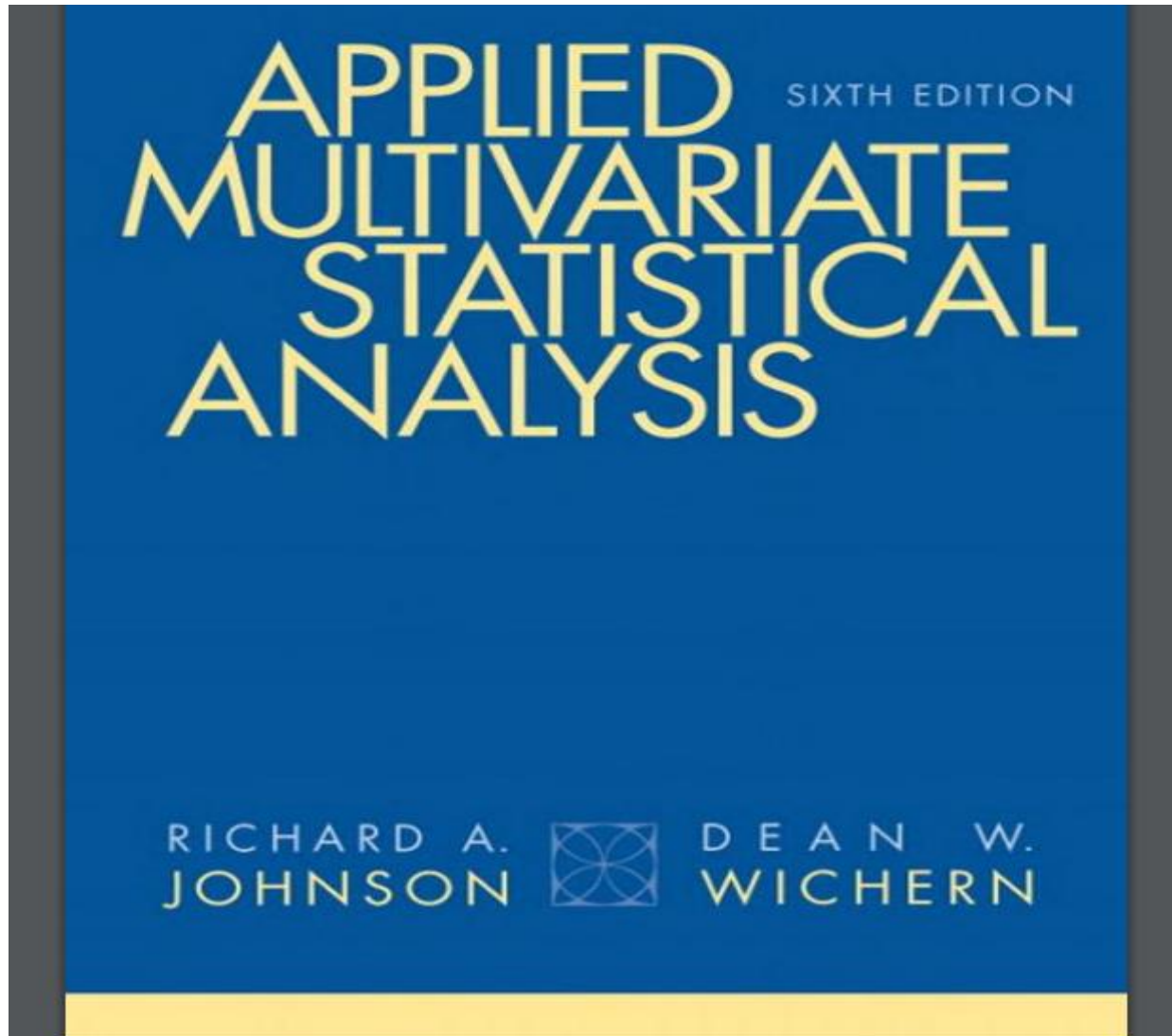


# Advanced Multivariate Methods



# Clustering, Distance Measures and Ordination- Similar Distance Between Objects of Unit of Analysis, i.e. Study Subjects

## 12.1 Introduction

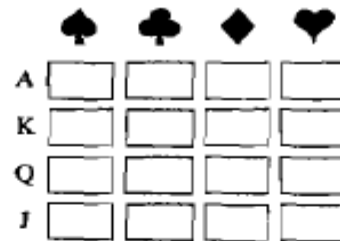
Rudimentary, exploratory procedures are often quite helpful in understanding the complex nature of multivariate relationships. For example, throughout this book, we have emphasized the value of data plots. In this chapter, we shall discuss some additional displays based on certain measures of distance 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.

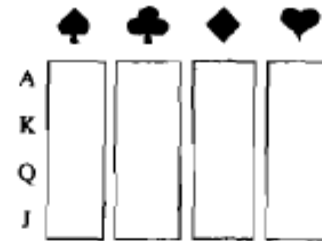
To illustrate the nature of the difficulty in defining a natural grouping, consider sorting the 16 face cards in an ordinary deck of playing cards into clusters of similar objects. Some groupings are illustrated in Figure 12.1. It is immediately clear that meaningful partitions depend on the definition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

# Grouping Face Cards: “Good” vs. “Bad”

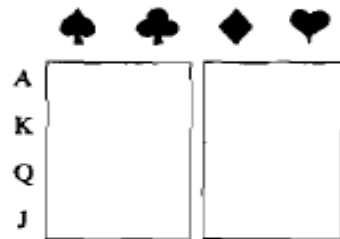
In most practical applications of cluster analysis, the investigator knows enough about the problem to distinguish “good” groupings from “bad” groupings. Why not enumerate all possible groupings and select the “best” ones for further study?



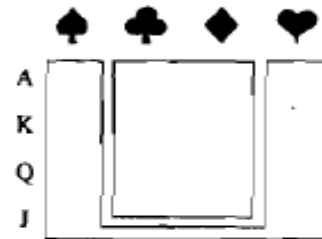
(a) Individual cards



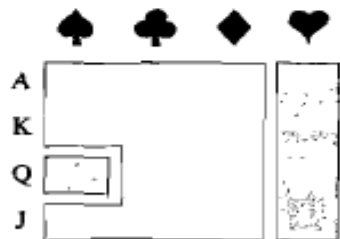
(b) Individual suits



(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.

# Algorithm for Sorting Objects into Groups: Association Between Objects Based on a Criterion Variable

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.

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.

# Similarity Measures: Distance and Similarity Coefficients for Pairs of Items

## 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),

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

# Similarity Measures: Distance and Similarity Coefficients for Pairs of Items

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.

# Measures of Distance or Dissimilarity

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

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

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

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.

# Clustering and Coding Using Binary or Dichotomous Variables

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.



# Large Distance Corresponds to Many Mismatches – Dissimilar Items

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

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

# Frequency Table of Matches and Mismatches

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 $k$			
	1	0	Totals	
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  $b = c = d = 1$ .

# Similarity Coefficients for Clustering Items

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

<b>Table 12.1</b> 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).]	

# Monotonicity and Clustering: Procedures Not Affected if Similarity is Changed If Relative Ordering Same or Unchanged

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  $(a_1 + d_1)/p \geq (a_{11} + d_{11})/p$ , we also have  $2(a_1 + d_1)/[2(a_1 + d_1) + b_1 + c_1] \geq 2(a_{11} + d_{11})/[2(a_{11} + d_{11}) + b_{11} + c_{11}]$ , 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: Calculating Values of Similarity Coefficient

**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

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

# Scores for Six Binary Variables

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

# Similarity Coefficient Gives Equal Weight to Matches

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. ■

# Properties of Distance: Nonnegative Definite Condition and Maximum Similarity Scales So That $\tilde{s}_{ii} = 1$

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

## 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		$a + c$	$b + d$	$n = a + b + c + d$	

# Contingency Table for Binary Variables: Calculate Product Moment Correlation

		Variable $k$		Totals	
		1	0		
Variable $i$	1	$a$	$b$	$a + b$	(12-10)
	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: Measuring Similarities of 11 Languages

**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, diereses, 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.

# Numerals One (1) to Ten (10) in Eleven Languages

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	negy	neljä
five	fem	fem	vijf	funf	cinq	cinco	cinque	piec	ot	viisi
six	seks	seks	zes	sechs	six	seis	sei	szesc	hat	kuusi
seven	sju	syv	zeven	sieben	sept	siete	sette	siedem	het	seitseman
eight	atte	otte	acht	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

# Concordant First Letters for Numbers in Eleven Languages

<b>Table 12.3 Concordant First Letters for Numbers in 11 Languages</b>											
	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	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. ■

# Hierarchical Clustering Methods

## 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.

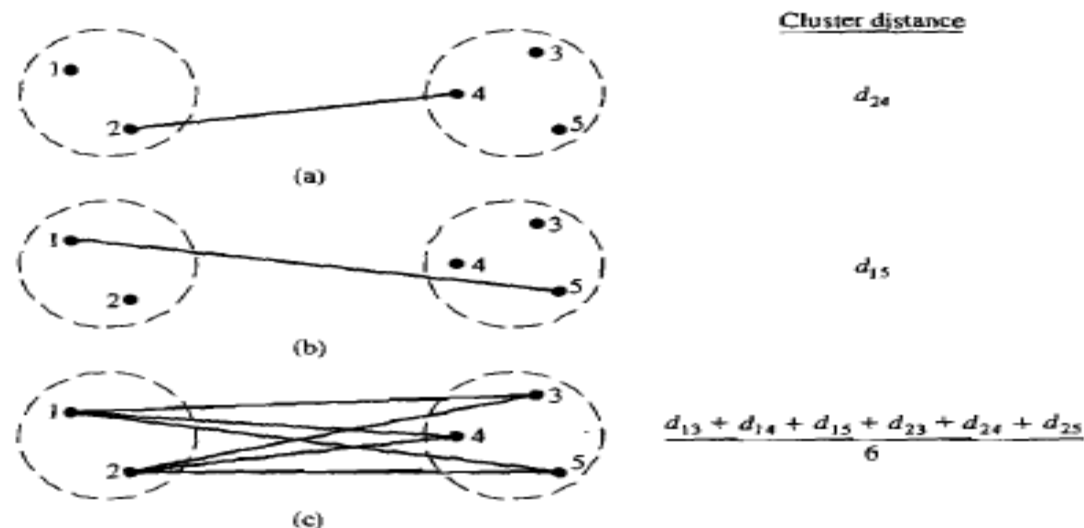
*Divisive hierarchical methods* work in the opposite direction. An initial single group of objects is divided into two subgroups such that the objects in one subgroup are “far from” the objects in the other. These subgroups are then further divided into dissimilar subgroups; the process continues until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

# Dendrogram: Displays Agglomerative and Divisive Methods

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.



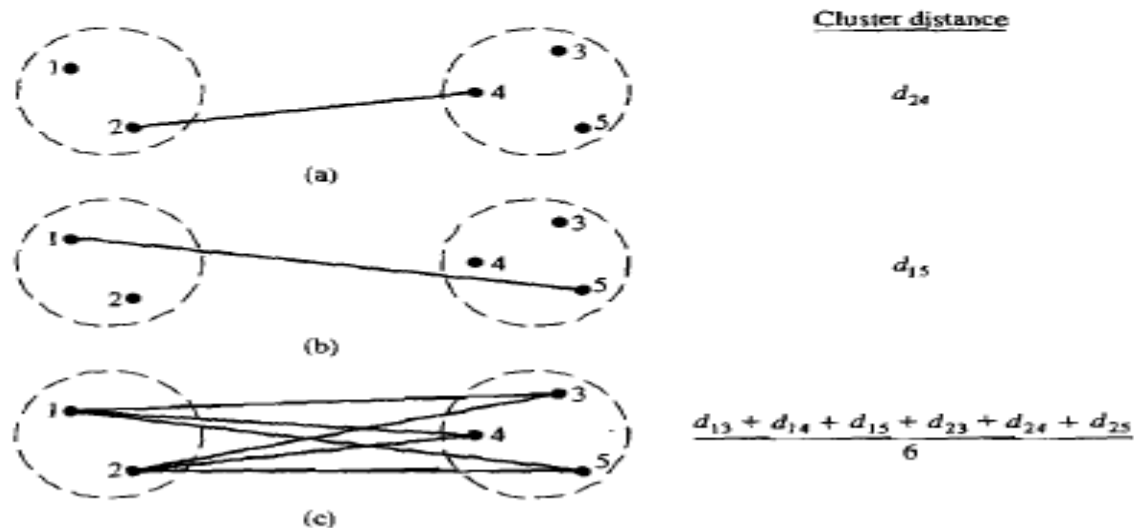
**Figure 12.2** Intercluster distance (dissimilarity) for (a) single linkage, (b) complete linkage, and (c) average linkage.

# Hierarchical Clustering: Linkage Methods (Suitable for Both Subjects and Variables)

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  $d_{UV}$ .



