

Cluster analysis

Based on H.C. Romesburg: *Cluster analysis for researchers*,
Lifetime Learning Publications, Belmont, CA, 1984
P.H.A. Sneath and R.R. Sokal: *Numericxal Taxonomy*, Freeman,
San Fransisco, CA, 1973

Jens C. Frisvad
BioCentrum-DTU

Biological data analysis and chemometrics

Two primary methods

- Cluster analysis (no projection)
 - Hierarchical clustering
 - Divisive clustering
 - Fuzzy clustering
- Ordination (projection)
 - Principal component analysis
 - Correspondence analysis
 - Multidimensional scaling

Advantages of cluster analysis

- Good for a quick overview of data
- Good if there are many groups in data
- Good if unusual similarity measures are needed
- Can be added on ordination plots (often as a minimum spanning tree, however)
- Good for the nearest neighbours, ordination better for the deeper relationships

Different clustering methods

- NCLAS: Agglomerative clustering by distance optimization
- HMCL: Agglomerative clustering by homogeneity optimization
- INFCL: Agglomerative clustering by information theory criteria
- MINGFC: Agglomerative clustering by global optimization
- ASSIN: Divisive monothetic clustering
- PARREL: Partitioning by global optimization
- FCM: Fuzzy c-means clustering
- MINSPAN: Minimum spanning tree
- REBLOCK: Block clustering (k-means clustering)

SAHN clustering

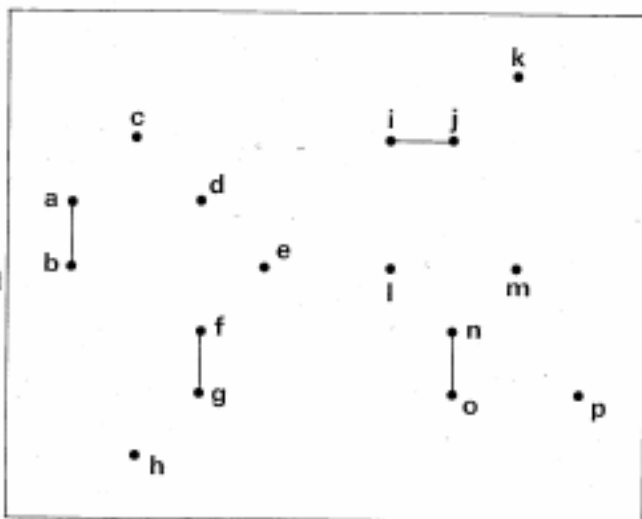
- Sequential agglomerative hierarchic nonoverlapping clustering

Single linkage

- Nearest neighbor, minimum method
- Close to minimum spanning tree
- Contracting space
- Chaining possible
- $\alpha_J = 0.5, \alpha_K = 0.5, \beta = 0, \gamma = -0.5$
- $U_{J,K} = \min U_{jk}$

$$U_{(J,K)L} = \alpha_J U_{J,L} + \alpha_K U_{K,L} + \beta U_{J,K} + \gamma |U_{J,L} - U_{K,L}|$$

