities is achievable, but difficult to display. We show a plot of the utility configuration obtained in  $q = 2$  dimensions in Figure 12.18. The axes lie along the sample principal components of the final scatter.

Although the stress for two dimensions is rather high (stress (2)  $\times 100\% = 19\%$ ), the distances between firms in Figure 12.18 are not wildly inconsistent with the clustering results presented earlier in this chapter. For example, the midwest utilities—Commonwealth Edison, Wisconsin Electric Power (WEPCO), Madison Gas and Electric (MG & E), and Northern States Power (NSP)—are close together (similar). Texas Utilities and Oklahoma Gas and Electric (Ok. G & E) are also very close together (similar). Other utilities tend to group according to geographical locations or similar environments.

The utilities cannot be positioned in two dimensions such that the interutility distances  $d_{ik}^{(2)}$  are entirely consistent with the original distances in Table 12.6. More flexibility for positioning the points is required, and this can only be obtained by introducing additional dimensions. ■

---

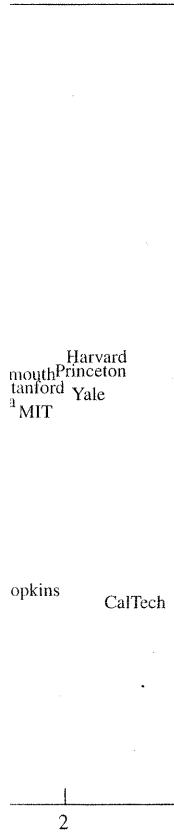
**Example 12.16 (Multidimensional scaling of universities)** Data related to 25 U.S. universities are given in Table 12.9 on page 729. (See Example 12.19.) These data give the average SAT score of entering freshmen, percent of freshmen in top



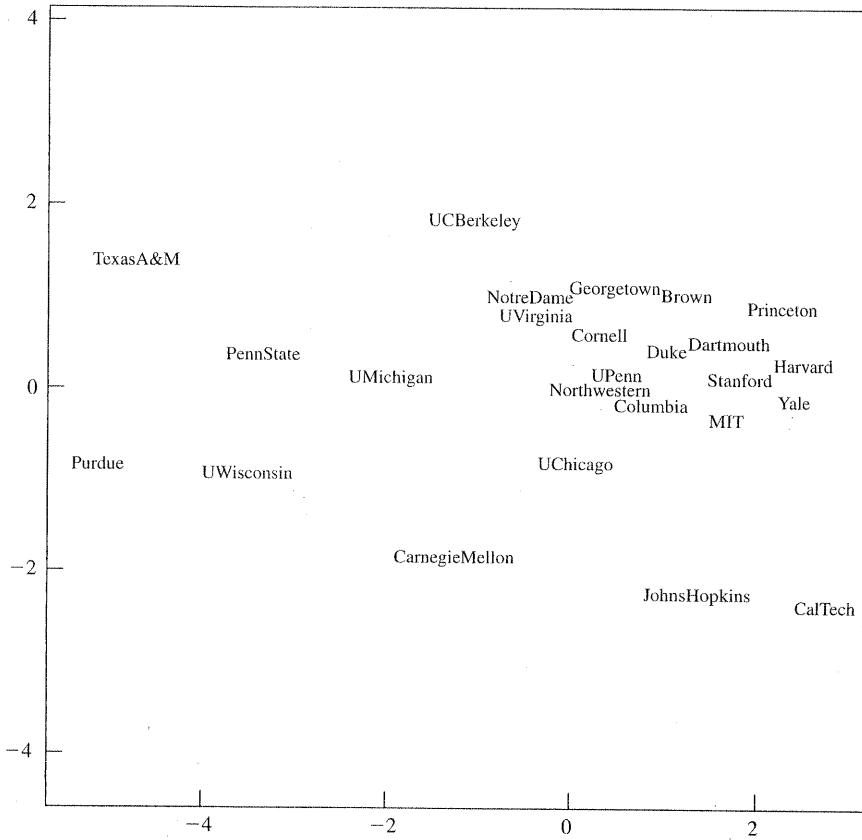
**Figure 12.19** A two-dimensional representation of universities produced by metric multidimensional scaling.

10% of high school class, percent of applicants accepted, student-faculty ratio, estimated annual expense, and graduation rate (%). A metric multidimensional scaling algorithm applied to the standardized university data gives the two-dimensional representation shown in Figure 12.19. Notice how the private universities cluster on the right of the plot while the large public universities are, generally, on the left. A nonmetric multidimensional scaling two-dimensional configuration is shown in Figure 12.20. For this example, the metric and nonmetric scaling representations are very similar, with the two dimensional stress value being approximately 10% for both scalings.

Classical metric scaling, or principal coordinate analysis, is equivalent to plotting the principal components. Different software programs choose the signs of the appropriate eigenvectors differently, so at first sight, two solutions may appear to be different. However, the solutions will coincide with a reflection of one or more of the axes. (See [26].)



duced by metric



**Figure 12.20** A two-dimensional representation of universities produced by nonmetric multidimensional scaling.

faculty ratio, estimated by metric. Two-dimensional representations of universities cluster together, generally, on the left. This representation is shown in Figure 12.20. The representations are approximately 10% equivalent to plotting the signs of the approximate distances. The signs of the approximate distances may appear to be the same for one or more of the universities.

To summarize, the key objective of multidimensional scaling procedures is a low-dimensional picture. Whenever multivariate data can be presented graphically in two or three dimensions, visual inspection can greatly aid interpretations.

When the multivariate observations are naturally numerical, and Euclidean distances in  $p$ -dimensions,  $d_{ik}^{(p)}$ , can be computed, we can seek a  $q < p$ -dimensional representation by minimizing

$$E = \left[ \sum_{i < k} \sum_{j < i} (d_{ij}^{(p)} - d_{ij}^{(q)})^2 / d_{ij}^{(p)} \right] \left[ \sum_{i < k} \sum_{j < i} d_{ij}^{(p)} \right]^{-1} \quad (12-27)$$

In this alternative approach, the Euclidean distances in  $p$  and  $q$  dimensions are compared directly. Techniques for obtaining low-dimensional representations by minimizing  $E$  are called *nonlinear mappings*.

The final goodness of fit of any low-dimensional representation can be depicted graphically by *minimal spanning trees*. (See [16] for a further discussion of these topics.)

## 12.7 Correspondence Analysis

Developed by the French, correspondence analysis is a graphical procedure for representing associations in a table of frequencies or counts. We will concentrate on a two-way table of frequencies or *contingency table*. If the contingency table has  $I$  rows and  $J$  columns, the plot produced by correspondence analysis contains two sets of points: A set of  $I$  points corresponding to the rows and a set of  $J$  points corresponding to the columns. The positions of the points reflect associations.

Row points that are close together indicate rows that have similar profiles (conditional distributions) across the columns. Column points that are close together indicate columns with similar profiles (conditional distributions) down the rows. Finally, row points that are close to column points represent combinations that occur more frequently than would be expected from an independence model—that is, a model in which the row categories are unrelated to the column categories.

The usual output from a correspondence analysis includes the “best” two-dimensional representation of the data, along with the coordinates of the plotted points, and a measure (called the *inertia*) of the amount of information retained in each dimension.

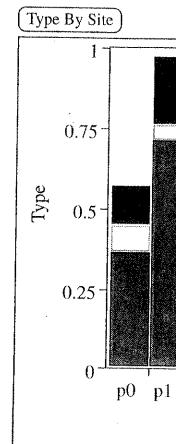
Before briefly discussing the algebraic development of correspondence analysis, it is helpful to illustrate the ideas we have introduced with an example.

**Example 12.17 (Correspondence analysis of archaeological data)** Table 12.8 contains the frequencies (counts) of  $J = 4$  different types of pottery (called potsherds) found at  $I = 7$  archaeological sites in an area of the American Southwest. If we divide the frequencies in each row (archaeological site) by the corresponding row total, we obtain a profile of types of pottery. The profiles for the different sites (rows) are shown in a bar graph in Figure 12.21(a). The widths of the bars are proportional to the total row frequencies. In general, the profiles are different; however, the profiles for sites P1 and P2 are similar, as are the profiles for sites P4 and P5.

The archaeological site profile for different types of pottery (columns) are shown in a bar graph in Figure 12.21(b). The site profiles are constructed using the

Site	Type				Total
	A	B	C	D	
P0	30	10	10	39	89
P1	53	4	16	2	75
P2	73	1	41	1	116
P3	20	6	1	4	31
P4	46	36	37	13	132
P5	45	6	59	10	120
P6	16	28	169	5	218
Total	283	91	333	74	781

Source: Data courtesy of M. J. Tretter.



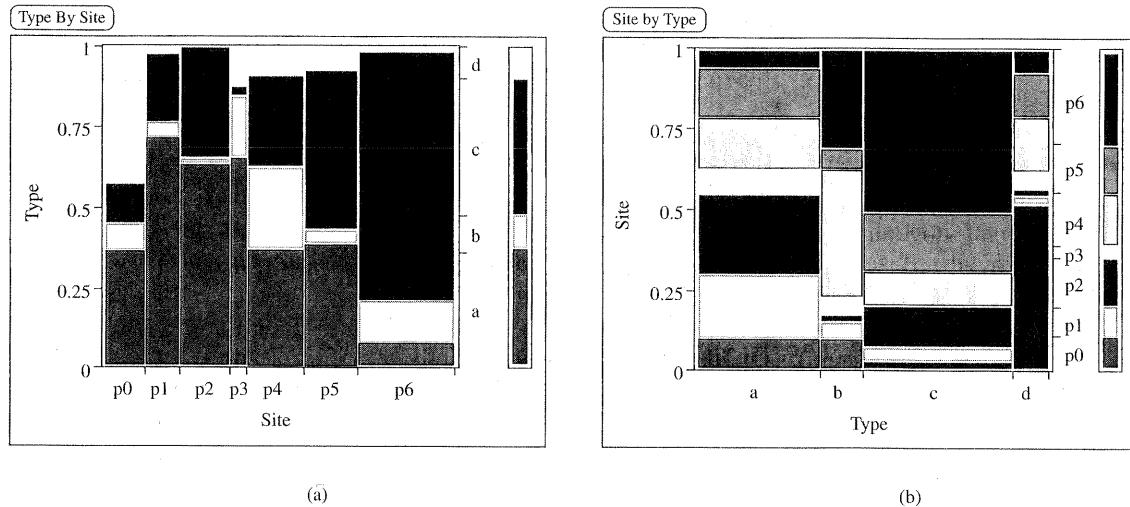


Figure 12.21 Site and pottery type profiles for the data in Table 12.8.

column totals. The bars in the figure appear to be quite different from one another. This suggests that the various types of pottery are not distributed over the archaeological sites in the same way.