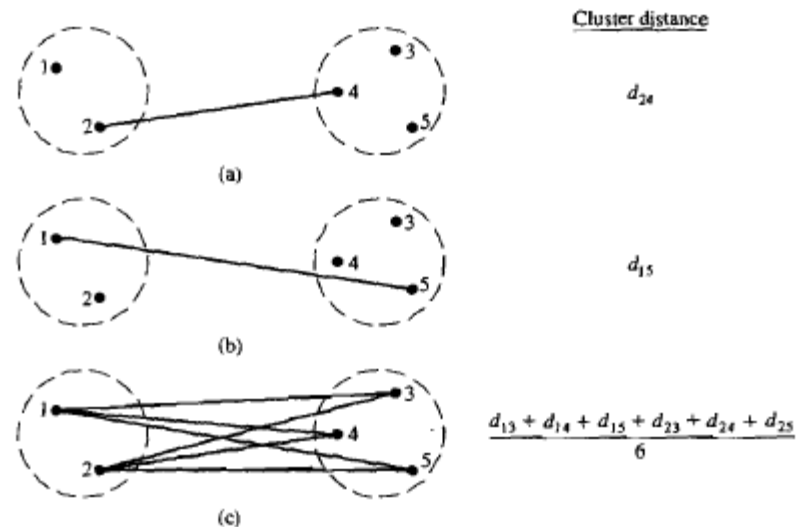
**Figure 12.2** Intercluster distance (dissimilarity) for (a) single linkage, (b) complete linkage, and (c) average linkage.



# Agglomerative Clustering Algorithm for N Objects (Sample Size $n = 30$ )

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.



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

# Single Linkage: Nearest Neighbor Connotes Smallest Distance or Largest Similarity

## 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. Dendrograms for some specific cases are considered in the following examples.

# Example: Clustering Using Single Linkage

**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:

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

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

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

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

$$d_{(35)1} = \min \{d_{31}, d_{51}\} = \min \{3, 11\} = 3$$

$$d_{(35)2} = \min \{d_{32}, d_{52}\} = \min \{7, 10\} = 7$$

$$d_{(35)4} = \min \{d_{34}, d_{54}\} = \min \{9, 8\} = 8$$

# Example: Clustering Using Single Linkage

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

$$\begin{array}{c} \begin{array}{c} (35) \\ 1 \\ 2 \\ 4 \end{array} \begin{array}{c} (35) \quad 1 \quad 2 \quad 4 \\ \left[ \begin{array}{cccc} 0 & & & \\ \textcircled{3} & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{array} \right] \end{array} \end{array}$$

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

$$d_{(135)2} = \min \{d_{(35)2}, d_{12}\} = \min \{7, 9\} = 7$$

$$d_{(135)4} = \min \{d_{(35)4}, d_{14}\} = \min \{8, 6\} = 6$$

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

$$\begin{array}{c} \begin{array}{c} (135) \\ 2 \\ 4 \end{array} \begin{array}{c} (135) \quad 2 \quad 4 \\ \left[ \begin{array}{ccc} 0 & & \\ 7 & 0 & \\ 6 & \textcircled{5} & 0 \end{array} \right] \end{array} \end{array}$$

# Minimum Nearest Neighbor Distance

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

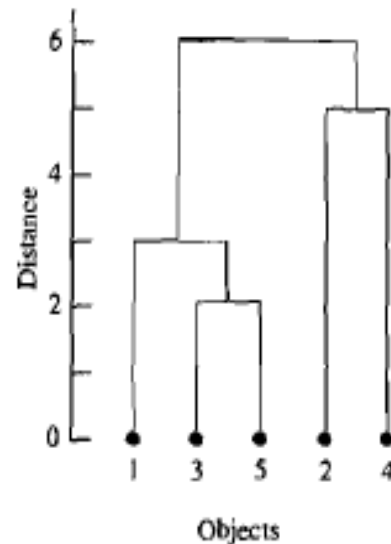
$$\begin{array}{cc} & \begin{array}{cc} (135) & (24) \end{array} \\ \begin{array}{c} (135) \\ (24) \end{array} & \left[ \begin{array}{cc} 0 & \\ \textcircled{6} & 0 \end{array} \right] \end{array}$$

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. ■

# Single Linkage Dendrogram for Distances Between Objects

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: Single Linkage Clustering of 11 Languages

**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 concordances from the perfect agreement 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	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

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

# Dendrogram: Single Linkage Clustering of 11 Languages

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,



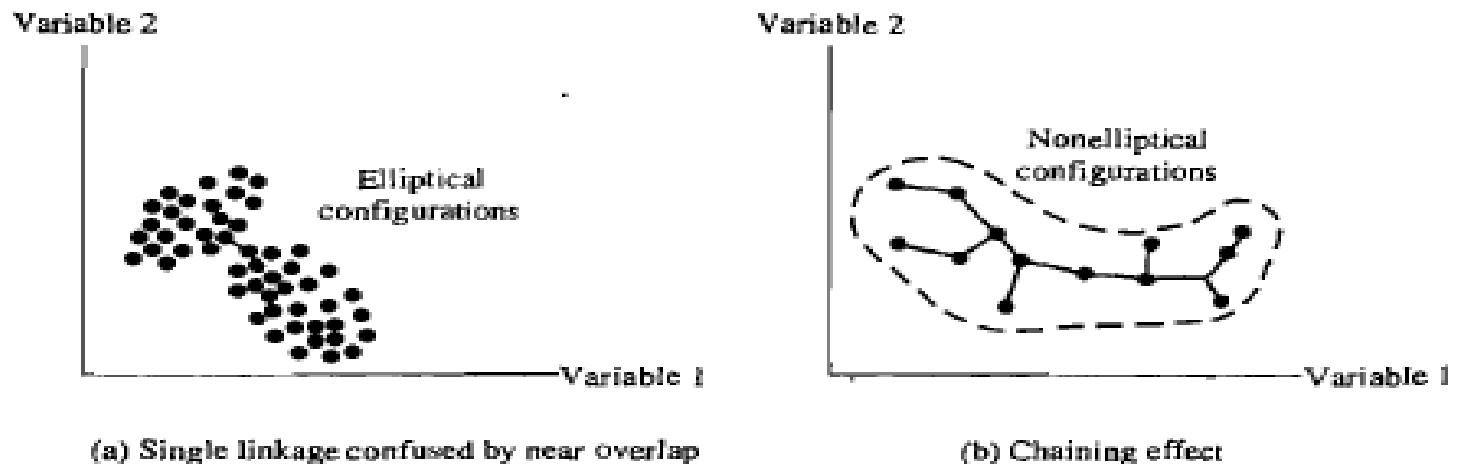
**Figure 12.4** Single linkage dendrograms for distances between numbers in 11 languages.

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. ■



# Single Linkage for Sub-Clusters: Delineates Nonellipsoidal

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 nonellipsoidal 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: Ensures Objects/Items In A Cluster Within Same Maximum or Minimum Distance

## 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

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.

# Clustering Using Complete Linkage

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

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

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

$$d_{(35)1} = \max\{d_{31}, d_{51}\} = \max\{3, 11\} = 11$$

$$d_{(35)2} = \max\{d_{32}, d_{52}\} = 10$$

$$d_{(35)4} = \max\{d_{34}, d_{54}\} = 9$$

and the modified distance matrix becomes

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

# Clustering Using Complete Linkage

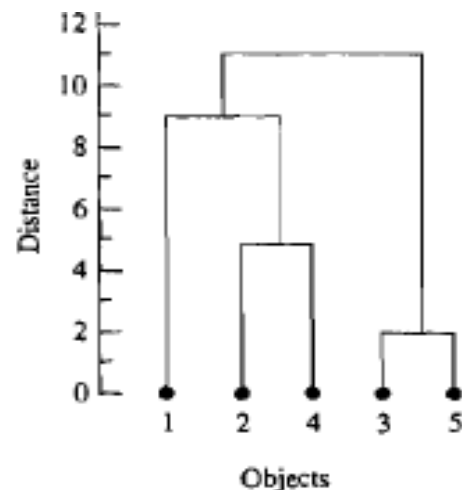
The next merger occurs between the most similar groups, 2 and 4, to give the cluster (24). At stage 3, we have

$$d_{(24)(35)} = \max \{d_{2(35)}, d_{4(35)}\} = \max \{10, 9\} = 10$$

$$d_{(24)1} = \max \{d_{21}, d_{41}\} = 9$$

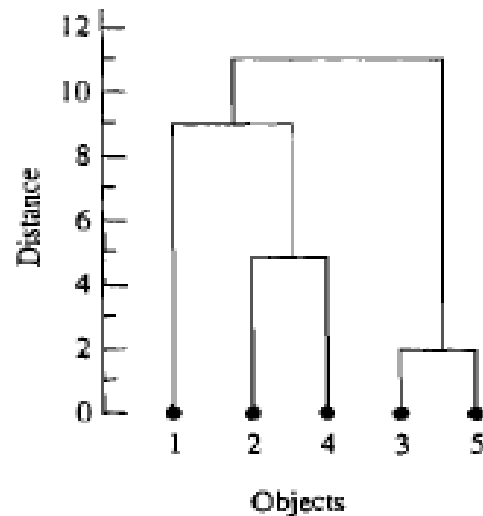
and the distance matrix

$$\begin{array}{c} (35) \quad (24) \quad 1 \\ \begin{array}{c} (35) \\ (24) \\ 1 \end{array} \left[ \begin{array}{ccc} 0 & & \\ 10 & 0 & \\ 11 & \textcircled{9} & 0 \end{array} \right] \end{array}$$



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

# Complete Linkage Dendrogram: Compare Allocation of Objects/Items



**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. ■

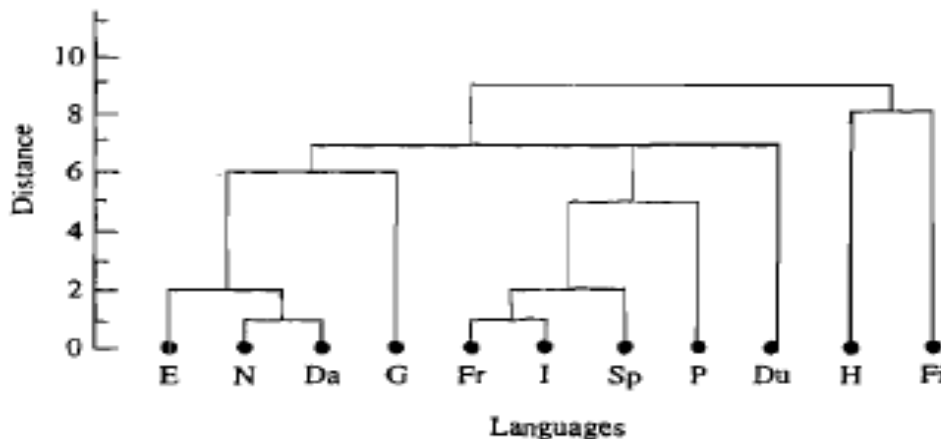
Comparing Figures 12.3 and 12.6, we see that the dendrograms for single linkage and complete linkage differ in the allocation of object 1 to previous groups.