Step 1
Link, $\Delta = 1$



Step 2
Link, $\Delta \leq \sqrt{2}$

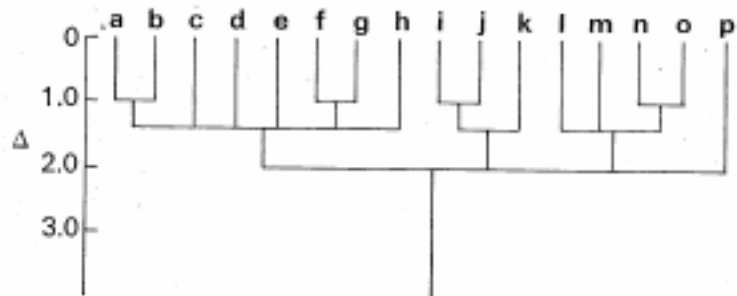
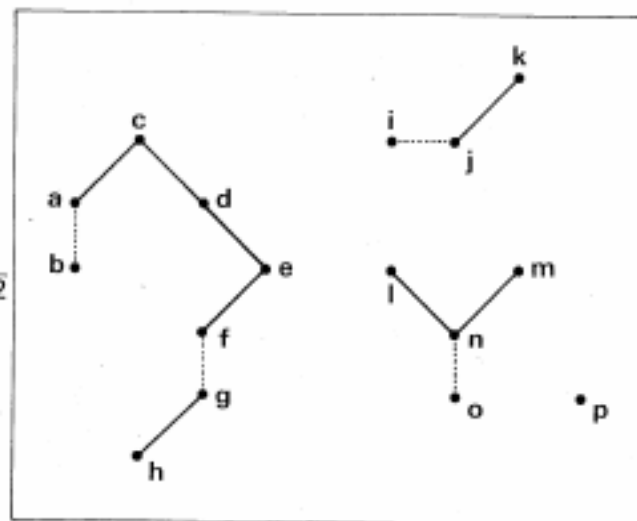
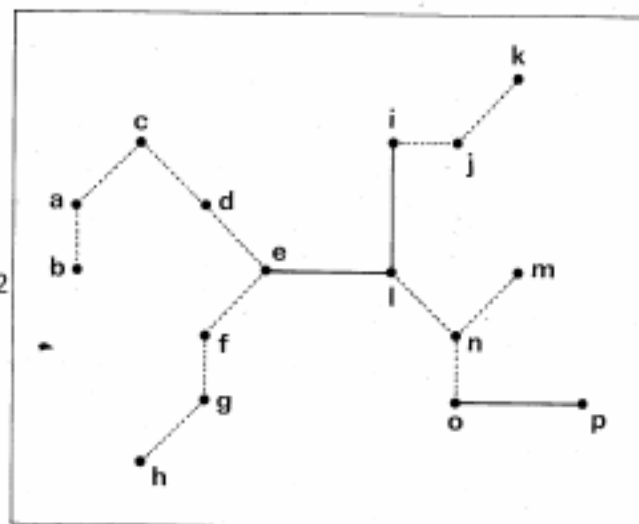


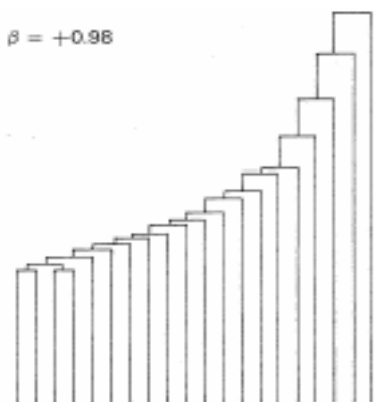
FIGURE 5-3

Single linkage clustering of the data in Table 5-1. For explanation, see text.