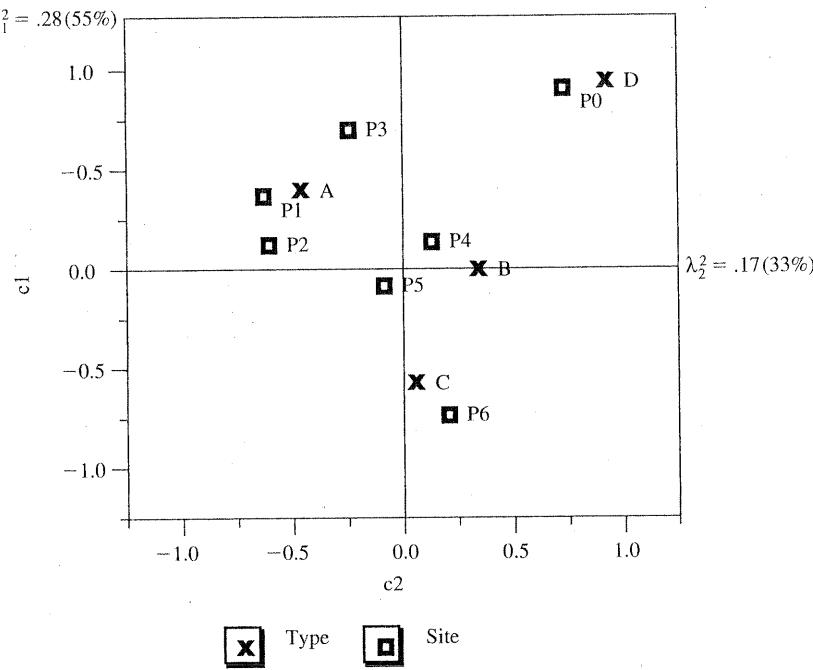
The two-dimensional plot from a correspondence analysis<sup>2</sup> of the pottery type–site data is shown in Figure 12.22.

The plot in Figure 12.22 indicates, for example, that sites P1 and P2 have similar pottery type profiles (the two points are close together), and sites P0 and P6 have very different profiles (the points are far apart). The individual points representing the types of pottery are spread out, indicating that their archaeological site profiles are quite different. These findings are consistent with the profiles pictured in Figure 12.21.

Notice that the points P0 and D are quite close together and separated from the remaining points. This indicates that pottery type D tends to be associated, almost exclusively, with site P0. Similarly, pottery type A tends to be associated with site P1 and, to lesser degrees, with sites P2 and P3. Pottery type B is associated with sites P4 and P5, and pottery type C tends to be associated, again, almost exclusively, with site P6. Since the archaeological sites represent different periods, these associations are of considerable interest to archaeologists.

The number  $\lambda_1^2 = .28$  at the end of the first coordinate axis in the two-dimensional plot is the *inertia* associated with the first dimension. This inertia is 55% of the total inertia. The inertia associated with the second dimension is  $\lambda_2^2 = .17$ , and the second dimension accounts for 33% of the total inertia. Together, the two dimensions account for  $55\% + 33\% = 88\%$  of the total inertia. Since, in this case, the data could be exactly represented in three dimensions, relatively little information (variation) is lost by representing the data in the two-dimensional plot of Figure 12.22. Equivalently, we may regard this plot as the best two-dimensional representation of the multidimensional scatter of row points and the multidimensional

<sup>2</sup>The JMP software was used for a correspondence analysis of the data in Table 12.8.



**Figure 12.22** A correspondence analysis plot of the pottery type–site data.

scatter of column points. The combined inertia of 88% suggests that the representation “fits” the data well.

In this example, the graphical output from a correspondence analysis shows the nature of the associations in the contingency table quite clearly.

### Algebraic Development of Correspondence Analysis

To begin, let  $\mathbf{X}$ , with elements  $x_{ij}$ , be an  $I \times J$  two-way table of unscaled frequencies or counts. In our discussion we take  $I > J$  and assume that  $\mathbf{X}$  is of full column rank  $J$ . The rows and columns of the contingency table  $\mathbf{X}$  correspond to different categories of two different characteristics. As an example, the array of frequencies of different pottery types at different archaeological sites shown in Table 12.8 is a contingency table with  $I = 7$  archaeological sites and  $J = 4$  pottery types.

If  $n$  is the total of the frequencies in the data matrix  $\mathbf{X}$ , we first construct a matrix of proportions  $\mathbf{P} = \{p_{ij}\}$  by dividing each element of  $\mathbf{X}$  by  $n$ . Hence

$$p_{ij} = \frac{x_{ij}}{n}, \quad i = 1, 2, \dots, I, \quad j = 1, 2, \dots, J, \quad \text{or} \quad \mathbf{P}_{(I \times J)} = \frac{1}{n} \mathbf{X}_{(I \times J)} \quad (12-28)$$

The matrix  $\mathbf{P}$  is called the *correspondence matrix*.

Next define the vectors of row and column sums  $\mathbf{r}$  and  $\mathbf{c}$  respectively, and the diagonal matrices  $\mathbf{D}_r$  and  $\mathbf{D}_c$  with the elements of  $\mathbf{r}$  and  $\mathbf{c}$  on the diagonals. Thus

$$\begin{aligned} r_i &= \sum_{j=1}^J p_{ij} = \sum_{j=1}^J \frac{x_{ij}}{n}, \quad i = 1, 2, \dots, I, \quad \text{or} \quad \underset{(I \times 1)}{\mathbf{r}} = \underset{(I \times J)(J \times 1)}{\mathbf{P} \mathbf{1}_J} \\ c_j &= \sum_{i=1}^I p_{ij} = \sum_{i=1}^I \frac{x_{ij}}{n}, \quad j = 1, 2, \dots, J, \quad \text{or} \quad \underset{(J \times 1)}{\mathbf{c}} = \underset{(J \times I)(I \times 1)}{\mathbf{P}' \mathbf{1}_I} \end{aligned} \quad (12-29)$$

where  $\mathbf{1}_J$  is a  $J \times 1$  and  $\mathbf{1}_I$  is a  $I \times 1$  vector of 1's and

$$\mathbf{D}_r = \text{diag}(r_1, r_2, \dots, r_I) \quad \text{and} \quad \mathbf{D}_c = \text{diag}(c_1, c_2, \dots, c_J) \quad (12-30)$$

We define the square root matrices

$$\begin{aligned} \mathbf{D}_r^{1/2} &= \text{diag}(\sqrt{r_1}, \dots, \sqrt{r_I}) \quad \mathbf{D}_r^{-1/2} = \text{diag}\left(\frac{1}{\sqrt{r_1}}, \dots, \frac{1}{\sqrt{r_I}}\right) \\ \mathbf{D}_c^{1/2} &= \text{diag}(\sqrt{c_1}, \dots, \sqrt{c_J}) \quad \mathbf{D}_c^{-1/2} = \text{diag}\left(\frac{1}{\sqrt{c_1}}, \dots, \frac{1}{\sqrt{c_J}}\right) \end{aligned} \quad (12-31)$$

for scaling purposes.

Correspondence analysis can be formulated as the weighted least squares problem to select  $\hat{\mathbf{P}} = \{\hat{p}_{ij}\}$ , a matrix of specified reduced rank, to minimize

$$\sum_{i=1}^I \sum_{j=1}^J \frac{(p_{ij} - \hat{p}_{ij})^2}{r_i c_j} = \text{tr}[(\mathbf{D}_r^{-1/2}(\mathbf{P} - \hat{\mathbf{P}})\mathbf{D}_c^{-1/2})(\mathbf{D}_r^{-1/2}(\mathbf{P} - \hat{\mathbf{P}})\mathbf{D}_c^{-1/2})'] \quad (12-32)$$

since  $(p_{ij} - \hat{p}_{ij})/\sqrt{r_i c_j}$  is the  $(i, j)$  element of  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \hat{\mathbf{P}})\mathbf{D}_c^{-1/2}$ .

As Result 12.1 demonstrates, the term  $\mathbf{rc}'$  is common to the approximation  $\hat{\mathbf{P}}$  whatever the  $I \times J$  correspondence matrix  $\mathbf{P}$ . The matrix  $\hat{\mathbf{P}} = \mathbf{rc}'$  can be shown to be the best rank 1 approximation to  $\mathbf{P}$ .

**Result 12.1.** The term  $\mathbf{rc}'$  is common to the approximation  $\hat{\mathbf{P}}$  whatever the  $I \times J$  correspondence matrix  $\mathbf{P}$ .

The reduced rank  $s$  approximation to  $\mathbf{P}$ , which minimizes the sum of squares (12-32), is given by

$$\mathbf{P} \doteq \sum_{k=1}^s \tilde{\lambda}_k (\mathbf{D}_r^{1/2} \tilde{\mathbf{u}}_k) (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)' = \mathbf{rc}' + \sum_{k=2}^s \tilde{\lambda}_k (\mathbf{D}_r^{1/2} \tilde{\mathbf{u}}_k) (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)'$$

where the  $\tilde{\lambda}_k$  are the singular values and the  $I \times 1$  vectors  $\tilde{\mathbf{u}}_k$  and the  $J \times 1$  vectors  $\tilde{\mathbf{v}}_k$  are the corresponding singular vectors of the  $I \times J$  matrix  $\mathbf{D}_r^{-1/2} \mathbf{P} \mathbf{D}_c^{-1/2}$ . The minimum value of (12-32) is  $\sum_{k=s+1}^J \tilde{\lambda}_k^2$ .

The reduced rank  $K > 1$  approximation to  $\mathbf{P} - \mathbf{rc}'$  is

$$\mathbf{P} - \mathbf{rc}' \doteq \sum_{k=1}^K \lambda_k (\mathbf{D}_r^{1/2} \mathbf{u}_k) (\mathbf{D}_c^{1/2} \mathbf{v}_k)' \quad (12-33)$$

$\lambda_2^2 = .17(33\%)$

pe-site data.

at the representa-

analysis shows the

■

of unscaled fre-  
ie that  $\mathbf{X}$  is of full  
>  $\mathbf{X}$  correspond to  
mple, the array of  
cal sites shown in  
es and  $J = 4$  pot-

rst construct a ma-  
. Hence

$$= \frac{1}{n} \underset{(I \times J)}{\mathbf{X}} \quad (12-28)$$

where the  $\lambda_k$  are the singular values and the  $I \times 1$  vectors  $\mathbf{u}_k$  and the  $J \times 1$  vectors  $\mathbf{v}_k$  are the corresponding singular vectors of the  $I \times J$  matrix  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$ .

Here  $\lambda_k = \tilde{\lambda}_{k+1}$ ,  $\mathbf{u}_k = \tilde{\mathbf{u}}_{k+1}$ , and  $\mathbf{v}_k = \tilde{\mathbf{v}}_{k+1}$  for  $k = 1, \dots, J-1$ .

**Proof.** We first consider a scaled version  $\mathbf{B} = \mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2}$  of the correspondence matrix  $\mathbf{P}$ . According to Result 2A.16, the best low rank =  $s$  approximation  $\hat{\mathbf{B}}$  to  $\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2}$  is given by the first  $s$  terms in the the singular-value decomposition

$$\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2} = \sum_{k=1}^J \tilde{\lambda}_k \tilde{\mathbf{u}}_k \tilde{\mathbf{v}}'_k \quad (12-34)$$

where

$$\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2}\tilde{\mathbf{v}}_k = \tilde{\lambda}_k \tilde{\mathbf{u}}_k \quad \tilde{\mathbf{u}}'_k \mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2} = \tilde{\lambda}_k \tilde{\mathbf{v}}'_k \quad (12-35)$$

and

$$|(\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2})(\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2})' - \tilde{\lambda}_k^2 \mathbf{I}| = 0 \quad \text{for } k = 1, \dots, J$$

The approximation to  $\mathbf{P}$  is then given by

$$\hat{\mathbf{P}} = \mathbf{D}_r^{1/2} \hat{\mathbf{B}} \mathbf{D}_c^{1/2} = \sum_{k=1}^s \tilde{\lambda}_k (\mathbf{D}_r^{1/2} \tilde{\mathbf{u}}_k) (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)'$$

and, by Result 2A.16, the error of approximation is  $\sum_{k=s+1}^J \tilde{\lambda}_k^2$ .