# Example: Complete Linkage Clustering of 11 Languages

**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.

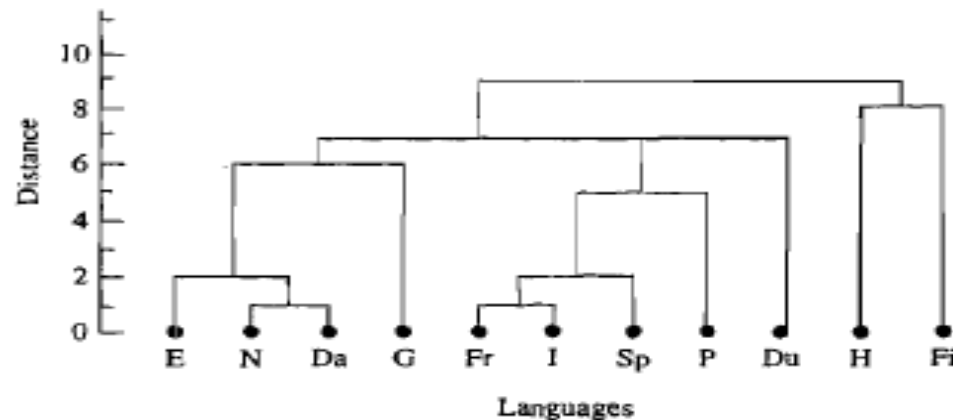
Comparing Figures 12.7 and 12.4, we see that both hierarchical methods yield the English–Norwegian–Danish and the French–Italian–Spanish language groups. Polish is merged with French–Italian–Spanish at an intermediate level. In addition, both methods merge Hungarian and Finnish only at the penultimate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

# Example: Complete Linkage Clustering of 11 Languages



**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. ■

# Clustering Variables Using Complete Linkage

**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.



# Public Utility Data (1975)

**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.

# Correlations Between Pairs of Six Variables ( n = 22 Public Utilities)

**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

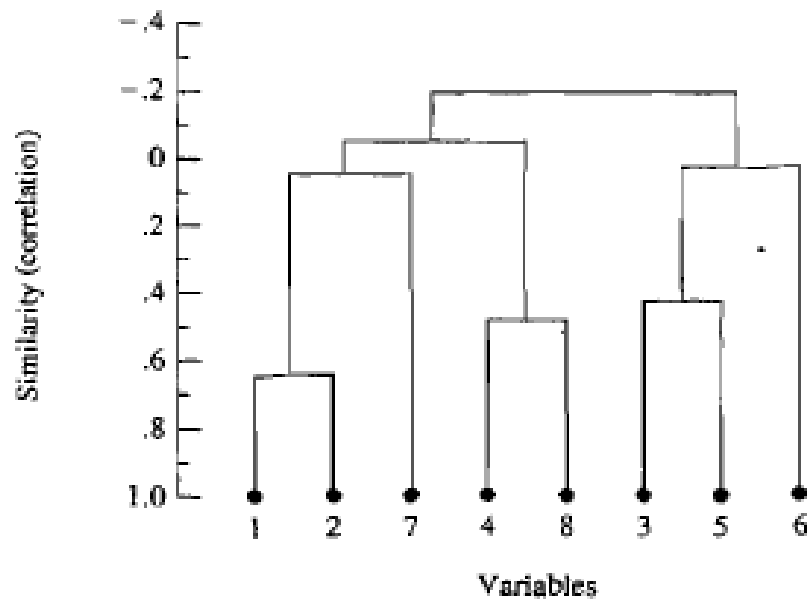
variables by the product-moment correlation coefficient. The correlation matrix is given in Table 12.5.

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. ■

# Complete Linkage Dendrogram for Similarities Among Eight Utility Variables

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 Formula

## 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.

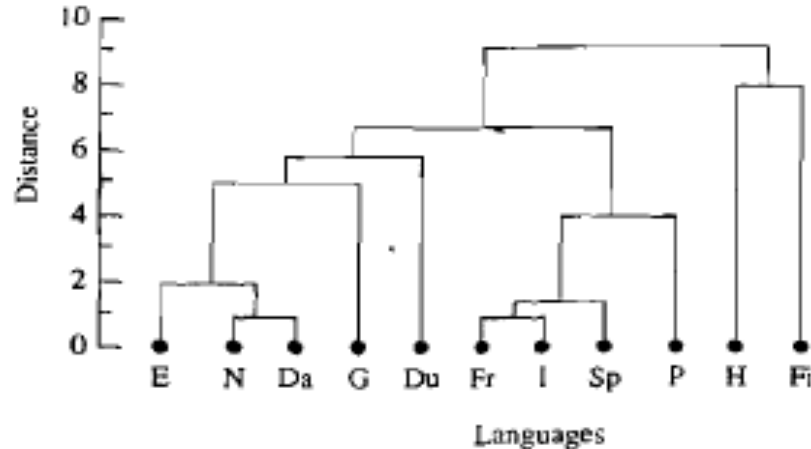
Again, the input to the average linkage algorithm may be distances or similarities, and the method can be used to group objects or variables. The average linkage algorithm proceeds in the manner of the general algorithm of (12-12). We begin by searching the distance matrix  $\mathbf{D} = \{d_{ik}\}$  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.

# Average Linkage Clustering of 11 Languages

**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. ■

# Average Linkage Clustering of Public Utilities (n = 22)

**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.

**Table 12.6** Distances Between 22 Utilities

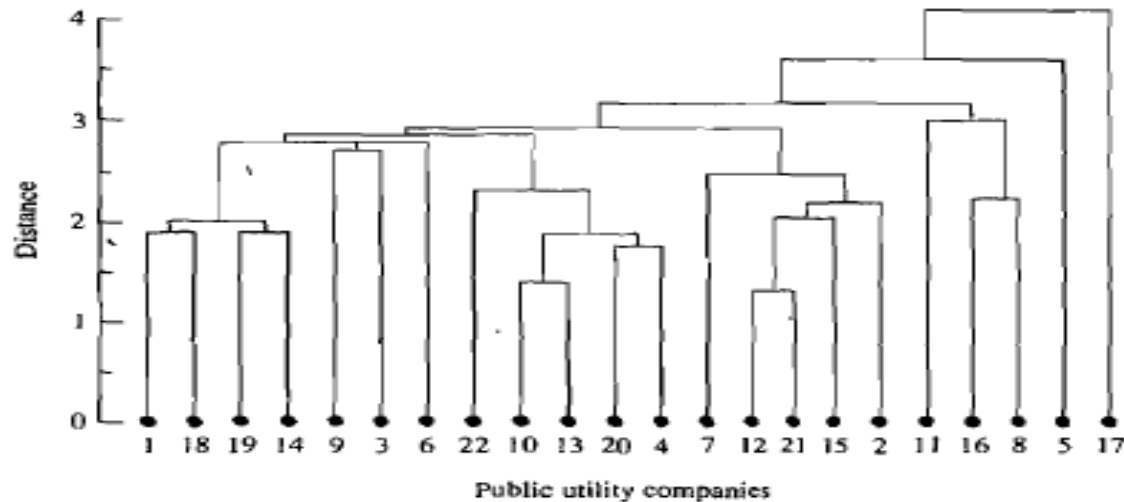
Firm no.	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
1	.00																					
2	3.10	.00																				
3	3.68	4.92	.00																			
4	2.46	2.16	4.11	.00																		
5	4.12	3.85	4.47	4.13	.00																	
6	3.61	4.22	2.99	3.20	4.60	.00																
7	3.90	3.45	4.22	3.97	4.60	3.35	.00															
8	2.74	3.89	4.99	3.69	5.16	4.91	4.36	.00														
9	3.25	3.96	2.75	3.75	4.49	3.73	2.80	3.59	.00													
10	3.10	2.71	3.93	1.49	4.05	3.83	4.51	3.67	3.57	.00												
11	3.49	4.79	5.90	4.86	6.46	6.00	6.00	3.46	5.18	5.08	.00											
12	3.22	2.43	4.03	3.50	3.60	3.74	1.66	4.06	2.74	3.94	5.21	.00										
13	3.96	3.43	4.39	2.58	4.76	4.55	5.01	4.14	3.66	1.41	5.31	4.50	.00									
14	2.11	4.32	2.74	3.23	4.82	3.47	4.91	4.34	3.82	3.61	4.32	4.34	4.39	.00								
15	2.59	2.50	5.16	3.19	4.26	4.07	2.93	3.85	4.11	4.26	4.74	2.33	5.10	4.24	.00							
16	4.03	4.84	5.26	4.97	5.82	5.84	5.04	2.20	3.63	4.53	3.43	4.62	4.41	5.17	5.18	.00						
17	4.40	3.62	6.36	4.89	5.63	6.10	4.58	5.43	4.90	5.48	4.75	3.50	5.61	5.56	3.40	5.56	.00					
18	1.88	2.90	2.72	2.65	4.34	2.85	2.95	3.24	2.43	3.07	3.95	2.45	3.78	2.30	3.00	3.97	4.43	.00				
19	2.41	4.63	3.18	3.46	5.13	2.58	4.52	4.11	4.11	4.13	4.52	4.41	5.01	1.88	4.03	5.23	6.09	2.47	.00			
20	3.17	3.00	3.73	1.82	4.39	2.91	3.54	4.09	2.95	2.05	5.35	3.43	2.23	3.74	3.78	4.82	4.87	2.92	3.90	.00		
21	3.45	2.32	5.09	3.88	3.64	4.63	2.68	3.98	3.74	4.36	4.88	1.38	4.94	4.93	2.10	4.57	3.10	3.19	4.97	4.15	.00	
22	2.51	2.42	4.11	2.58	3.77	4.03	4.00	3.24	3.21	2.56	3.44	3.00	2.74	3.51	3.35	3.46	3.63	2.55	3.97	2.62	3.01	.00

# Distances Between n = 22 Utilities

**Table 12.6** Distances Between 22 Utilities

Firm no.	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
1	.00																					
2	3.10	.00																				
3	3.68	4.92	.00																			
4	2.46	2.16	4.11	.00																		
5	4.12	3.85	4.47	4.13	.00																	
6	3.61	4.22	2.99	3.20	4.60	.00																
7	3.90	3.45	4.22	3.97	4.60	3.35	.00															
8	2.74	3.89	4.99	3.69	5.16	4.91	4.36	.00														
9	3.25	3.96	2.75	3.75	4.49	3.73	2.80	3.59	.00													
10	3.10	2.71	3.93	1.49	4.05	3.83	4.51	3.67	3.57	.00												
11	3.49	4.79	5.90	4.86	6.46	6.00	6.00	3.46	5.18	5.08	.00											
12	3.22	2.43	4.03	3.50	3.60	3.74	1.66	4.06	2.74	3.94	5.21	.00										
13	3.96	3.43	4.39	2.58	4.76	4.55	5.01	4.14	3.66	1.41	5.31	4.50	.00									
14	2.11	4.32	2.74	3.23	4.82	3.47	4.91	4.34	3.82	3.61	4.32	4.34	4.39	.00								
15	2.59	2.50	5.16	3.19	4.26	4.07	2.93	3.85	4.11	4.26	4.74	2.33	5.10	4.24	.00							
16	4.03	4.84	5.26	4.97	5.82	5.84	5.04	2.20	3.63	4.53	3.43	4.62	4.41	5.17	5.18	.00						
17	4.40	3.62	6.36	4.89	5.63	6.10	4.58	5.43	4.90	5.48	4.75	3.50	5.61	5.56	3.40	5.56	.00					
18	1.88	2.90	2.72	2.65	4.34	2.85	2.95	3.24	2.43	3.07	3.95	2.45	3.78	2.30	3.00	3.97	4.43	.00				
19	2.41	4.63	3.18	3.46	5.13	2.58	4.52	4.11	4.11	4.13	4.52	4.41	5.01	1.88	4.03	5.23	6.09	2.47	.00			
20	3.17	3.00	3.73	1.82	4.39	2.91	3.54	4.09	2.95	2.05	5.35	3.43	2.23	3.74	3.78	4.82	4.87	2.92	3.90	.00		
21	3.45	2.32	5.09	3.88	3.64	4.63	2.68	3.98	3.74	4.36	4.88	1.38	4.94	4.93	2.10	4.57	3.10	3.19	4.97	4.15	.00	
22	2.51	2.42	4.11	2.58	3.77	4.03	4.00	3.24	3.21	2.56	3.44	3.00	2.74	3.51	3.35	3.46	3.63	2.55	3.97	2.62	3.01	.00

# Average Linkage Dendrogram for Distances Between $n = 22$ Utilities

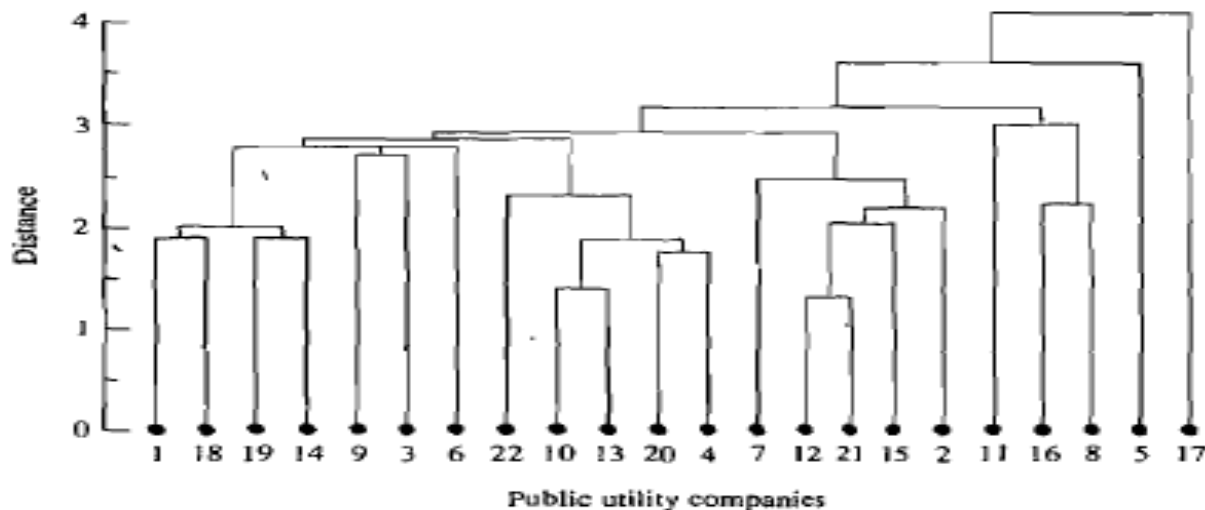


**Figure 12.10** Average linkage dendrogram for distances between 22 public utility companies.

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.