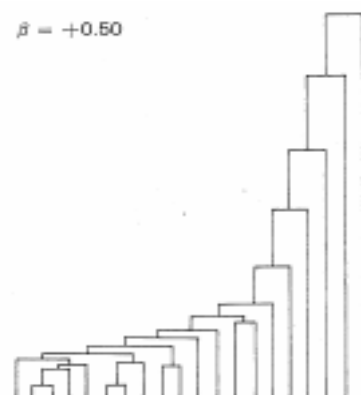
Step 3
Link, $\Delta \leq 2$



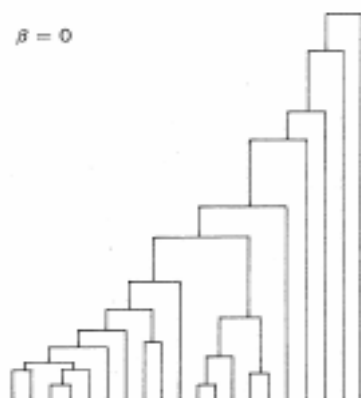
$\beta = +0.98$



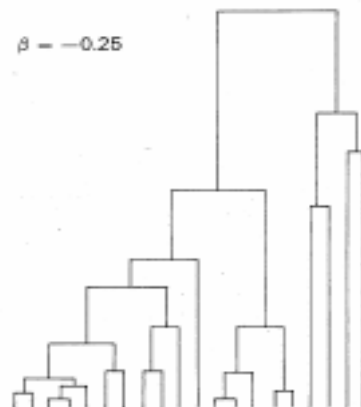
$\beta = +0.50$



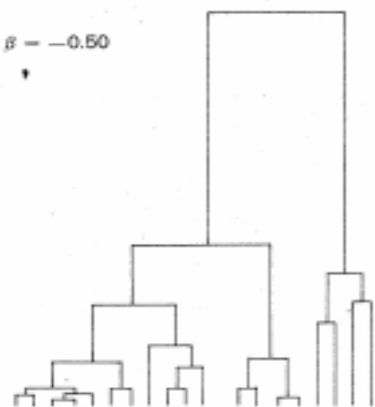
$\beta = 0$



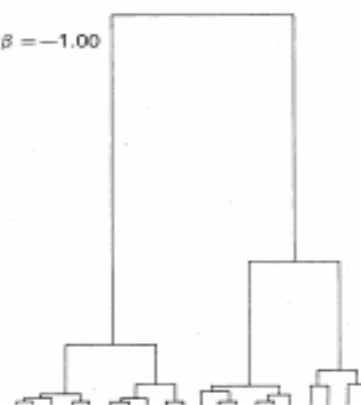
$\beta = -0.25$



$\beta = -0.50$



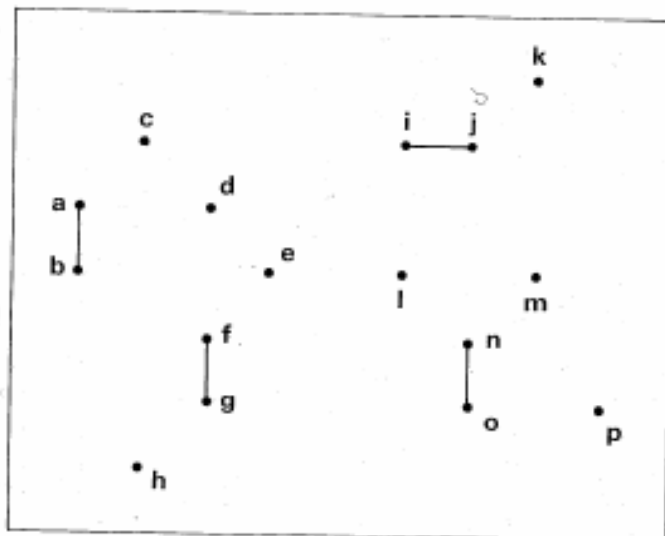
$\beta = -1.00$



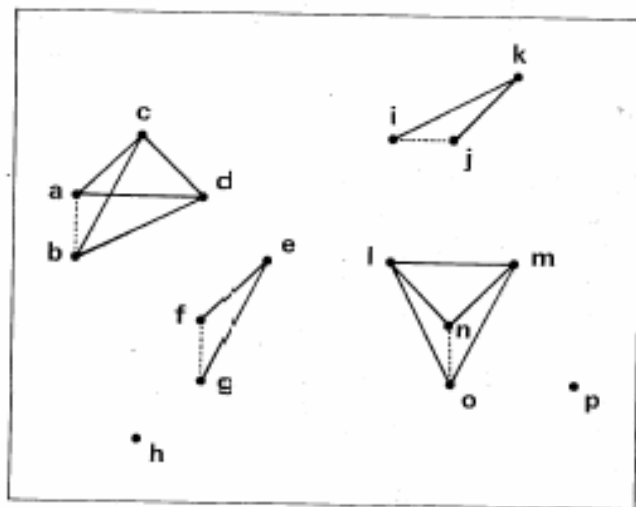
Complete linkage

- Furthest neighbor, maximum method
- Dilating space
- $\alpha_J = 0.5, \alpha_K = 0.5, \beta = 0, \gamma = 0.5$
- $U_{J,K} = \max U_{jk}$

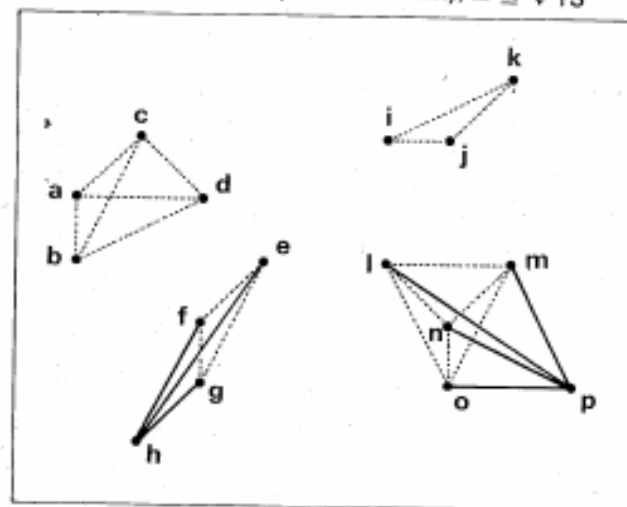
Step 1
 Level of cohesion
 (maximum link), $\Delta = 1$