Whatever the correspondence matrix  $\mathbf{P}$ , the term  $\mathbf{rc}'$  always provides a (the best) rank one approximation. This corresponds to the assumption of independence of the rows and columns. To see this, let  $\tilde{\mathbf{u}}_1 = \mathbf{D}_r^{1/2}\mathbf{1}_I$  and  $\tilde{\mathbf{v}}_1 = \mathbf{D}_c^{1/2}\mathbf{1}_J$ , where  $\mathbf{1}_I$  is a  $I \times 1$  and  $\mathbf{1}_J$  a  $J \times 1$  vector of 1's. We verify that (12-35) holds for these choices.

$$\begin{aligned} \tilde{\mathbf{u}}'_1 (\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2}) &= (\mathbf{D}_r^{1/2}\mathbf{1}_I)' (\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2}) \\ &= \mathbf{1}'_I \mathbf{P} \mathbf{D}_c^{-1/2} = \mathbf{c}' \mathbf{D}_c^{-1/2} \\ &= [\sqrt{c_1}, \dots, \sqrt{c_J}] = (\mathbf{D}_c^{1/2}\mathbf{1}_J)' = \tilde{\mathbf{v}}'_1 \end{aligned}$$

and

$$\begin{aligned} (\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2})\tilde{\mathbf{v}}_1 &= (\mathbf{D}_r^{-1/2}\mathbf{PD}_c^{-1/2})(\mathbf{D}_c^{1/2}\mathbf{1}_J) \\ &= \mathbf{D}_r^{-1/2}\mathbf{P}\mathbf{1}_J = \mathbf{D}_r^{-1/2}\mathbf{r} \\ &= \begin{bmatrix} \sqrt{r_1} \\ \vdots \\ \sqrt{r_I} \end{bmatrix} = \mathbf{D}_r^{1/2}\mathbf{1}_I = \tilde{\mathbf{u}}_1 \end{aligned}$$

That is,

$$(\tilde{\mathbf{u}}_1, \tilde{\mathbf{v}}_1) = (\mathbf{D}_r^{1/2}\mathbf{1}_I, \mathbf{D}_c^{1/2}\mathbf{1}_J) \quad (12-36)$$

are singular vectors associated with singular value  $\tilde{\lambda}_1 = 1$ . For any correspondence matrix,  $\mathbf{P}$ , the common term in every expansion is

$$\mathbf{D}_r^{1/2}\mathbf{u}_1 \mathbf{v}'_1 \mathbf{D}_c^{1/2} = \mathbf{D}_r \mathbf{1}_I \mathbf{1}'_J \mathbf{D}_c = \mathbf{rc}'$$

the  $J \times 1$  vectors  
 $\frac{1}{2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$ .  
 $- 1$ .

the correspondence  
 approximation  $\hat{\mathbf{B}}$  to  
 decomposition

(12-34)

$\tilde{v}'_k$  (12-35)

$1, \dots, J$

provides a (the  
 n of independence  
 $\frac{1}{2}\mathbf{1}_J$ , where  $\mathbf{1}_J$  is a  
 r these choices.

(12-36)

by correspondence

Therefore, we have established the first approximation and (12-34) can always be expressed as

$$\mathbf{P} = \mathbf{rc}' + \sum_{k=2}^J \tilde{\lambda}_k (\mathbf{D}_r^{1/2} \tilde{\mathbf{u}}_k) (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)'$$

Because of the common term, the problem can be rephrased in terms of  $\mathbf{P} - \mathbf{rc}'$  and its scaled version  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$ . By the orthogonality of the singular vectors of  $\mathbf{D}_r^{-1/2}\mathbf{P}\mathbf{D}_c^{-1/2}$ , we have  $\tilde{\mathbf{u}}'_k(\mathbf{D}_r^{1/2}\mathbf{1}_I) = 0$  and  $\tilde{\mathbf{v}}'_k(\mathbf{D}_c^{1/2}\mathbf{1}_J) = 0$ , for  $k > 1$ , so

$$\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2} = \sum_{k=2}^J \tilde{\lambda}_k \tilde{\mathbf{u}}_k \tilde{\mathbf{v}}'_k$$

is the singular-value decomposition of  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$  in terms of the singular values and vectors obtained from  $\mathbf{D}_r^{-1/2}\mathbf{P}\mathbf{D}_c^{-1/2}$ . Converting to singular values and vectors  $\lambda_k$ ,  $\mathbf{u}_k$ , and  $\mathbf{v}_k$  from  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$  only amounts to changing  $k$  to  $k - 1$  so  $\lambda_k = \lambda_{k+1}$ ,  $\mathbf{u}_k = \tilde{\mathbf{u}}_{k+1}$ , and  $\mathbf{v}_k = \tilde{\mathbf{v}}_{k+1}$  for  $k = 1, \dots, J - 1$ .

In terms of the singular value decomposition for  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$ , the expansion for  $\mathbf{P} - \mathbf{rc}'$  takes the form

$$\mathbf{P} - \mathbf{rc}' = \sum_{k=1}^{J-1} \lambda_k (\mathbf{D}_r^{1/2} \mathbf{u}_k) (\mathbf{D}_c^{1/2} \mathbf{v}_k)' \quad (12-37)$$

The best rank  $K$  approximation to  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$  is given by  $\sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}'_k$ . Then, the best approximation to  $\mathbf{P} - \mathbf{rc}'$  is

$$\mathbf{P} - \mathbf{rc}' \doteq \sum_{k=1}^K \lambda_k (\mathbf{D}_r^{1/2} \mathbf{u}_k) (\mathbf{D}_c^{1/2} \mathbf{v}_k)' \quad (12-38)$$

**Remark.** Note that the vectors  $\mathbf{D}_r^{1/2} \mathbf{u}_k$  and  $\mathbf{D}_c^{1/2} \mathbf{v}_k$  in the expansion (12-38) of  $\mathbf{P} - \mathbf{rc}'$  need not have length 1 but satisfy the scaling

$$\begin{aligned} (\mathbf{D}_r^{1/2} \mathbf{u}_k)' \mathbf{D}_r^{-1} (\mathbf{D}_r^{1/2} \mathbf{u}_k) &= \mathbf{u}'_k \mathbf{u}_k = 1 \\ (\mathbf{D}_c^{1/2} \mathbf{v}_k)' \mathbf{D}_c^{-1} (\mathbf{D}_c^{1/2} \mathbf{v}_k) &= \mathbf{v}'_k \mathbf{v}_k = 1 \end{aligned}$$

Because of this scaling, the expansions in Result 12.1 have been called a generalized singular-value decomposition.

Let  $\Lambda$ ,  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_J]$  and  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_J]$  be the matrices of singular values and vectors obtained from  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2}$ . It is usual in correspondence analysis to plot the first two or three columns of  $\mathbf{F} = \mathbf{D}_r^{-1}(\mathbf{D}_r^{1/2}\mathbf{U})\Lambda$  and  $\mathbf{G} = \mathbf{D}_c^{-1}(\mathbf{D}_c^{1/2}\mathbf{V})\Lambda$  or  $\lambda_k \mathbf{D}_r^{-1/2} \mathbf{u}_k$  and  $\lambda_k \mathbf{D}_c^{-1/2} \mathbf{v}_k$  for  $k = 1, 2$ , and maybe 3.

The joint plot of the coordinates in  $\mathbf{F}$  and  $\mathbf{G}$  is called a *symmetric map* (see Greenacre [13]) since the points representing the rows and columns have the same normalization, or scaling, along the dimensions of the solution. That is, the geometry for the row points is identical to the geometry for the column points.

**Example 12.18 (Calculations for correspondence analysis)** Consider the  $3 \times 2$  contingency table

	B1	B2	Total
A1	24	12	36
A2	16	48	64
A3	60	40	100
	100	100	200

The correspondence matrix is

$$\mathbf{P} = \begin{bmatrix} .12 & .06 \\ .08 & .24 \\ .30 & .20 \end{bmatrix}$$

with marginal totals  $\mathbf{c}' = [.5, .5]$  and  $\mathbf{r}' = [.18, .32, .50]$ . The negative square root matrices are

$$\mathbf{D}_r^{-1/2} = \text{diag}(\sqrt{2}/.6, \sqrt{2}/.8, \sqrt{2}) \quad \mathbf{D}_c^{-1/2} = \text{diag}(\sqrt{2}, \sqrt{2})$$

Then

$$\mathbf{P} - \mathbf{rc}' = \begin{bmatrix} .12 & .06 \\ .08 & .24 \\ .30 & .20 \end{bmatrix} - \begin{bmatrix} .18 \\ .32 \\ .50 \end{bmatrix} \begin{bmatrix} .5 & .5 \end{bmatrix} = \begin{bmatrix} .03 & -.03 \\ -.08 & .08 \\ .05 & -.05 \end{bmatrix}$$

The scaled version of this matrix is

$$\begin{aligned} \mathbf{A} &= \mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2} = \begin{bmatrix} \frac{\sqrt{2}}{.6} & 0 & 0 \\ 0 & \frac{\sqrt{2}}{.8} & 0 \\ 0 & 0 & \sqrt{2} \end{bmatrix} \begin{bmatrix} .03 & -.03 \\ -.08 & .08 \\ .05 & -.05 \end{bmatrix} \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{2} \end{bmatrix} \\ &= \begin{bmatrix} 0.1 & -0.1 \\ -0.2 & 0.2 \\ 0.1 & -0.1 \end{bmatrix} \end{aligned}$$

Since  $I > J$ , the square of the singular values and the  $\mathbf{v}_i$  are determined from

$$\mathbf{A}'\mathbf{A} = \begin{bmatrix} .1 & -.2 & .1 \\ -.1 & .2 & -.1 \end{bmatrix} \begin{bmatrix} .1 & -.1 \\ -.2 & .2 \\ .1 & -.1 \end{bmatrix} = \begin{bmatrix} .06 & -.06 \\ -.06 & .06 \end{bmatrix}$$

sider the  $3 \times 2$

It is easily checked that  $\lambda_1^2 = .12$ ,  $\lambda_2^2 = 0$ , since  $J - 1 = 1$ , and that

$$\mathbf{v}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$$

Further,

$$\mathbf{A}\mathbf{A}' = \begin{bmatrix} .1 & -.1 \\ -.2 & .2 \\ .1 & -.1 \end{bmatrix} \begin{bmatrix} .1 & -.2 & .1 \\ -.1 & .2 & -.1 \end{bmatrix} = \begin{bmatrix} .02 & -.04 & .02 \\ -.04 & .08 & -.04 \\ .02 & -.04 & .02 \end{bmatrix}$$