# Average Linkage Dendrogram for Distances Between $n = 22$ Utilities



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

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: Error Sum of Squares (ESS) Criterion

## 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  $ESS_k$  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  $ESS_k$  or  $ESS = ESS_1 + ESS_2 + \dots + ESS_K$ . 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,  $ESS_k = 0$ ,  $k = 1, 2, \dots, N$ , so  $ESS = 0$ . At the other extreme, when all the clusters are combined in a single group of  $N$  items, the value of ESS is given by

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

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.

# Dendrogram: Vertical Axis Gives Values ESS At Which Merges Occur

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

# Clustering Pure Malt Scotch Whiskies

**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.

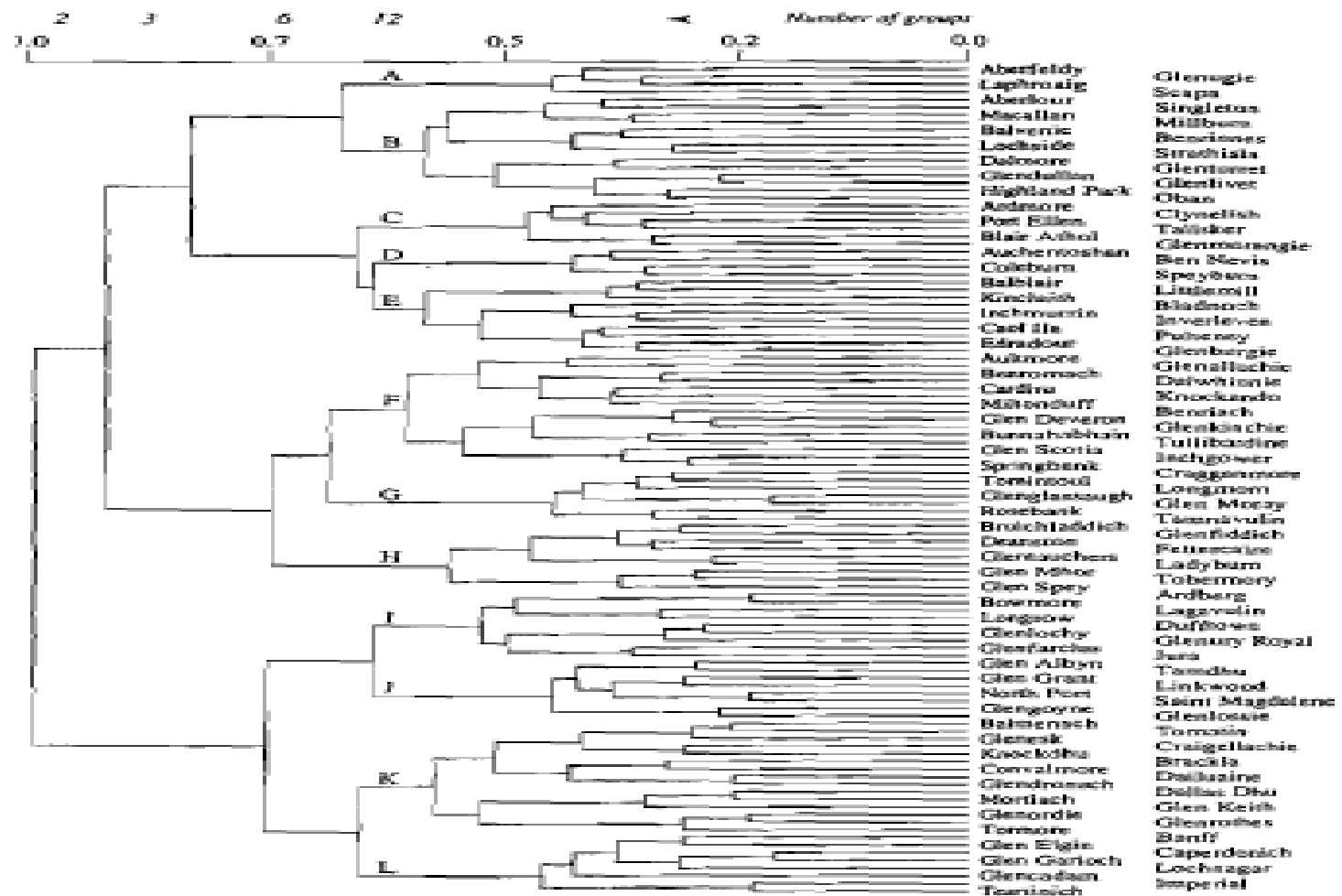
Weighted similarity coefficients  $\{s_{ik}\}$  were created from binary variables representing the presence or absence of characteristics. The resulting "distances," defined as  $\{d_{ik} = 1 - s_{ik}\}$ , were used with Ward's method to group the 109 pure (single-) malt Scotch whiskies. The resulting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

# Clustering $n = 109$ Pure Malt Scotch Whiskies

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

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

# Dendrogram: Clustering n = 109 Pure Malt Scotch Whiskies



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

# Note: Hierarchical Clustering No Provision for Reallocation of Objects Incorrectly Grouped at Early Stage

## 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.

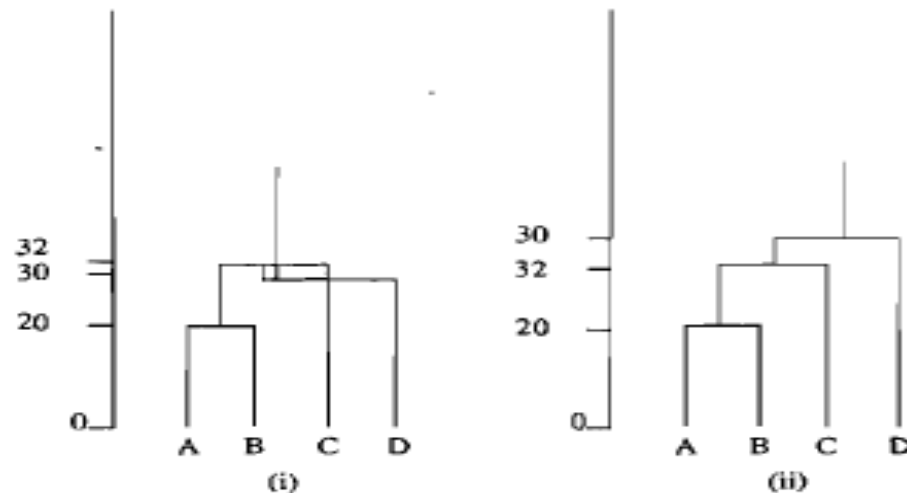
For a particular problem, it is a good idea to try several clustering methods and, within a given method, a couple different ways of assigning distances (similarities). If the outcomes from the several methods are (roughly) consistent with one another, perhaps a case for “natural” groupings can be advanced.

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 Ties in Similarity or Distance Matrix Can Produce Multiple Dendrograms or Multiple Clustering Solutions

Common values (ties) in the similarity or distance matrix can produce multiple solutions to a hierarchical clustering problem. That is, the dendrograms 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 (dendrograms) can be properly interpreted and different groupings (dendrograms) compared to assess their overlap. A further discussion of this issue appears in [27].

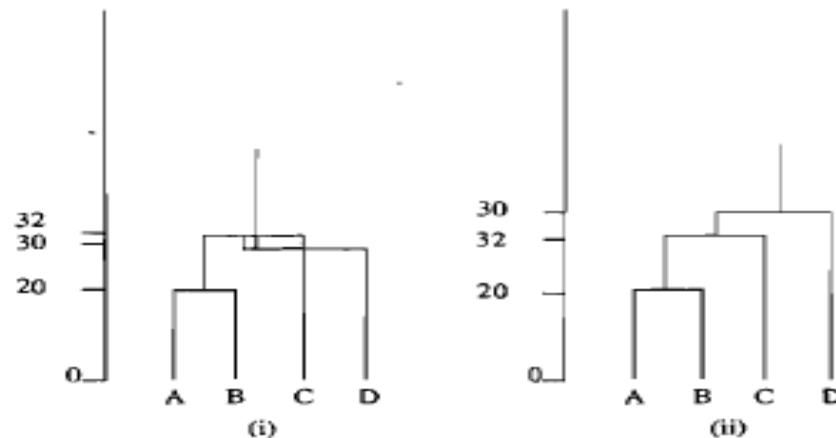
Some data sets and hierarchical clustering methods can produce *inversions*. (See [27].) An inversion occurs when an object joins an existing cluster at a smaller distance (greater similarity) than that of a previous consolidation. An inversion is represented two different ways in the following diagram:





# Inversions: Dendrogram With Crossover or Non-monotonic Scale

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:



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.

# Nonhierarchical Clustering Methods

## 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.

# Non-Hierarchical Methods: 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.

# Example: Clustering with the K-Means Method

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

# K-Means Based on Euclidean Distance

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:

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

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

# Clustering Centroids Calculated from Squared Distances

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

$$d^2(A, (AB)) = (5 - 2)^2 + (3 - 2)^2 = 10$$

if  $A$  is not moved

$$d^2(A, (CD)) = (5 + 1)^2 + (3 + 2)^2 = 61$$

$$d^2(A, (B)) = (5 + 1)^2 + (3 - 1)^2 = 40$$

if  $A$  is moved to the  $(CD)$  group

$$d^2(A, (ACD)) = (5 - 1)^2 + (3 + .33)^2 = 27.09$$

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

$$d^2(B, (AB)) = (-1 - 2)^2 + (1 - 2)^2 = 10$$

if  $B$  is not moved

$$d^2(B, (CD)) = (-1 + 1)^2 + (1 + 2)^2 = 9$$

$$d^2(B, (A)) = (-1 - 5)^2 + (1 - 3)^2 = 40$$

if  $B$  is moved to the  $(CD)$  group

$$d^2(B, (BCD)) = (-1 + 1)^2 + (1 + 1)^2 = 4$$

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.

# Contingency Table to Check for Squared Distances from Group Centroids

We check  $C$  for reassignment.

$$d^2(C, (A)) = (1 - 5)^2 + (-2 - 3)^2 = 41 \quad \text{if } C \text{ is not moved}$$

$$d^2(C, (BCD)) = (1 + 1)^2 + (-2 + 1)^2 = 5$$

$$d^2(C, (AC)) = (1 - 3)^2 + (-2 - 5)^2 = 10.25 \quad \text{if } C \text{ is moved to the } (A) \text{ group}$$

$$d^2(C, (BD)) = (1 + 2)^2 + (-2 + 5)^2 = 11.25$$

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 reassignments take place and the final  $K = 2$  clusters are  $(A)$  and  $(BCD)$ .

For the final clusters, we have

Cluster	Squared distances to group centroids			
	Item			
	$A$	$B$	$C$	$D$
$A$	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  $(BCD)$ :  $4 + 5 + 5 = 14$

# Results: Seven Possibilities for K=2 Clusters

For the final clusters, we have

Cluster	Squared distances to group centroids			
	Item			
	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
<i>A</i>	0	40	41	89
( <i>BCD</i> )	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

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

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*)



# Check Stability: Generate Table Cluster Centroids (Means) and Within Cluster Variances to Delineate Group Differences

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: K-means Clustering of n =22 Public Utilities

**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

	1	2	3	4
1	0			
2	3.08	0		
3	3.29	3.56	0	
4	3.05	2.84	3.18	0

# Example: K-means Clustering of $n = 22$ Public Utilities

$K = 5$

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).

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

# Using Univariate F-ratios To Test Between Cluster Variability and Within Cluster Variability

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 (FUELC) 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. ■

# Outliers Will Influence K-Means: Thus Mitigate Leverage Points Via Transformation

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.
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.

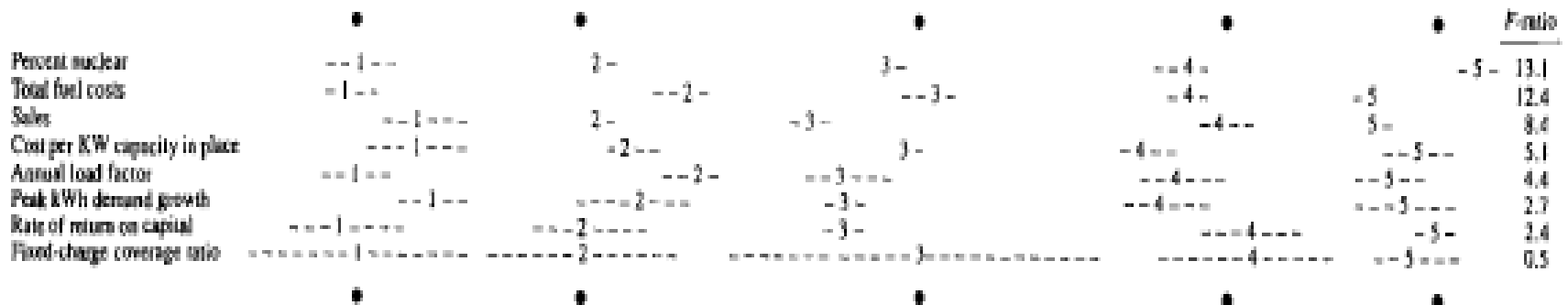
In cases in which a single run of the algorithm requires the user to specify  $K$ , it is always a good idea to rerun the algorithm for several choices.

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

# Cluster Profiles (K=5)

## Ordered by F-Ratio

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.

# Clustering Based on Statistical Methods

## 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*

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

where each  $p_k \geq 0$  and  $\sum_{k=1}^K p_k = 1$ . This distribution  $f_{\text{Mix}}(\mathbf{x})$  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.

# Clusters Generated in Ellipsoidal Shape With Heaviest Concentration Near Center

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

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

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



# Inferences Based on Likelihood for N Objects or Subjects and Fitted Number of Clusters Where Proportions, Mean Vectors and Covariance Matrices Unknown

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

$$\begin{aligned} L(p_1, \dots, p_K, \mu_1, \Sigma_1, \dots, \mu_K, \Sigma_K) &= \prod_{j=1}^N f_{Mix}(\mathbf{x}_j | \mu_1, \Sigma_1, \dots, \mu_K, \Sigma_K) \\ &= \prod_{j=1}^N \left( \sum_{k=1}^K p_k \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x}_j - \mu_k)' \Sigma_k^{-1} (\mathbf{x}_j - \mu_k) \right) \right) \end{aligned} \quad (12-18)$$

where the proportions  $p_1, \dots, p_K$ , the mean vectors  $\mu_1, \dots, \mu_K$ , and the covariance matrices  $\Sigma_1, \dots, \Sigma_K$  are unknown. The measurements for different objects are treated as independent and identically distributed observations from the mixture distribution.

# Statistical Models Not the Same as Single Linkage, Complete Linkage or Average Linkage Because Use Maximum Likelihood

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  $\Sigma_k$  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  $\hat{p}_1, \dots, \hat{p}_K, \hat{\mu}_1, \hat{\Sigma}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_K$  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

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

# Statistical Method of Clustering: Use Maximum Likelihood and Akaike Information Criterion (AIC) to Calculate Model difference

A good approach to selecting a model is to first obtain the maximum likelihood estimates  $\hat{p}_1, \dots, \hat{p}_K, \hat{\mu}_1, \hat{\Sigma}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_K$  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

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

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

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

# Bayesian Information Criterion (BIC)

## Similar But Uses Logarithm of No. Parameters in Penalty Function

covariances. For the Akaike information criterion (AIC), the penalty is  $2N \times$  (number of parameters) so

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

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

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

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
$\Sigma_k = \eta \mathbf{I}$	$K(p + 1)$	$\ln L_{\max} - 2 \ln(N) K(p + 1)$
$\Sigma_k = \eta_k \mathbf{I}$	$K(p + 2) - 1$	$\ln L_{\max} - 2 \ln(N) (K(p + 2) - 1)$
$\Sigma_k = \eta_k \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$	$K(p + 2) + p - 1$	$\ln L_{\max} - 2 \ln(N) (K(p + 2) + p - 1)$

# Mixture Model Can Be Applied With Certain Software, i.e. Hierarchical and Statistical EM Algorithm and BIC Criterion

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  $1, 2, \dots, K$ . 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: Model Based Clustering of Types of Irises Iris Data

**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

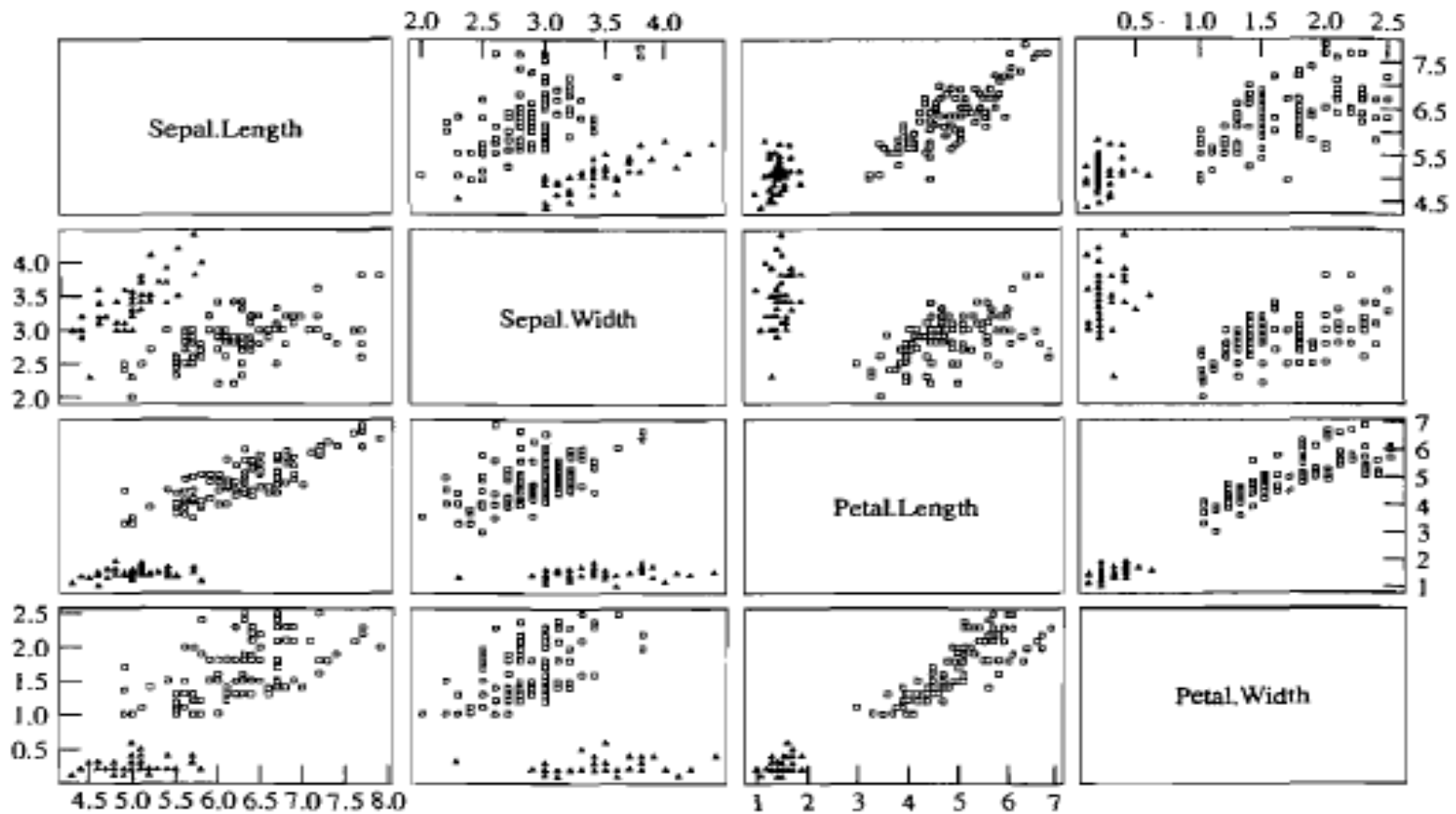
$$\mu_1 = \begin{bmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} 5.90 \\ 2.75 \\ 4.40 \\ 1.43 \end{bmatrix}, \quad \mu_3 = \begin{bmatrix} 6.85 \\ 3.07 \\ 5.73 \\ 2.07 \end{bmatrix},$$

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  $\text{BIC} = -705.1$ .

# Multiple Scatter Plots for $K = 3$ Clusters for Iris Data with Four Variables



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

# Using BIC Criterion with $K=2$ Groups of Clusters and Covariance Structures, Best Choice is Two Groups With 33% in First Cluster and 67% in Second Cluster

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

$$\mu_1 = \begin{bmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{bmatrix}, \quad \mu_2 = \begin{bmatrix} 6.26 \\ 2.87 \\ 4.91 \\ 1.68 \end{bmatrix}$$

and the two estimated covariance matrices are

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

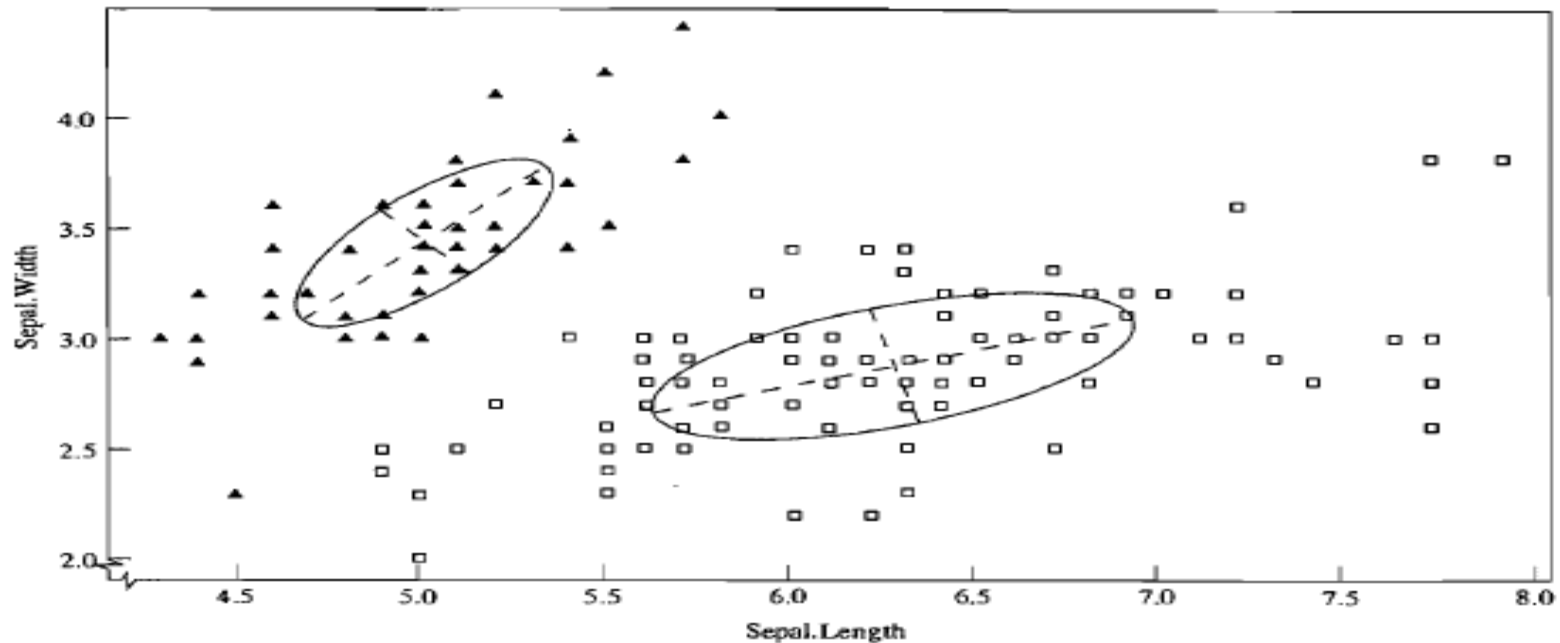
Essentially, two species of Iris have been put in the same cluster as the projected view of the scatter plot of the sepal measurements in Figure 12.14 shows. ■



# Multidimensional Scaling and Clustering

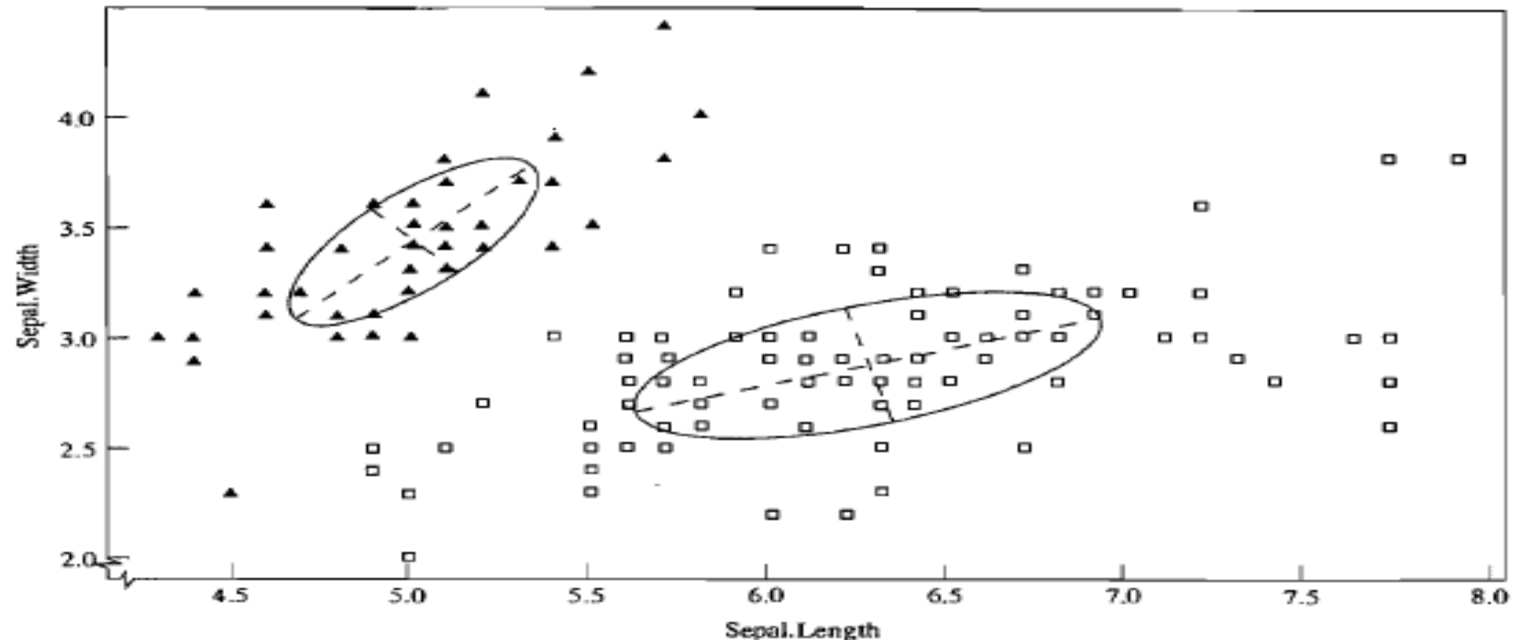
## 12.6 Multidimensional Scaling

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** Scatter plot of sepal measurements for best model.