A computer calculation confirms that the single nonzero eigenvalue is  $\lambda_1^2 = .12$ , so that the singular value has absolute value  $\lambda_1 = .2\sqrt{3}$  and, as you can easily check,

$$\mathbf{u}_1 = \begin{bmatrix} \frac{1}{\sqrt{6}} \\ -\frac{2}{\sqrt{6}} \\ \frac{1}{\sqrt{6}} \end{bmatrix}$$

The expansion of  $\mathbf{P} - \mathbf{rc}'$  is then the single term

$$\lambda_1(\mathbf{D}_r^{1/2}\mathbf{u}_1)(\mathbf{D}_c^{1/2}\mathbf{v}_1)'$$

$$= \sqrt{.12} \begin{bmatrix} \frac{.6}{\sqrt{2}} & 0 & 0 \\ 0 & \frac{.8}{\sqrt{2}} & 0 \\ 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} \\ -\frac{2}{\sqrt{6}} \\ \frac{1}{\sqrt{6}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$= \sqrt{.12} \begin{bmatrix} \frac{.3}{\sqrt{3}} \\ -\frac{.8}{\sqrt{3}} \\ \frac{.5}{\sqrt{3}} \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{-1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} = \begin{bmatrix} .03 & -.03 \\ -.08 & .08 \\ .05 & -.05 \end{bmatrix} \quad \text{check}$$

There is only one pair of vectors to plot

$$\lambda_1 \mathbf{D}_r^{1/2} \mathbf{u}_1 = \sqrt{.12} \begin{bmatrix} \frac{.6}{\sqrt{2}} & 0 & 0 \\ 0 & \frac{.8}{\sqrt{2}} & 0 \\ 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} \\ -\frac{2}{\sqrt{6}} \\ \frac{1}{\sqrt{6}} \end{bmatrix} = \sqrt{.12} \begin{bmatrix} \frac{.3}{\sqrt{3}} \\ -\frac{.8}{\sqrt{3}} \\ \frac{.5}{\sqrt{3}} \end{bmatrix}$$

and

$$\lambda_1 \mathbf{D}_c^{1/2} \mathbf{v}_1 = \sqrt{.12} \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} = \sqrt{.12} \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}$$

There is a second way to define contingency analysis. Following Greenacre [13], we call the preceding approach the *matrix approximation method* and the approach to follow the *profile approximation method*. We illustrate the profile approximation method using the row profiles; however, an analogous solution results if we were to begin with the column profiles.

Algebraically, the row profiles are the rows of the matrix  $\mathbf{D}_r^{-1} \mathbf{P}$ , and contingency analysis can be defined as the approximation of the row profiles by points in a low-dimensional space. Consider approximating the row profiles by the matrix  $\mathbf{P}^*$ . Using the square-root matrices  $\mathbf{D}_r^{1/2}$  and  $\mathbf{D}_c^{1/2}$  defined in (12-31), we can write

$$(\mathbf{D}_r^{-1} \mathbf{P} - \mathbf{P}^*) \mathbf{D}_c^{-1/2} = \mathbf{D}_r^{-1/2} (\mathbf{D}_r^{-1/2} \mathbf{P} - \mathbf{D}_r^{1/2} \mathbf{P}^*) \mathbf{D}_c^{-1/2}$$

and the least squares criterion (12-32) can be written, with  $p_{ij}^* = \hat{p}_{ij}/r_i$ , as

$$\begin{aligned} \sum_i \sum_j \frac{(p_{ij} - \hat{p}_{ij})^2}{r_i c_j} &= \sum_i r_i \sum_j \frac{(p_{ij}/r_i - p_{ij}^*)^2}{c_j} \\ &= \text{tr} [\mathbf{D}_r^{1/2} \mathbf{D}_r^{1/2} (\mathbf{D}_r^{-1} \mathbf{P} - \mathbf{P}^*) \mathbf{D}_c^{-1/2} \mathbf{D}_c^{-1/2} (\mathbf{D}_r^{-1} \mathbf{P} - \mathbf{P}^*)'] \\ &= \text{tr} [\mathbf{D}_r^{1/2} (\mathbf{D}_r^{-1/2} \mathbf{P} - \mathbf{D}_r^{1/2} \mathbf{P}^*) \mathbf{D}_c^{-1/2} \mathbf{D}_c^{-1/2} (\mathbf{D}_r^{-1/2} \mathbf{P} - \mathbf{D}_r^{1/2} \mathbf{P}^*)' \mathbf{D}_r^{-1/2}] \\ &= \text{tr} [[(\mathbf{D}_r^{-1/2} \mathbf{P} - \mathbf{D}_r^{1/2} \mathbf{P}^*) \mathbf{D}_c^{-1/2}] [(\mathbf{D}_r^{-1/2} \mathbf{P} - \mathbf{D}_r^{1/2} \mathbf{P}^*) \mathbf{D}_c^{-1/2}]] \quad (12-39) \end{aligned}$$

Minimizing the last expression for the trace in (12-39) is precisely the first minimization problem treated in the proof of Result 12.1. By (12-34),  $\mathbf{D}_r^{-1/2} \mathbf{P} \mathbf{D}_c^{-1/2}$  has the singular-value decomposition

$$\mathbf{D}_r^{-1/2} \mathbf{P} \mathbf{D}_c^{-1/2} = \sum_{k=1}^J \tilde{\lambda}_k \tilde{\mathbf{u}}_k \tilde{\mathbf{v}}_k' \quad (12-40)$$

The best rank  $K$  approximation is obtained by using the first  $K$  terms of this expansion. Since, by (12-39), we have  $\mathbf{D}_r^{-1/2} \mathbf{P} \mathbf{D}_c^{-1/2}$  approximated by  $\mathbf{D}_r^{1/2} \mathbf{P}^* \mathbf{D}_c^{-1/2}$ , we left

multiply by  $\mathbf{D}_r^{-1/2}$  and right multiply by  $\mathbf{D}_c^{1/2}$  to obtain the generalized singular-value decomposition

$$\mathbf{D}_r^{-1}\mathbf{P} = \sum_{k=1}^J \tilde{\lambda}_k \mathbf{D}_r^{-1/2} \tilde{\mathbf{u}}_k (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)' \quad (12-41)$$

where, from (12-36),  $(\tilde{\mathbf{u}}_1, \tilde{\mathbf{v}}_1) = (\mathbf{D}_r^{1/2} \mathbf{1}_I, \mathbf{D}_c^{1/2} \mathbf{1}_J)$  are singular vectors associated with singular value  $\lambda_1 = 1$ . Since  $\mathbf{D}_r^{-1/2}(\mathbf{D}_r^{1/2} \mathbf{1}_I) = \mathbf{1}_I$  and  $(\mathbf{D}_c^{1/2} \mathbf{1}_J)' \mathbf{D}_c^{1/2} = \mathbf{c}'$ , the leading term in the decomposition (12-41) is  $\mathbf{1}_I \mathbf{c}'$ .

Consequently, in terms of the singular values and vectors from  $\mathbf{D}_r^{-1/2} \mathbf{P} \mathbf{D}_c^{-1/2}$ , the reduced rank  $K < J$  approximation to the row profiles  $\mathbf{D}_r^{-1}\mathbf{P}$  is

$$\mathbf{P}^* = \mathbf{1}_I \mathbf{c}' + \sum_{k=2}^K \tilde{\lambda}_k \mathbf{D}_r^{-1/2} \tilde{\mathbf{u}}_k (\mathbf{D}_c^{1/2} \tilde{\mathbf{v}}_k)' \quad (12-42)$$

In terms of the singular values and vectors  $\lambda_k, \mathbf{u}_k$  and  $\mathbf{v}_k$  obtained from  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}') \mathbf{D}_c^{-1/2}$ , we can write

$$\mathbf{P}^* - \mathbf{1}_I \mathbf{c}' = \sum_{k=1}^{K-1} \lambda_k \mathbf{D}_r^{-1/2} \mathbf{u}_k (\mathbf{D}_c^{1/2} \mathbf{v}_k)'$$

(Row profiles for the archaeological data in Table 12.8 are shown in Figure 12.21 on page 717.)

## Inertia

Total inertia is a measure of the variation in the count data and is defined as the weighted sum of squares

$$\text{tr}[\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}') \mathbf{D}_c^{-1/2}(\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}') \mathbf{D}_c^{-1/2})'] = \sum_i \sum_j \frac{(p_{ij} - r_i c_j)^2}{r_i c_j} = \sum_{k=1}^{J-1} \lambda_k^2 \quad (12-43)$$

where the  $\lambda_k$  are the singular values obtained from the singular-value decomposition of  $\mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}') \mathbf{D}_c^{-1/2}$  (see the proof of Result 12.1).<sup>3</sup>

The inertia associated with the best reduced rank  $K < J$  approximation to the centered matrix  $\mathbf{P} - \mathbf{rc}'$  (the  $K$ -dimensional solution) has inertia  $\sum_{k=1}^K \lambda_k^2$ . The residual inertia (variation) not accounted for by the rank  $K$  solution is equal to the sum of squares of the remaining singular values:  $\lambda_{K+1}^2 + \lambda_{K+2}^2 + \dots + \lambda_{J-1}^2$ . For plots, the inertia associated with dimension  $k, \lambda_k^2$ , is ordinarily displayed along the  $k$ th coordinate axis, as in Figure 12.22 for  $k = 1, 2$ .

<sup>3</sup>Total inertia is related to the chi-square measure of association in a two-way contingency table,  $\chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$ . Here  $O_{ij} = x_{ij}$  is the observed frequency and  $E_{ij}$  is the expected frequency for the  $ij$ th cell. In our context, if the row variable is independent of (unrelated to) the column variable,  $E_{ij} = n r_i c_j$ , and

$$\text{Total inertia} = \sum_{i=1}^I \sum_{j=1}^J \frac{(p_{ij} - r_i c_j)^2}{r_i c_j} = \frac{\chi^2}{n}$$

## Interpretation in Two Dimensions

Since the inertia is a measure of the data table's total variation, how do we interpret a large value for the proportion  $(\lambda_1^2 + \lambda_2^2)/\sum_{k=1}^{J-1} \lambda_k^2$ ? Geometrically, we say that the associations in the centered data are well represented by points in a plane, and this best approximating plane accounts for nearly all the variation in the data beyond that accounted for by the rank 1 solution (independence model). Algebraically, we say that the approximation

$$\mathbf{P} = \mathbf{rc}' \doteq \lambda_1 \mathbf{u}_1 \mathbf{v}_1' + \lambda_2 \mathbf{u}_2 \mathbf{v}_2'$$

is very good or, equivalently, that

$$\mathbf{P} \doteq \mathbf{rc}' + \lambda_1 \mathbf{u}_1 \mathbf{v}_1' + \lambda_2 \mathbf{u}_2 \mathbf{v}_2'$$

## Final Comments

Correspondence analysis is primarily a graphical technique designed to represent associations in a low-dimensional space. It can be regarded as a scaling method, and can be viewed as a complement to other methods such as multidimensional scaling (Section 12.6) and biplots (Section 12.8). Correspondence analysis also has links to principal component analysis (Chapter 8) and canonical correlation analysis (Chapter 10). The book by Greenacre [14] is one choice for learning more about correspondence analysis.