# Scatter Plot of Iris Sepal Measurements for Best Model



**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.

# Clustering Method: Matrix Multidimensional Scaling or Principal Component Analysis

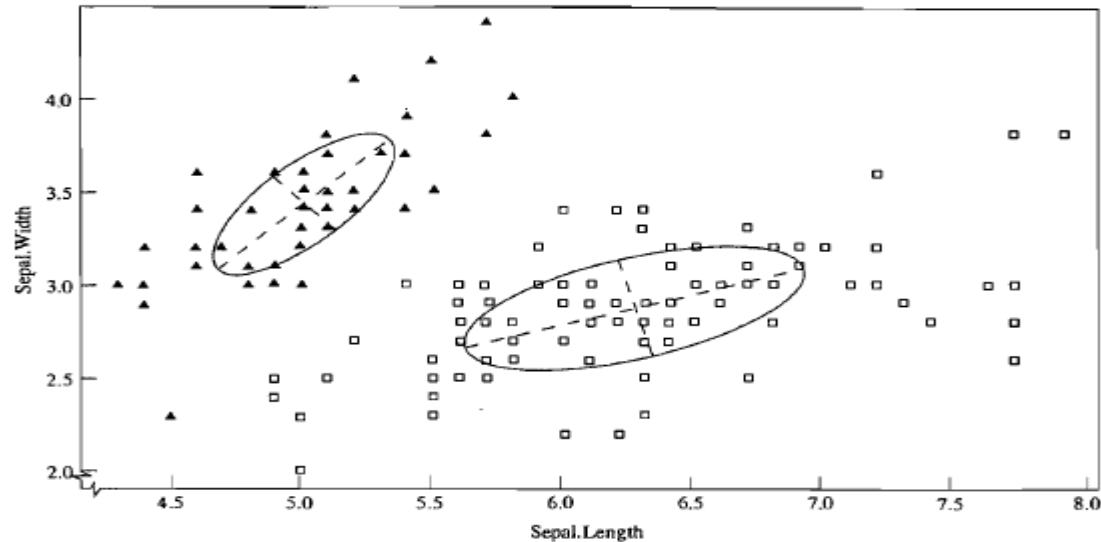


Figure 12.14 Scatter plot of sepal measurements for best model.

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*.

# Multidimensional Scaling: Basic Algorithm

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

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

# Goodness of Fit: Ascending Ordering of Monotonic Relationship Between Initial Similarities and Distances

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

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

The  $\hat{d}_{ik}^{(q)}$ 's in the stress formula are numbers known to satisfy (12-22); that is, they are monotonically related to the similarities. The  $\hat{d}_{ik}^{(q)}$ '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  $d_{ik}^{(q)}$ '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:

<i>Stress</i>	<i>Goodness of fit</i>	(12-24)
20%	Poor	
10%	Fair	
5%	Good	
2.5%	Excellent	
0%	Perfect	

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

# Diminish Stress: i.e. Another Measure of Discrepancy

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 SStress, replaces the  $d_{ik}$ 's and  $\hat{d}_{ik}$ 's in (12-23) by their squares and is given by

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

The value of SStress 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.

# Look for Elbow in Stress-Dimensionality Plot: Distances Ordered from Largest to Smallest

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 SStress (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. (12-26)

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.

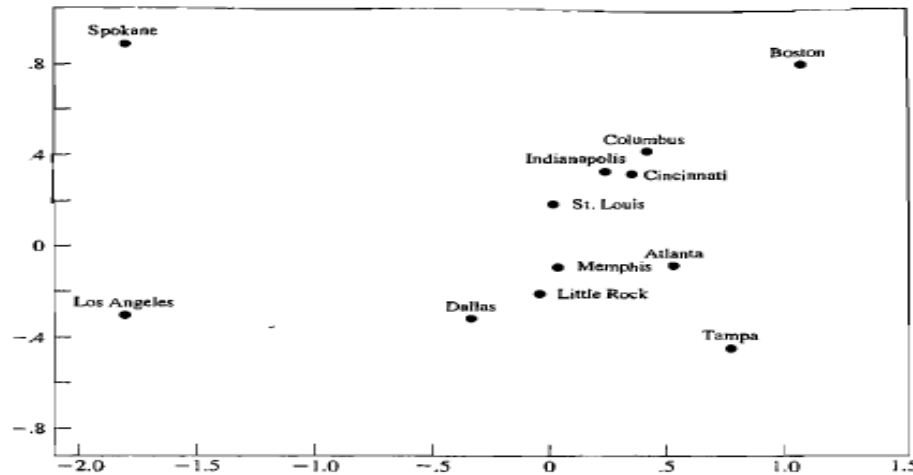
# Multidimensional Scaling of $n = 12$ US Cities: Distance Between Pairs

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

	Atlanta (1)	Boston (2)	Cincinnati (3)	Columbus (4)	Dallas (5)	Indianapolis (6)	Little Rock (7)	Los Angeles (8)	Memphis (9)	St. Louis (10)	Spokane (11)	Tampa (12)
(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



# Multidimensional Scaling $q = 2$



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

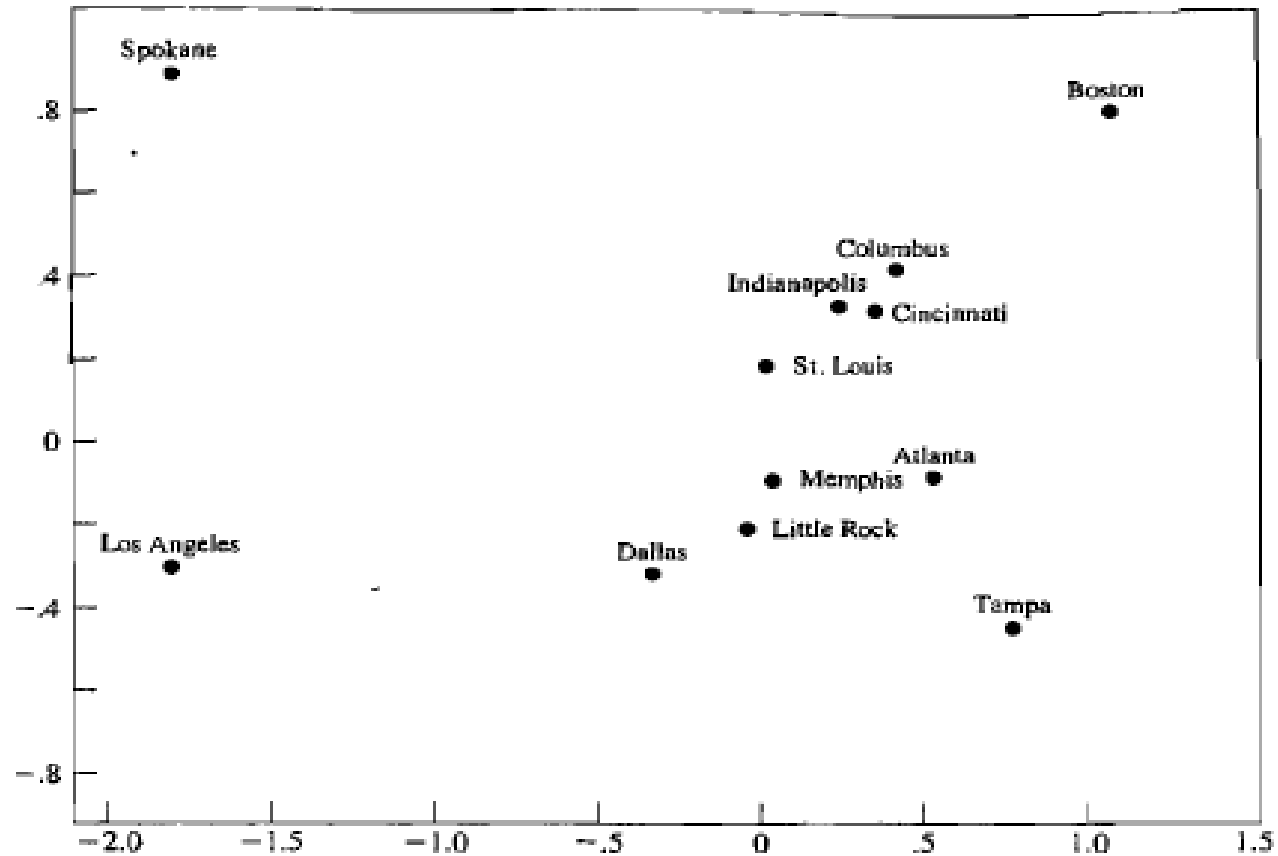
Since the cities naturally lie in a two-dimensional space (a nearly level part of the curved surface of the earth), it is not surprising that multidimensional scaling with  $q = 2$  will locate these items about as they occur on a map. Note that if the distances in the table are ordered from largest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

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. ■

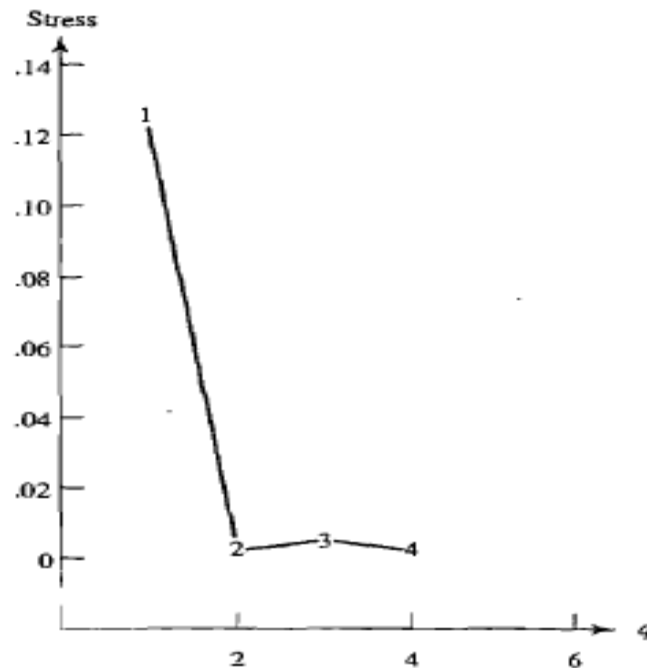
# Geometrical Representation of Cities Produced by Multidimensional Scaling



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

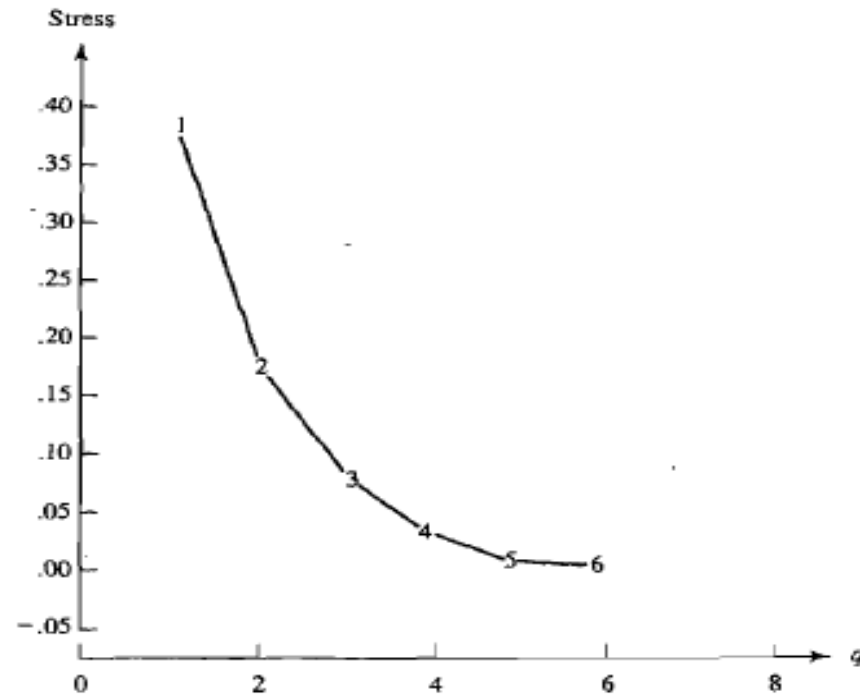
# Example: Multidimensional Scaling of Public Utilities

**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.

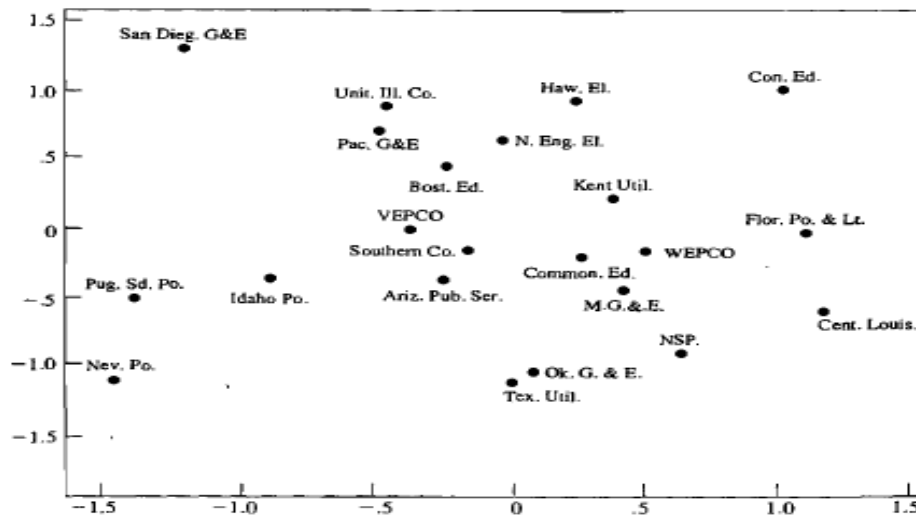
# Stress Function For Distances Between Utilities: No Sharp Elbow to Determine No. of Clusters



**Figure 12.17** Stress function for distances between utilities.

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.

# Multidimensional Scaling: Geometrical Representation of Public Utilities



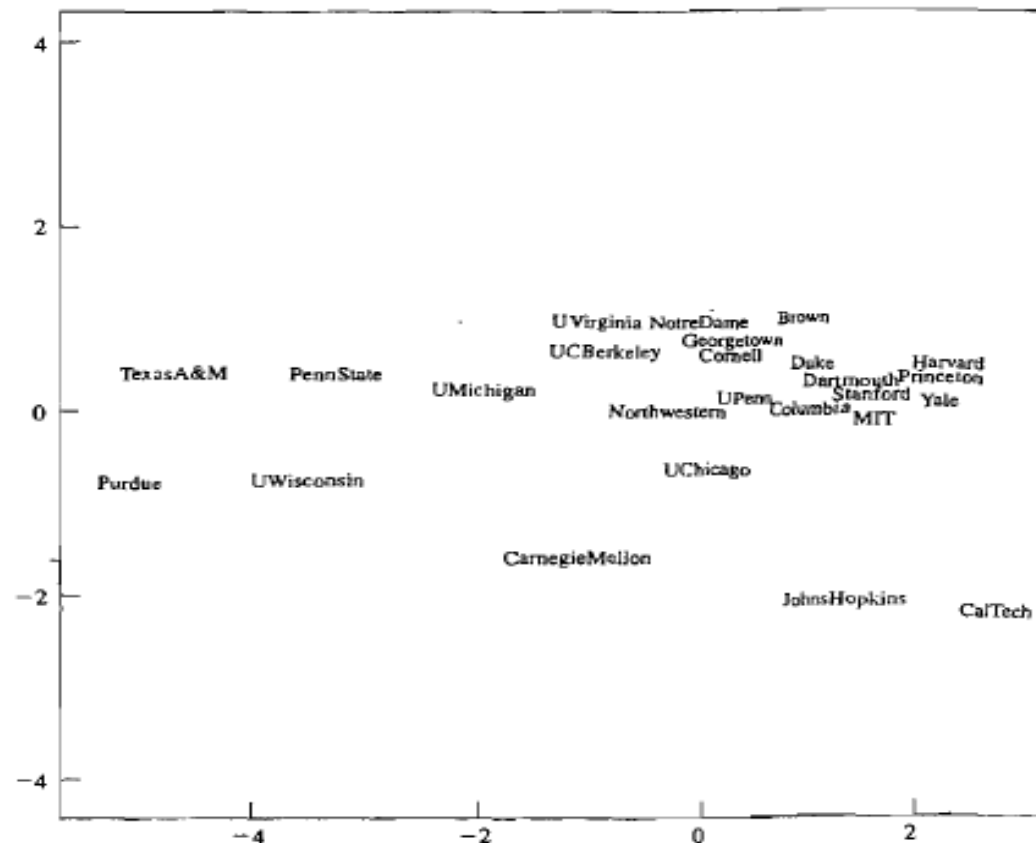
**Figure 12.18** A geometrical representation of utilities produced by multidimensional scaling.

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. ■

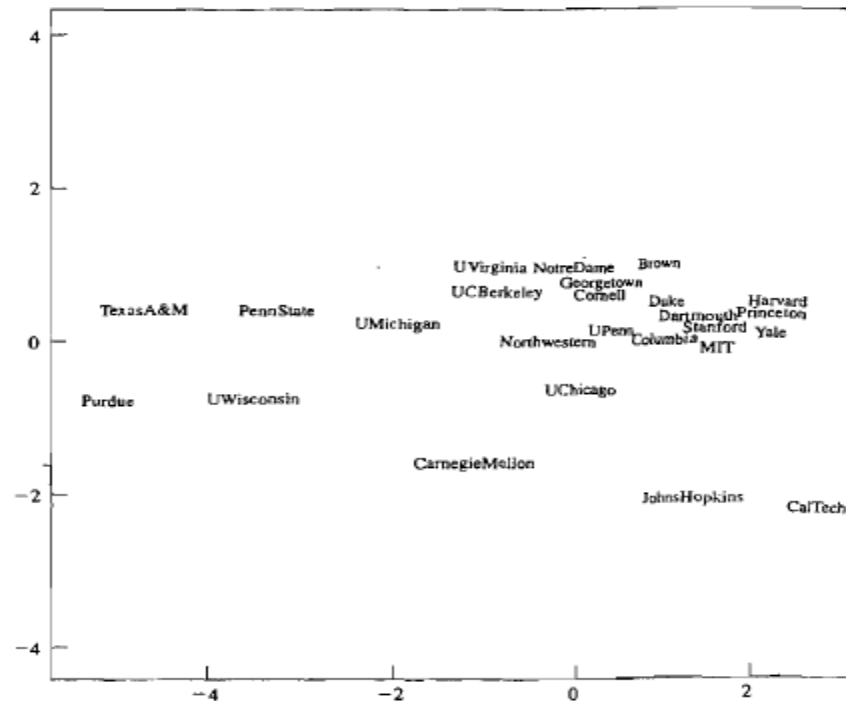
# Example: Multidimensional Scaling of Universities

**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.

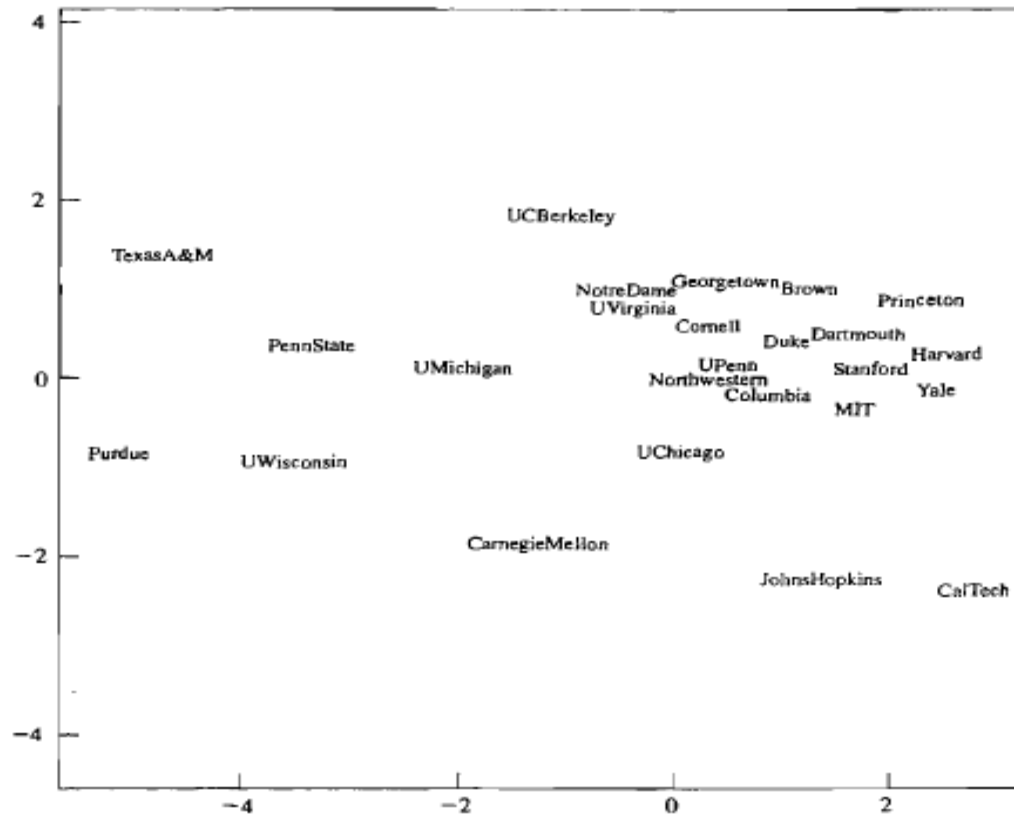
# Multidimensional Scaling of Universities: Private vs. Public Clusters



**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. ■

# Non Metric Multidimensional Scaling Similar to Non Metric With Stress Value Apx. 10%



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

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. ■



# Multidimensional Scaling: Euclidean Distances $p$ and $q$ Dimensions Compared Directly

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].)

4. Plot minimum stress ( $q$ ) versus  $q$  and choose the best number of dimensions,  $q^*$ , from an examination of this plot. (12-26)

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

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

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.)

# Biplots for Viewing Sampling Units and Variables: Graphical Representation of $n \times p$ Data Matrix

## 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.

# Biplot of n=22 Public Utilities Based on Variables Annual Load Factor, etc.

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-

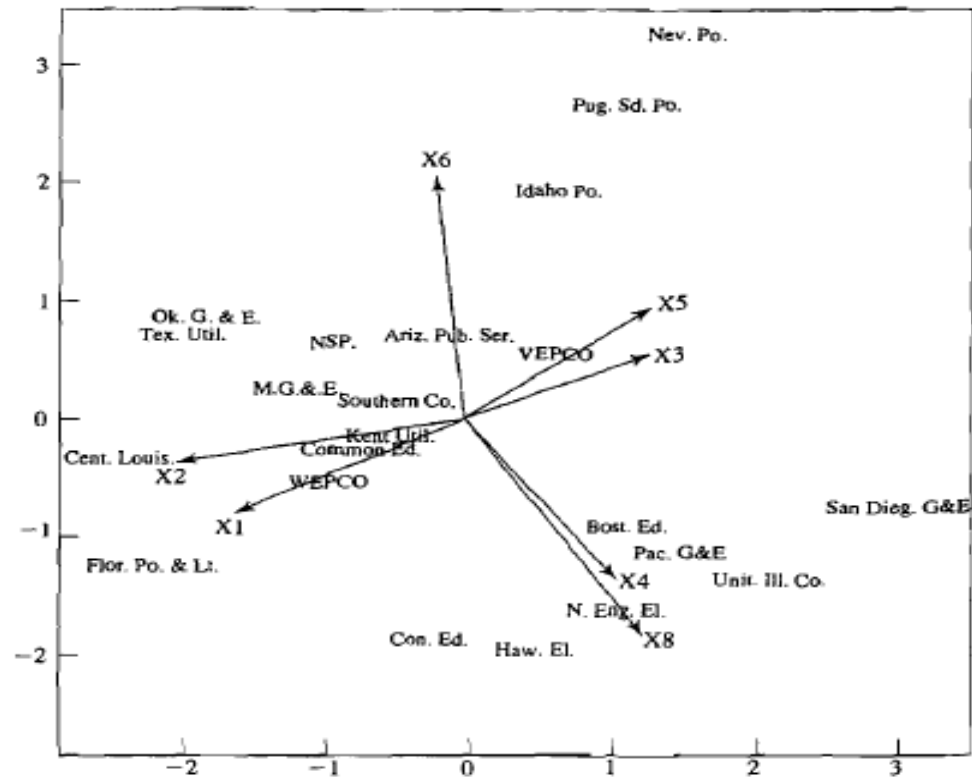


Figure 12.23 A biplot of the data on public utilities.