## 12.8 Biplots for Viewing Sampling Units and Variables

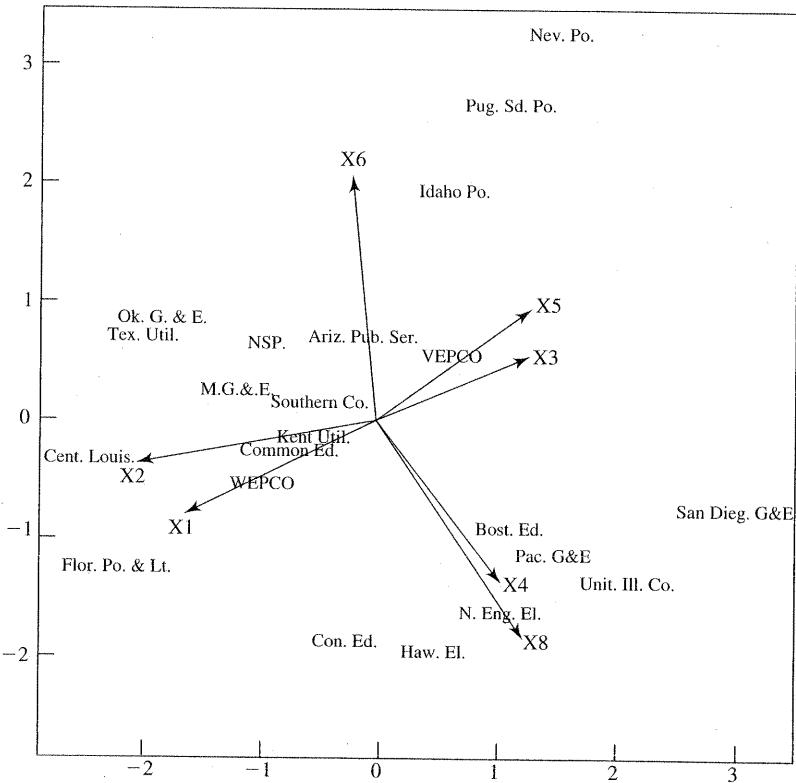
A *biplot* is a graphical representation of the information in an  $n \times p$  data matrix. The *bi-* refers to the two kinds of information contained in a data matrix. The information in the rows pertains to samples or sampling units and that in the columns pertains to variables.

When there are only two variables, scatter plots can represent the information on both the sampling units and the variables in a single diagram. This permits the visual inspection of the position of one sampling unit relative to another and the relative importance of each of the two variables to the position of any unit.

With several variables, one can construct a matrix array of scatter plots, but there is no one single plot of the sampling units. On the other hand, a two-dimensional plot of the sampling units can be obtained by graphing the first two principal components, as in Section 8.4. The idea behind biplots is to add the information about the variables to the principal component graph.

Figure 12.23 gives an example of a biplot for the public utilities data in Table 12.4.

You can see how the companies group together and which variables contribute to their positioning within this representation. For instance,  $X_4$  = annual load factor and  $X_8$  = total fuel costs are primarily responsible for the grouping of the mostly coastal companies in the lower right. The two variables  $X_1$  = fixed-



**Figure 12.23** A biplot of the data on public utilities.

w do we interpret  
, we say that the  
a plane, and this  
the data beyond  
Algebraically, we

igned to represent  
iling method, and  
nensional scaling  
s also has links to  
relation analysis  
ning more about

bles

$\times p$  data matrix.  
matrix. The infor-  
at in the columns

t the information  
his permits the vi-  
ther and the rela-  
unit.

of scatter plots,  
her hand, a two-  
ching the first two  
to add the infor-

utilities data in  
ch variables con-  
nce,  $X_4$  = annual  
r the grouping of  
ables  $X_1$  = fixed-

charge ratio and  $X_2$  = rate of return on capital put the Florida and Louisiana companies together.

### Constructing Biplots

The construction of a biplot proceeds from the sample principal components.

According to Result 8A.1, the best two-dimensional approximation to the data matrix  $\mathbf{X}$  approximates the  $j$ th observation  $\mathbf{x}_j$  in terms of the sample values of the first two principal components. In particular,

$$\mathbf{x}_j \doteq \bar{\mathbf{x}} + \hat{y}_{j1}\hat{\mathbf{e}}_1 + \hat{y}_{j2}\hat{\mathbf{e}}_2 \quad (12-44)$$

where  $\hat{\mathbf{e}}_1$  and  $\hat{\mathbf{e}}_2$  are the first two eigenvectors of  $\mathbf{S}$  or, equivalently, of  $\mathbf{X}'_c\mathbf{X}_c = (n-1)\mathbf{S}$ . Here  $\mathbf{X}_c$  denotes the mean corrected data matrix with rows  $(\mathbf{x}_j - \bar{\mathbf{x}})'$ . The eigenvectors determine a plane, and the coordinates of the  $j$ th unit (row) are the pair of values of the principal components,  $(\hat{y}_{j1}, \hat{y}_{j2})$ .

To include the information on the variables in this plot, we consider the pair of eigenvectors  $(\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2)$ . These eigenvectors are the coefficient vectors for the first two sample principal components. Consequently, each row of the matrix  $\hat{\mathbf{E}} = [\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2]$

positions a variable in the graph, and the magnitudes of the coefficients (the coordinates of the variable) show the weightings that variable has in each principal component. The positions of the variables in the plot are indicated by a vector. Usually, statistical computer programs include a multiplier so that the lengths of all of the vectors can be suitably adjusted and plotted on the same axes as the sampling units. Units that are close to a variable likely have high values on that variable. To interpret a new point  $\mathbf{x}_0$ , we plot its principal components  $\hat{\mathbf{E}}'(\mathbf{x}_0 - \bar{\mathbf{x}})$ .

A direct approach to obtaining a biplot starts from the singular value decomposition (see Result 2A.15), which first expresses the  $n \times p$  mean corrected matrix  $\mathbf{X}_c$  as

$$\mathbf{X}_c = \underset{(n \times p)}{\mathbf{U}} \underset{(p \times p)}{\Lambda} \underset{(p \times p)}{\mathbf{V}'} \quad (12-45)$$

where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  and  $\mathbf{V}$  is an orthogonal matrix whose columns are the eigenvectors of  $\mathbf{X}_c' \mathbf{X}_c = (n-1)\mathbf{S}$ . That is,  $\mathbf{V} = \hat{\mathbf{E}} = [\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2, \dots, \hat{\mathbf{e}}_p]$ . Multiplying (12-45) on the right by  $\hat{\mathbf{E}}$ , we find

$$\mathbf{X}_c \hat{\mathbf{E}} = \mathbf{U} \Lambda \quad (12-46)$$

where the  $j$ th row of the left-hand side,

$$[(\mathbf{x}_j - \bar{\mathbf{x}})' \hat{\mathbf{e}}_1, (\mathbf{x}_j - \bar{\mathbf{x}})' \hat{\mathbf{e}}_2, \dots, (\mathbf{x}_j - \bar{\mathbf{x}})' \hat{\mathbf{e}}_p] = [\hat{y}_{j1}, \hat{y}_{j2}, \dots, \hat{y}_{jp}]$$

is just the value of the principal components for the  $j$ th item. That is,  $\mathbf{U} \Lambda$  contains all of the values of the principal components, while  $\mathbf{V} = \hat{\mathbf{E}}$  contains the coefficients that define the principal components.

The best rank 2 approximation to  $\mathbf{X}_c$  is obtained by replacing  $\Lambda$  by  $\Lambda^* = \text{diag}(\lambda_1, \lambda_2, 0, \dots, 0)$ . This result, called the Eckart–Young theorem, was established in Result 8.A.1. The approximation is then

$$\mathbf{X}_c \doteq \mathbf{U} \Lambda^* \mathbf{V}' = [\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2] \begin{bmatrix} \hat{\mathbf{e}}'_1 \\ \hat{\mathbf{e}}'_2 \end{bmatrix} \quad (12-47)$$

where  $\hat{\mathbf{y}}_1$  is the  $n \times 1$  vector of values of the first principal component and  $\hat{\mathbf{y}}_2$  is the  $n \times 1$  vector of values of the second principal component.

In the biplot, each *row* of the data matrix, or item, is represented by the point located by the pair of values of the principal components. The *i*th *column* of the data matrix, or variable, is represented as an arrow from the origin to the point with coordinates  $(e_{1i}, e_{2i})$ , the entries in the *i*th column of the second matrix  $[\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2]'$  in the approximation (12-47). This scale may not be compatible with that of the principal components, so an arbitrary multiplier can be introduced that adjusts all of the vectors by the same amount.

The idea of a biplot, to represent both units and variables in the same plot, extends to canonical correlation analysis, multidimensional scaling, and even more complicated nonlinear techniques. (See [12].)

ents (the coordinates of principal components) as a vector. Usually, we have gths of all of the sampling units. To interpret the variable.

lar value decomposition mean corrected

(12-45)

e columns are the  
,  $\hat{e}_p]$ . Multiplying

(12-46)

$\dots, \hat{y}_{jp}]$

s,  $\mathbf{U}\mathbf{A}$  contains all  
s the coefficients

replacing  $\mathbf{A}$  by  
theorem, was es-

(12-47)

ment and  $\hat{y}_2$  is the

ed by the point located in the second column of the data matrix. The point with coordinates  $[\hat{e}_1, \hat{e}_2]'$  in the plot of the principal components all of the vectors

the same plot, ex-  
g, and even more

**Example 12.19 (A biplot of universities and their characteristics)** Table 12.9 gives the data on some universities for certain variables used to compare or rank major universities. These variables include  $X_1$  = average SAT score of new freshmen,  $X_2$  = percentage of new freshmen in top 10% of high school class,  $X_3$  = percentage of applicants accepted,  $X_4$  = student-faculty ratio,  $X_5$  = estimated annual expenses and  $X_6$  = graduation rate (%).

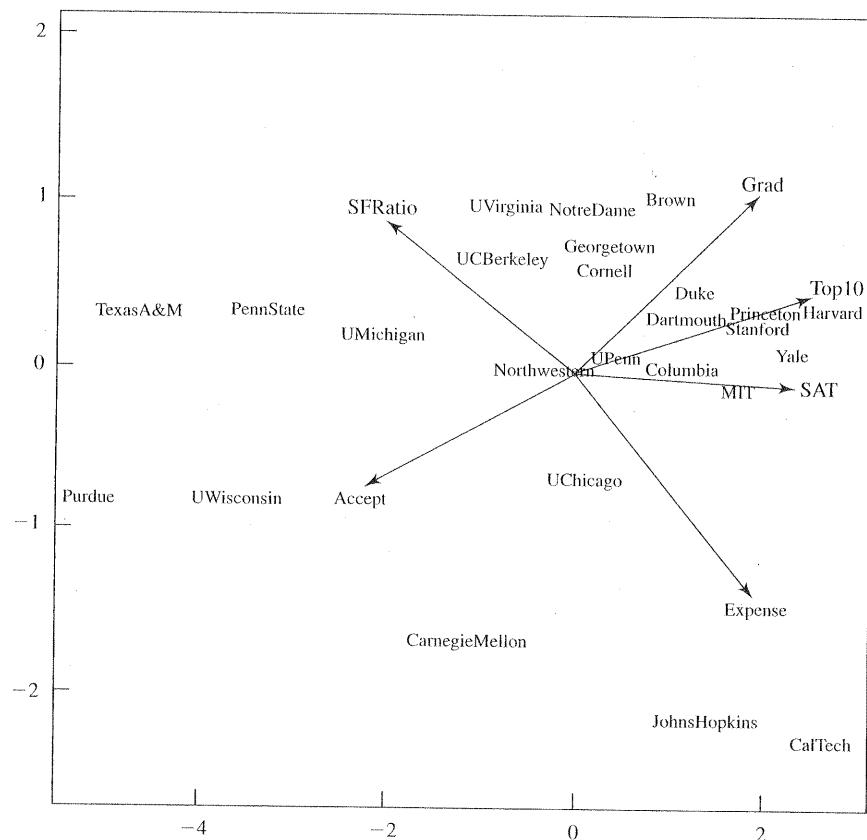
Because two of the variables, SAT and Expenses, are on a much different scale from that of the other variables, we standardize the data and base our biplot on the matrix of standardized observations  $\mathbf{z}_j$ . The biplot is given in Figure 12.24 on page 730.

Notice how Cal Tech and Johns Hopkins are off by themselves; the variable Expense is mostly responsible for this positioning. The large state universities in our sample are to the left in the biplot, and most of the private schools are on the right.

**Table 12.9** Data on Universities

University	SAT	Top10	Accept	SFRatio	Expenses	Grad
Harvard	14.00	91	14	11	39.525	97
Princeton	13.75	91	14	8	30.220	95
Yale	13.75	95	19	11	43.514	96
Stanford	13.60	90	20	12	36.450	93
MIT	13.80	94	30	10	34.870	91
Duke	13.15	90	30	12	31.585	95
CalTech	14.15	100	25	6	63.575	81
Dartmouth	13.40	89	23	10	32.162	95
Brown	13.10	89	22	13	22.704	94
JohnsHopkins	13.05	75	44	7	58.691	87
UChicago	12.90	75	50	13	38.380	87
UPenn	12.85	80	36	11	27.553	90
Cornell	12.80	83	33	13	21.864	90
Northwestern	12.60	85	39	11	28.052	89
Columbia	13.10	76	24	12	31.510	88
NotreDame	12.55	81	42	13	15.122	94
UVirginia	12.25	77	44	14	13.349	92
Georgetown	12.55	74	24	12	20.126	92
CarnegieMellon	12.60	62	59	9	25.026	72
UMichigan	11.80	65	68	16	15.470	85
UCBerkeley	12.40	95	40	17	15.140	78
UWisconsin	10.85	40	69	15	11.857	71
PennState	10.81	38	54	18	10.185	80
Purdue	10.05	28	90	19	9.066	69
TexasA&M	10.75	49	67	25	8.704	67

Source: *U.S. News & World Report*, September 18, 1995, p. 126.



**Figure 12.24** A biplot of the data on universities.

Large values for the variables SAT, Top10, and Grad are associated with the private school group. Northwestern lies in the middle of the biplot.

A newer version of the biplot, due to Gower and Hand [12], has some advantages. Their biplot, developed as an extension of the scatter plot, has features that make it easier to interpret.

- The two axes for the principal components are suppressed.
- An axis is constructed for each variable and a scale is attached.

As in the original biplot, the  $i$ -th item is located by the corresponding pair of values of the first two principal components

$$(\hat{y}_{1i}, \hat{y}_{2i}) = ((\mathbf{x}_i - \bar{\mathbf{x}})' \hat{\mathbf{e}}_1, (\mathbf{x}_i - \bar{\mathbf{x}})' \hat{\mathbf{e}}_2)$$

where  $\hat{\mathbf{e}}_1$  and where  $\hat{\mathbf{e}}_2$  are the first two eigenvectors of  $\mathbf{S}$ . The scales for the principal components are not shown on the graph.

In addition the arrows for the variables in the original biplot are replaced by axes that extend in both directions and that have scales attached. As was the case with the arrows, the axis for the  $i$ -th variable is determined by the  $i$ -th row of  $\mathbf{E} = [\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2]$ .

To begin, we let  $\mathbf{u}_i$  the vector with 1 in the  $i$ -th position and 0's elsewhere. Then an arbitrary  $p \times 1$  vector  $\mathbf{x}$  can be expressed as

$$\mathbf{x} = \sum_{i=1}^p x_i \mathbf{u}_i$$

and, by Definition 2.A.12, its projection onto the space of the first two eigenvectors has coefficient vector

$$\hat{\mathbf{E}}' \mathbf{x} = \sum_{i=1}^p x_i (\hat{\mathbf{E}}' \mathbf{u}_i)$$

so the contribution of the  $i$ -th variable to the vector sum is  $x_i (\hat{\mathbf{E}}' \mathbf{u}_i) = x_i [e_{1i}, e_{2i}]'$ . The two entries  $e_{1i}$  and  $e_{2i}$  in the  $i$ -th row of  $\hat{\mathbf{E}}$  determine the direction of the axis for the  $i$ -th variable.

The projection vector of the sample mean  $\bar{\mathbf{x}} = \sum_{i=1}^p \bar{x}_i \mathbf{u}_i$

$$\hat{\mathbf{E}}' \bar{\mathbf{x}} = \sum_{i=1}^p \bar{x}_i (\hat{\mathbf{E}}' \mathbf{u}_i)$$

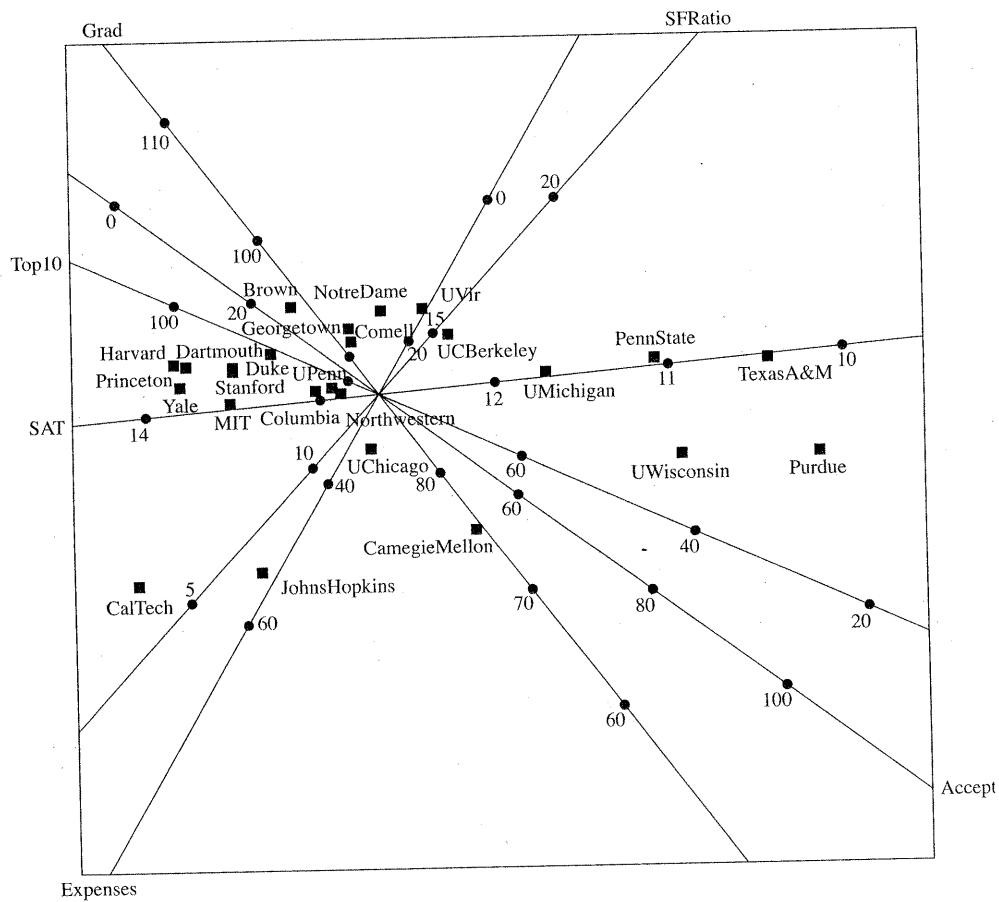
is the origin of the biplot. Every  $\mathbf{x}$  can also be written as  $\mathbf{x} = \bar{\mathbf{x}} + (\mathbf{x} - \bar{\mathbf{x}})$  and its projection vector has two components

$$\sum_{i=1}^p \bar{x}_i (\hat{\mathbf{E}}' \mathbf{u}_i) + \sum_{i=1}^p (x_i - \bar{x}_i) (\hat{\mathbf{E}}' \mathbf{u}_i)$$

Starting from the origin, the points in the direction  $w[e_{1i}, e_{2i}]'$  are plotted for  $w = 0, \pm 1, \pm 2, \dots$ . This provides a scale for the mean centered variable  $x_i - \bar{x}_i$ . It defines the distance in the biplot for a change of one unit in  $x_i$ . But, the origin for the  $i$ -th variable corresponds to  $w = 0$  because the term  $\bar{x}_i (\hat{\mathbf{E}}' \mathbf{u}_i)$  was ignored. The axis label needs to be translated so that the value  $\bar{x}_i$  is at the origin of the biplot. Since  $\bar{x}_i$  is typically not an integer (or another nice number), an integer (or other nice number) closest to it can be chosen and the scale translated appropriately. Computer software simplifies this somewhat difficult task.

The scale allows us to visually interpolate the position of  $x_i [e_{1i}, e_{2i}]'$  in the biplot. The scales predict the values of a variable, not give its exact value, as they are based on a two dimensional approximation.

**Example 12.20 (An alternative biplot for the university data)** We illustrate this newer biplot with the university data in Table 12.9. The alternative biplot with an axis for each variable is shown in Figure 12.25. Compared with Figure 12.24, the software reversed the direction of the first principal component. Notice, for example, that expenses and student faculty ratio separate Cal Tech and Johns Hopkins from the other universities. Expenses for Cal Tech and Johns Hopkins can be seen to be about 57 thousand a year, and the student faculty ratios are in the single digits. The large state universities, on the right hand side of the plot, have relatively high student faculty ratios, above 20, relatively low SAT scores of entering freshman, and only about 50% or fewer of their entering students in the top 10% of their high school class. The scaled axes on the newer biplot are more informative than the arrows in the original biplot.



**Figure 12.25** An alternative biplot of the data on universities.

See le Roux and Gardner [23] for more examples of this alternative biplot and references to appropriate special purpose statistical software.

## 12.9 Procrustes Analysis: A Method for Comparing Configurations

Starting with a given  $n \times n$  matrix of distances  $\mathbf{D}$ , or similarities  $\mathbf{S}$ , that relate  $n$  objects, two or more configurations can be obtained using different techniques. The possible methods include both metric and nonmetric multidimensional scaling. The question naturally arises as to how well the solutions coincide. Figures 12.19 and 12.20 in Example 12.16 respectively give the metric multidimensional scaling (principal coordinate analysis) and nonmetric multidimensional scaling solutions for the data on universities. The two configurations appear to be quite similar, but a quantitative measure would be useful. A numerical comparison of two configurations, obtained by moving one configuration so that it aligns best with the other, is called *Procrustes analysis*, after the innkeeper Procrustes, in Greek mythology, who would either stretch or lop off customers' limbs so they would fit his bed.