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

# Constructing Biplots: Proceeds from Principal Components Eigenvalues and Eigenvectors

## 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  $\hat{\mathbf{E}}'(\mathbf{x}_0 - \bar{\mathbf{x}})$ .

# Biplot Construction Uses Orthogonality Like Principal Components and Factoring

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

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

where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  and  $\mathbf{V}$  is an orthogonal matrix whose columns are the eigenvectors of  $\mathbf{X}_c' \mathbf{X}_c = (n-1)\hat{\mathbf{S}}$ . 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

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

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,  $\mathbf{U} \mathbf{\Lambda}$  contains all of the values of the principal components, while  $\mathbf{V} = \hat{\mathbf{E}}$  contains the coefficients that define the principal components.

# In Biplot, Each Row of Data Matrix or Item Represented by Point Located by Pair of Values of 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.

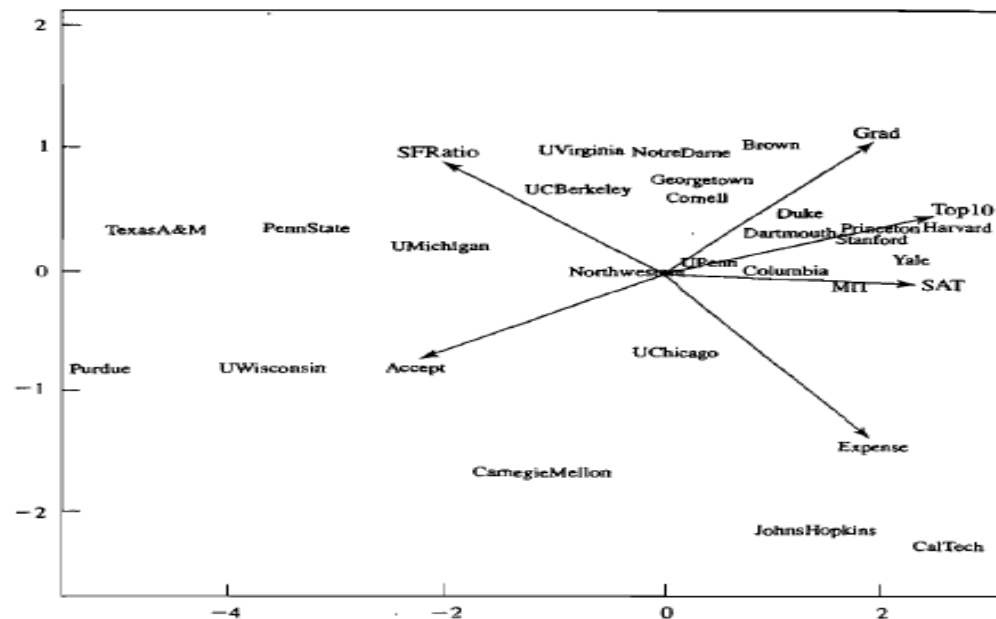
In the biplot, each *row* of the data matrix, or item, is represented by the point located by the pair of values of the principal components. The *ith column* of the data matrix, or variable, is represented as an *arrow* from the origin to the point with coordinates  $(e_{1i}, e_{2i})$ , the entries in the *ith 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].)

# Example: Biplot of Universities and Their Characteristics, Attributes or Predictors

**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  $z_i$ . The biplot is given in Figure 12.24 on page 7



**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. ■

# Biplot: Clustering/Division/Separation of Public vs. Private Universities

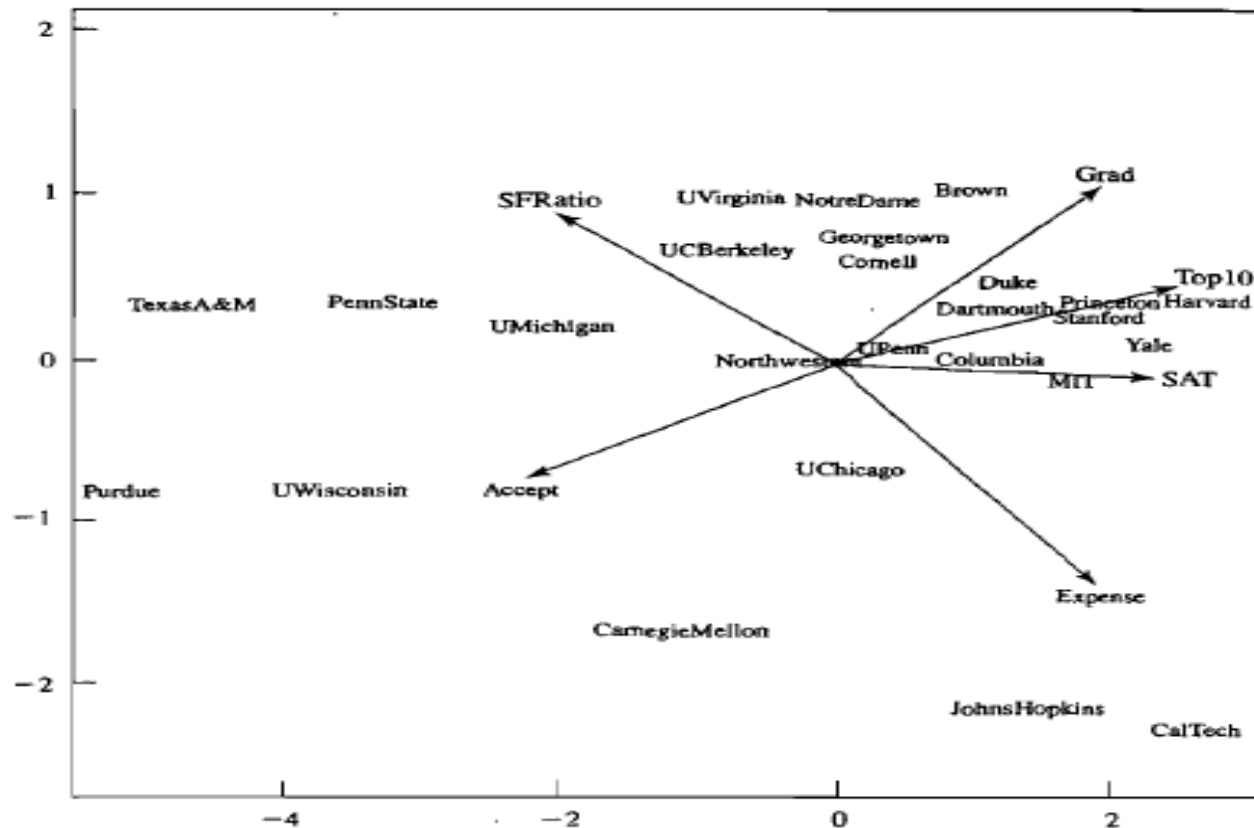
Notice how Cal Tech and Johns Hopkins are off by themselves; the variable Expense is mostly responsible for this positioning. The large state universities in our sample are to the left in the biplot, and most of the private schools are on the right.

<b>Table 12.9 Data on Universities</b>						
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.



# Biplot of Universities: Large Values for SAT, Top 10 and Grade Rate Variables Associated with Private School Group



**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. ■

# Biplot (Gower and Hand) Version: Principal Component Axes Repressed; Constructed Axis for Each Variables and Scale is Attached

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  $\hat{\mathbf{E}} = [\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2]$ .

# Biplot (Gower and Hand): Projection Vector of Sample Mean Origin of the Biplot

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

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

# Biplot (Gower and Hand): Provides a Scale for Mean-Centered Variable and Thus Interpolate Position of Eigenvectors in Biplot

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

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

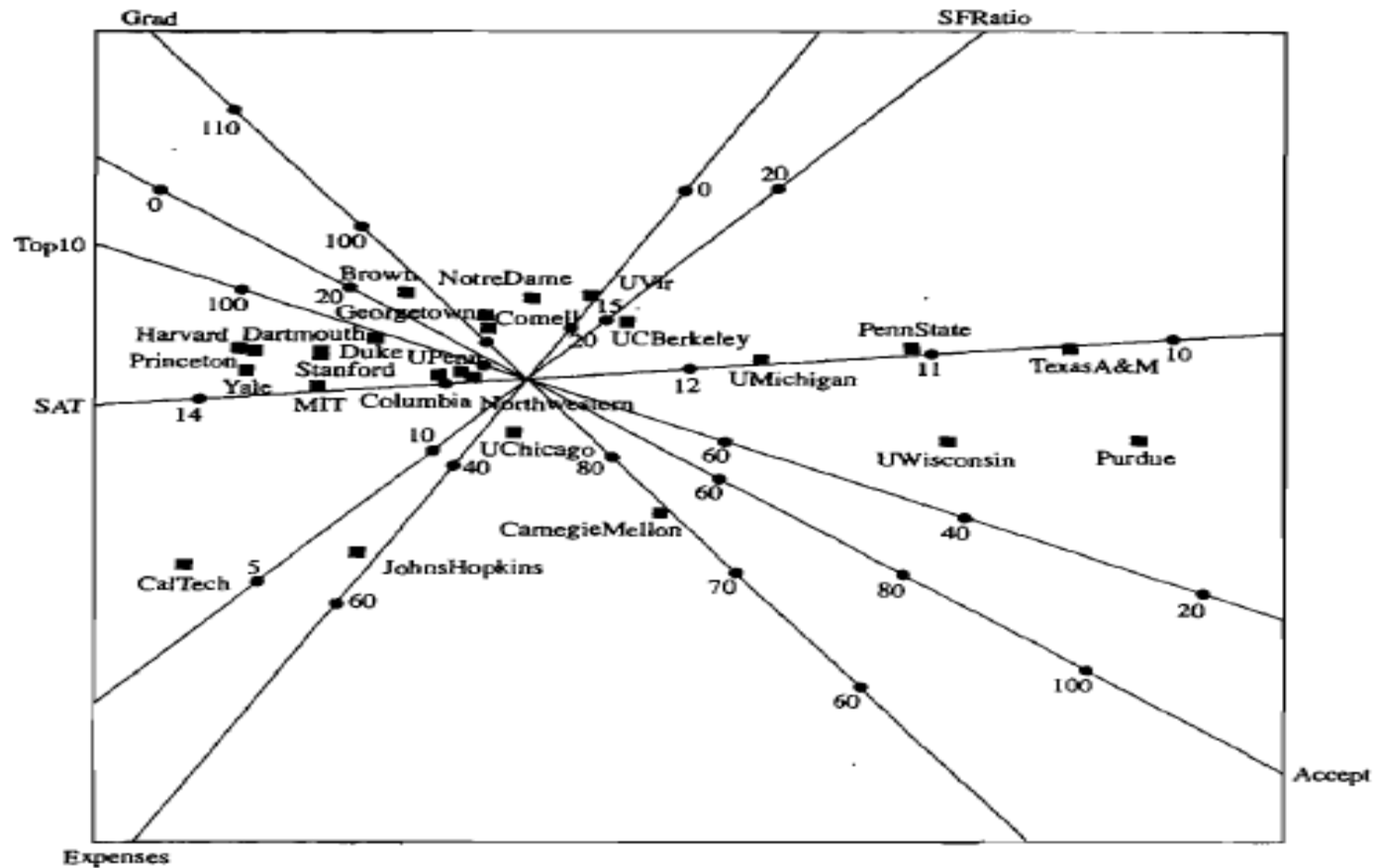
Starting from the origin, the points in the direction  $w[e_{1i}, e_{2i}]'$  are plotted for  $w = 0, \pm 1, \pm 2, \dots$ . This provides a scale for the mean centered variable  $x_i - \bar{x}_i$ . It defines the distance in the biplot for a change of one unit in  $x_i$ . But, the origin for the  $i$ -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.

# Biplot (Gower and Hand) for University Data: Reverses the Direction 1<sup>st</sup> Principal Component

**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. ■

# Biplot (Gower and Hand): Scaled Axes More Informative Than Previous 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.