## Constructing the Procrustes Measure of Agreement

Suppose the  $n \times p$  matrix  $\mathbf{X}^*$  contains the coordinates of the  $n$  points obtained for plotting with technique 1 and the  $n \times q$  matrix  $\mathbf{Y}^*$  contains the coordinates from technique 2, where  $q \leq p$ . By adding columns of zeros to  $\mathbf{Y}^*$ , if necessary, we can assume that  $\mathbf{X}^*$  and  $\mathbf{Y}^*$  both have the same dimension  $n \times p$ . To determine how compatible the two configurations are, we move, say, the second configuration to match the first by shifting each point by the same amount and rotating or reflecting the configuration about the coordinate axes.<sup>4</sup>

Mathematically, we translate by a vector  $\mathbf{b}$  and multiply by an orthogonal matrix  $\mathbf{Q}$  so that the coordinates of the  $j$ th point  $\mathbf{y}_j$  are transformed to

$$\mathbf{Q}\mathbf{y}_j + \mathbf{b}$$

The vector  $\mathbf{b}$  and orthogonal matrix  $\mathbf{Q}$  are then varied to order to minimize the sum, over all  $n$  points, of squared distances

$$d_j^2(\mathbf{x}_j, \mathbf{Q}\mathbf{y}_j + \mathbf{b}) = (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b})'(\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b}) \quad (12-48)$$

between  $\mathbf{x}_j$  and the transformed coordinates  $\mathbf{Q}\mathbf{y}_j + \mathbf{b}$  obtained for the second technique. We take, as a measure of fit, or agreement, between the two configurations, the residual sum of squares

$$PR^2 = \min_{\mathbf{Q}, \mathbf{b}} \sum_{j=1}^n (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b})'(\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b}) \quad (12-49)$$

The next result shows how to evaluate this Procrustes residual sum of squares measure of agreement and determines the *Procrustes rotation* of  $\mathbf{Y}^*$  relative to  $\mathbf{X}^*$ .

**Result 12.2** Let the  $n \times p$  configurations  $\mathbf{X}^*$  and  $\mathbf{Y}^*$  both be centered so that all columns have mean zero. Then

$$\begin{aligned} PR^2 &= \sum_{j=1}^n \mathbf{x}'_j \mathbf{x}_j + \sum_{j=1}^n \mathbf{y}'_j \mathbf{y}_j - 2 \sum_{i=1}^p \lambda_i \\ &= \text{tr}[\mathbf{X}^* \mathbf{X}^{*\prime}] + \text{tr}[\mathbf{Y}^* \mathbf{Y}^{*\prime}] - 2 \text{tr}[\Lambda] \end{aligned} \quad (12-50)$$

where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  and the minimizing transformation is

$$\hat{\mathbf{Q}} = \sum_{i=1}^p \mathbf{v}_i \mathbf{u}_i' = \mathbf{V} \mathbf{U}' \quad \hat{\mathbf{b}} = \mathbf{0} \quad (12-51)$$

<sup>4</sup>Sibson [30] has proposed a numerical measure of the agreement between two configurations, given by the coefficient

$$\gamma = 1 - \frac{[\text{tr}(\mathbf{Y}^* \mathbf{X}^* \mathbf{X}^{*\prime} \mathbf{Y}^*)^{1/2}]^2}{\text{tr}(\mathbf{X}^* \mathbf{X}^*) \text{tr}(\mathbf{Y}^* \mathbf{Y}^*)}$$

For identical configurations,  $\gamma = 0$ . If necessary,  $\gamma$  can be computed after a Procrustes analysis has been completed.

Here  $\Lambda$ ,  $\mathbf{U}$ , and  $\mathbf{V}$  are obtained from the singular-value decomposition

$$\sum_{j=1}^n \mathbf{y}_j \mathbf{x}'_j = \mathbf{Y}^{*'} \mathbf{X}^* = \mathbf{U} \Lambda \mathbf{V}'$$

**Proof.** Because the configurations are centered to have zero means ( $\sum_{j=1}^n \mathbf{x}_j = \mathbf{0}$  and  $\sum_{j=1}^n \mathbf{y}_j = \mathbf{0}$ ), we have

$$\sum_{j=1}^n (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b})' (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j - \mathbf{b}) = \sum_{j=1}^n (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j)' (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j) + n\mathbf{b}'\mathbf{b}$$

The last term is nonnegative, so the best fit occurs for  $\hat{\mathbf{b}} = \mathbf{0}$ . Consequently, we need only consider

$$PR^2 = \min_{\mathbf{Q}} \sum_{j=1}^n (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j)' (\mathbf{x}_j - \mathbf{Q}\mathbf{y}_j) = \sum_{j=1}^n \mathbf{x}'_j \mathbf{x}_j + \sum_{j=1}^n \mathbf{y}'_j \mathbf{y}_j - 2 \max_{\mathbf{Q}} \sum_{j=1}^n \mathbf{x}'_j \mathbf{Q} \mathbf{y}_j$$

Using  $\mathbf{x}'_j \mathbf{Q} \mathbf{y}_j = \text{tr}[\mathbf{Q} \mathbf{y}_j \mathbf{x}'_j]$ , we find that the expression being maximized becomes

$$\sum_{j=1}^n \mathbf{x}'_j \mathbf{Q} \mathbf{y}_j = \sum_{j=1}^n \text{tr}[\mathbf{Q} \mathbf{y}_j \mathbf{x}'_j] = \text{tr}\left[\mathbf{Q} \left(\sum_{j=1}^n \mathbf{y}_j \mathbf{x}'_j\right)\right]$$

By the singular-value decomposition,

$$\sum_{j=1}^n \mathbf{y}_j \mathbf{x}'_j = \mathbf{Y}^{*'} \mathbf{X}^* = \mathbf{U} \Lambda \mathbf{V}' = \sum_{j=1}^p \lambda_i \mathbf{u}_i \mathbf{v}'_i$$

where  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p]$  and  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p]$  are  $p \times p$  orthogonal matrices. Consequently,

$$\sum_{j=1}^n \mathbf{x}'_j \mathbf{Q} \mathbf{y}_j = \text{tr}\left[\mathbf{Q} \left(\sum_{i=1}^p \lambda_i \mathbf{u}_i \mathbf{v}'_i\right)\right] = \sum_{i=1}^p \lambda_i \text{tr}[\mathbf{Q} \mathbf{u}_i \mathbf{v}'_i]$$

The variable quantity in the  $i$ th term

$$\text{tr}[\mathbf{Q} \mathbf{u}_i \mathbf{v}'_i] = \mathbf{v}'_i \mathbf{Q} \mathbf{u}_i$$

has an upper bound of 1 as can be seen by applying the Cauchy-Schwarz inequality (2-48) with  $\mathbf{b} = \mathbf{Q} \mathbf{v}_i$  and  $\mathbf{d} = \mathbf{u}_i$ . That is, since  $\mathbf{Q}$  is orthogonal,

$$\mathbf{v}'_i \mathbf{Q} \mathbf{u}_i \leq \sqrt{\mathbf{v}'_i \mathbf{Q} \mathbf{Q}' \mathbf{v}_i} \sqrt{\mathbf{u}'_i \mathbf{u}_i} = \sqrt{\mathbf{v}'_i \mathbf{v}_i} \times 1 = 1$$

ition

$$\text{means } \left( \sum_{j=1}^n \mathbf{x}_j = \mathbf{0} \right)$$

$$\mathbf{Q}\mathbf{y}_j) + n\mathbf{b}'\mathbf{b}$$

equently, we need

$$2 \max_{\mathbf{Q}} \sum_{j=1}^n \mathbf{x}_j' \mathbf{Q} \mathbf{y}_j$$

lized becomes

ogonal matrices.

Each of these  $p$  terms can be maximized by the same choice  $\mathbf{Q} = \mathbf{VU}'$ . With this choice,

$$\mathbf{v}_i' \mathbf{Q} \mathbf{u}_i = \mathbf{v}_i' \mathbf{VU}' \mathbf{u}_i = [0, \dots, 0, 1, 0, \dots, 0] \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = 1$$

Therefore,

$$-2 \max_{\mathbf{Q}} \sum_{j=1}^n \mathbf{x}_j' \mathbf{Q} \mathbf{y}_j = -2(\lambda_1 + \lambda_2 + \dots + \lambda_p)$$

Finally, we verify that  $\mathbf{Q}\mathbf{Q}' = \mathbf{VU}'\mathbf{UV}' = \mathbf{V}\mathbf{I}_p\mathbf{V}' = \mathbf{I}_p$ , so  $\mathbf{Q}$  is a  $p \times p$  orthogonal matrix, as required. ■

---

**Example 12.21 (Procrustes analysis of the data on universities)** Two configurations, produced by metric and nonmetric multidimensional scaling, of data on universities are given Example 12.16. The two configurations appear to be quite close. There is a two-dimensional array of coordinates for each of the two scaling methods. Initially, the sum of squared distances is

$$\sum_{j=1}^{25} (\mathbf{x}_j - \mathbf{y}_j)' (\mathbf{x}_j - \mathbf{y}_j) = 3.862$$

A computer calculation gives

$$\mathbf{U} = \begin{bmatrix} -.9990 & .0448 \\ .0448 & .9990 \end{bmatrix} \quad \mathbf{V} = \begin{bmatrix} -1.0000 & .0076 \\ .0076 & 1.0000 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 114.9439 & 0.000 \\ 0.000 & 21.3673 \end{bmatrix}$$

According to Result 12.2, to better align these two solutions, we multiply the nonmetric scaling solution by the orthogonal matrix

$$\hat{\mathbf{Q}} = \sum_{i=1}^2 \mathbf{v}_i \mathbf{u}_i = \mathbf{VU}' = \begin{bmatrix} .9993 & -.0372 \\ .0372 & .9993 \end{bmatrix}$$

This corresponds to clockwise rotation of the nonmetric solution by about 2 degrees. After rotation, the sum of squared distances, 3.862, is reduced to the Procrustes measure of fit

$$PR^2 = \sum_{j=1}^{25} \mathbf{x}_j' \mathbf{x}_j + \sum_{j=1}^{25} \mathbf{y}_j' \mathbf{y}_j - 2 \sum_{j=1}^2 \lambda_i = 3.673$$

**Example 12.22 (Procrustes analysis and additional ordinations of data on forests)**

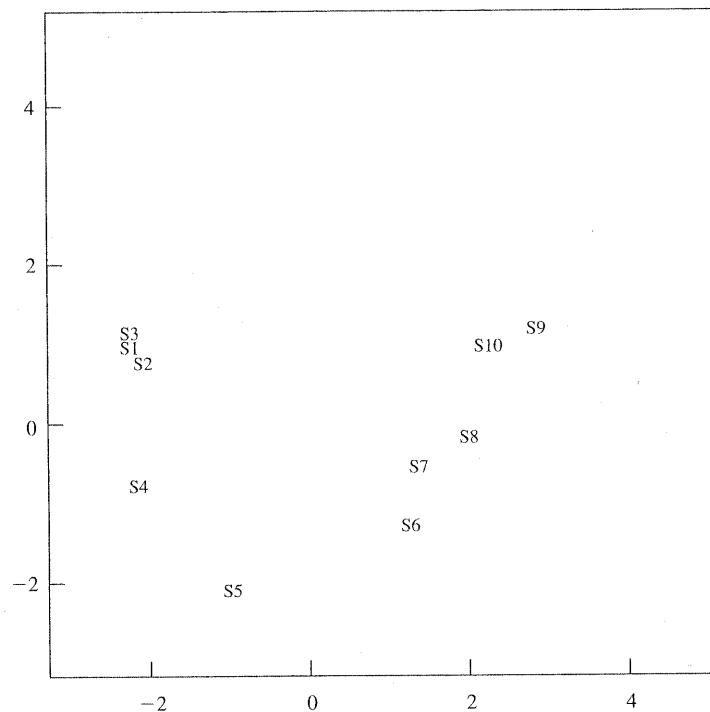
Data were collected on the populations of eight species of trees growing on ten upland sites in southern Wisconsin. These data are shown in Table 12.10.

The metric, or principal coordinate, solution and nonmetric multidimensional scaling solution are shown in Figures 12.26 and 12.27.

**Table 12.10** Wisconsin Forest Data

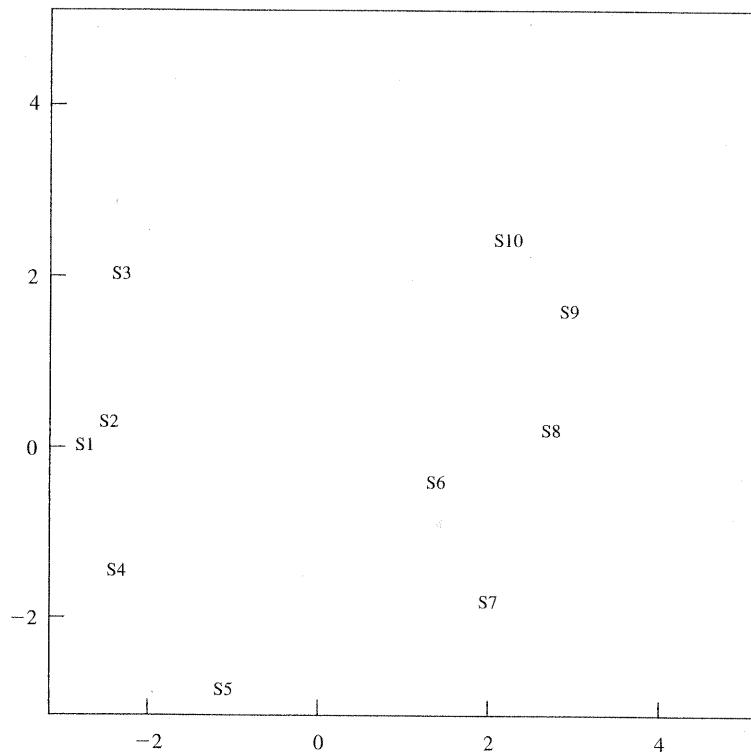
Tree	Site									
	1	2	3	4	5	6	7	8	9	10
BurOak	9	8	3	5	6	0	5	0	0	0
BlackOak	8	9	8	7	0	0	0	0	0	0
WhiteOak	5	4	9	9	7	7	4	6	0	2
RedOak	3	4	0	6	9	8	7	6	4	3
AmericanElm	2	2	4	5	6	0	5	0	2	5
Basswood	0	0	0	0	2	7	6	6	7	6
Ironwood	0	0	0	0	0	0	7	4	6	5
SugarMaple	0	0	0	0	0	5	4	8	8	9

Source: See [24].

**Figure 12.26** Metric multidimensional scaling of the data on forests.

(data on forests)  
s growing on ten  
12.10.  
multidimensional

3	9	10
0	0	
0	0	
0	2	
4	3	
2	5	
7	6	
6	5	
8	9	



**Figure 12.27** Nonmetric multidimensional scaling of the data on forests.

Using the coordinates of the points in Figures 12.26 and 12.27, we obtain the initial sum of squared distances for fit:

$$\sum_{j=1}^{10} (\mathbf{x}_j - \mathbf{y}_j)' (\mathbf{x}_j - \mathbf{y}_j) = 8.547$$

A computer calculation gives

$$\mathbf{U} = \begin{bmatrix} -.9833 & -.1821 \\ -.1821 & .9833 \end{bmatrix} \quad \mathbf{V} = \begin{bmatrix} -1.0000 & -.0001 \\ -.0001 & 1.0000 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 43.3748 & 0.0000 \\ 0.0000 & 14.9103 \end{bmatrix}$$

According to Result 12.2, to better align these two solutions, we multiply the non-metric scaling solution by the orthogonal matrix

$$\hat{\mathbf{Q}} = \sum_{i=1}^2 \mathbf{v}_i \mathbf{u}_i' = \mathbf{V} \mathbf{U}' = \begin{bmatrix} .9833 & .1821 \\ -.1821 & .9833 \end{bmatrix}$$

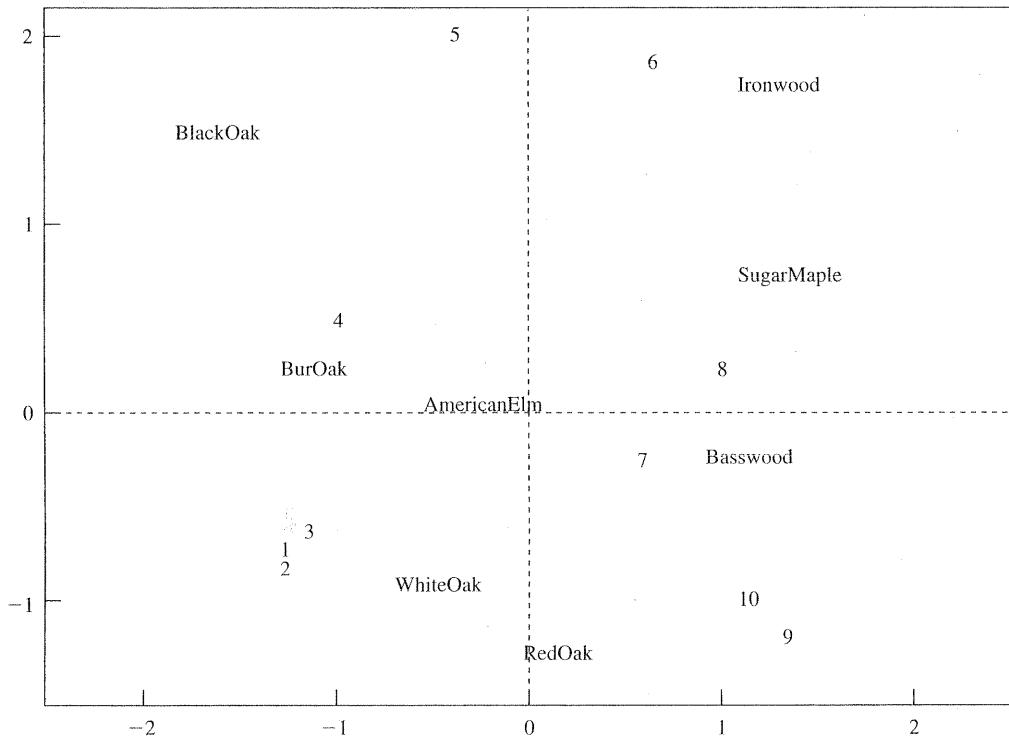
This corresponds to clockwise rotation of the nonmetric solution by about 10 degrees. After rotation, the sum of squared distances, 8.547, is reduced to the Procrustes measure of fit

$$PR^2 = \sum_{j=1}^{10} \mathbf{x}_j' \mathbf{x}_j + \sum_{j=1}^{10} \mathbf{y}_j' \mathbf{y}_j - 2 \sum_{i=1}^2 \lambda_i = 6.599$$

We note that the sampling sites seem to fall along a curve in both pictures. This could lead to a one-dimensional *nonlinear ordination* of the data. A quadratic or other curve could be fit to the points. By adding a scale to the curve, we would obtain a one-dimensional ordination.

It is informative to view the Wisconsin forest data when both sampling units and variables are shown. A correspondence analysis applied to the data produces the plot in Figure 12.28. The biplot is shown in Figure 12.29.

All of the plots tell similar stories. Sites 1–5 tend to be associated with species of oak trees, while sites 7–10 tend to be associated with basswood, ironwood, and sugar maples. American elm trees are distributed over most sites, but are more closely associated with the lower numbered sites. There is almost a continuum of sites distinguished by the different species of trees. ■



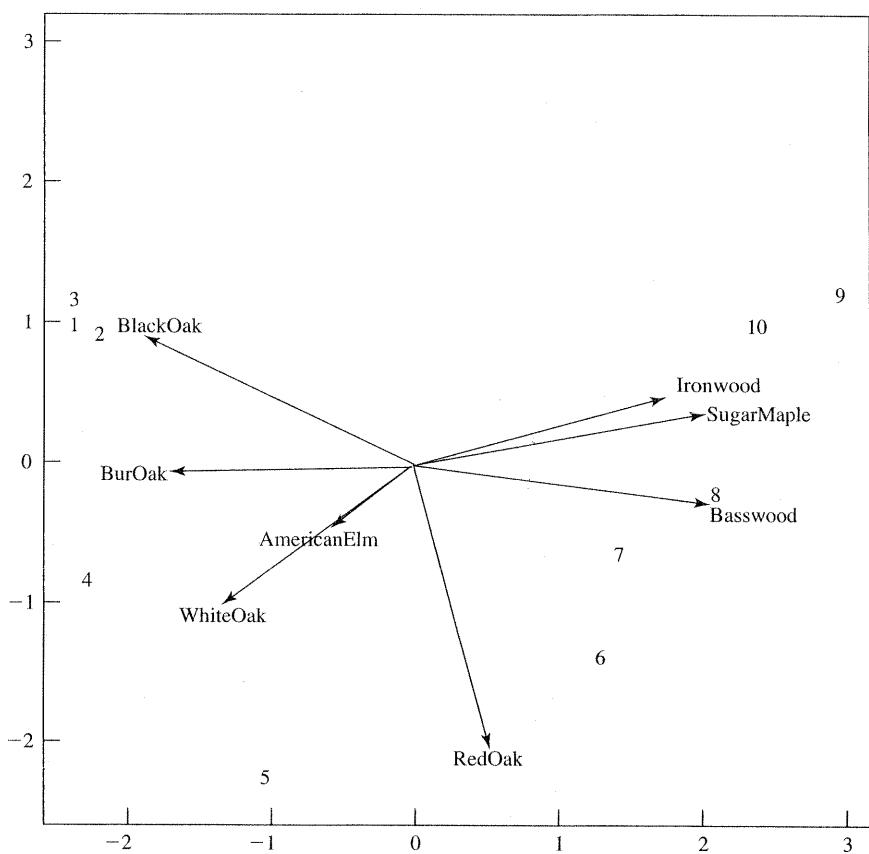
**Figure 12.28** The correspondence analysis plot of the data on forests.

about 10 degrees.  
to the Procrustes

oth pictures. This  
a. A quadratic or  
curve, we would

ampling units and  
ata produces the

ed with species of  
nwood, and sugar  
are more closely  
ontinuum of sites



**Figure 12.29** The biplot of the data on forests.