Step 2
 Level of cohesion
 (maximum link), $\Delta \leq \sqrt{5}$



Step 3
 Level of cohesion
 (maximum link), $\Delta \leq \sqrt{13}$

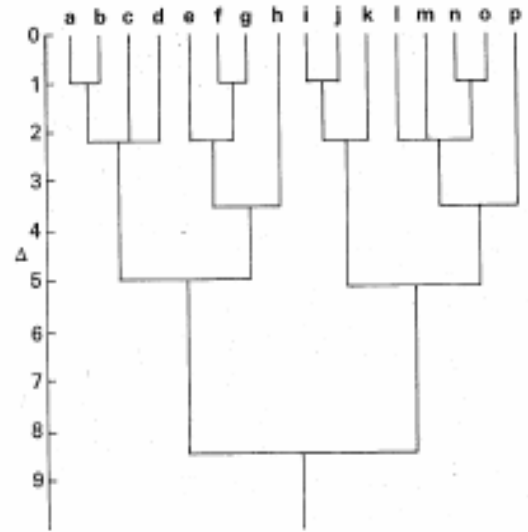
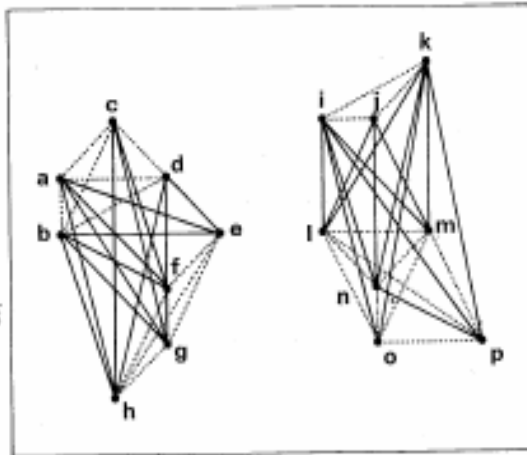


Step 4

Levels of cohesion
(maximum links)

Left cluster, $\Delta \leq 5$

Right cluster, $\Delta \leq \sqrt{26}$

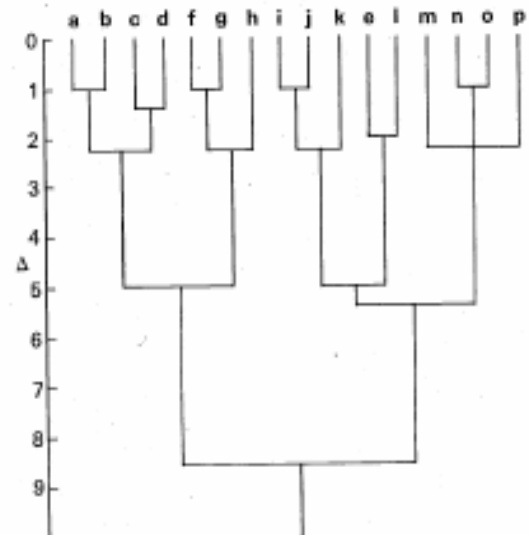
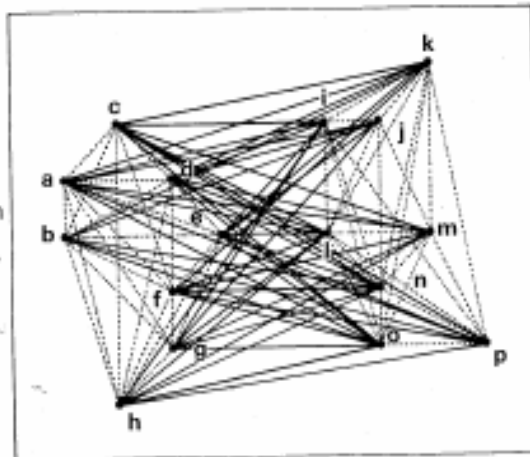


Phenogram A

Step 5

Level of cohesion
(maximum link),

$\Delta \leq \sqrt{73}$



Phenogram B

Average linkage

- Arithmetic average
 - Unweighted: UPGMA (group average)
 - Weighted: WPGMA
- Centroid
 - Unweighted centroid (Centroid)
 - Weighted centroid (Median)

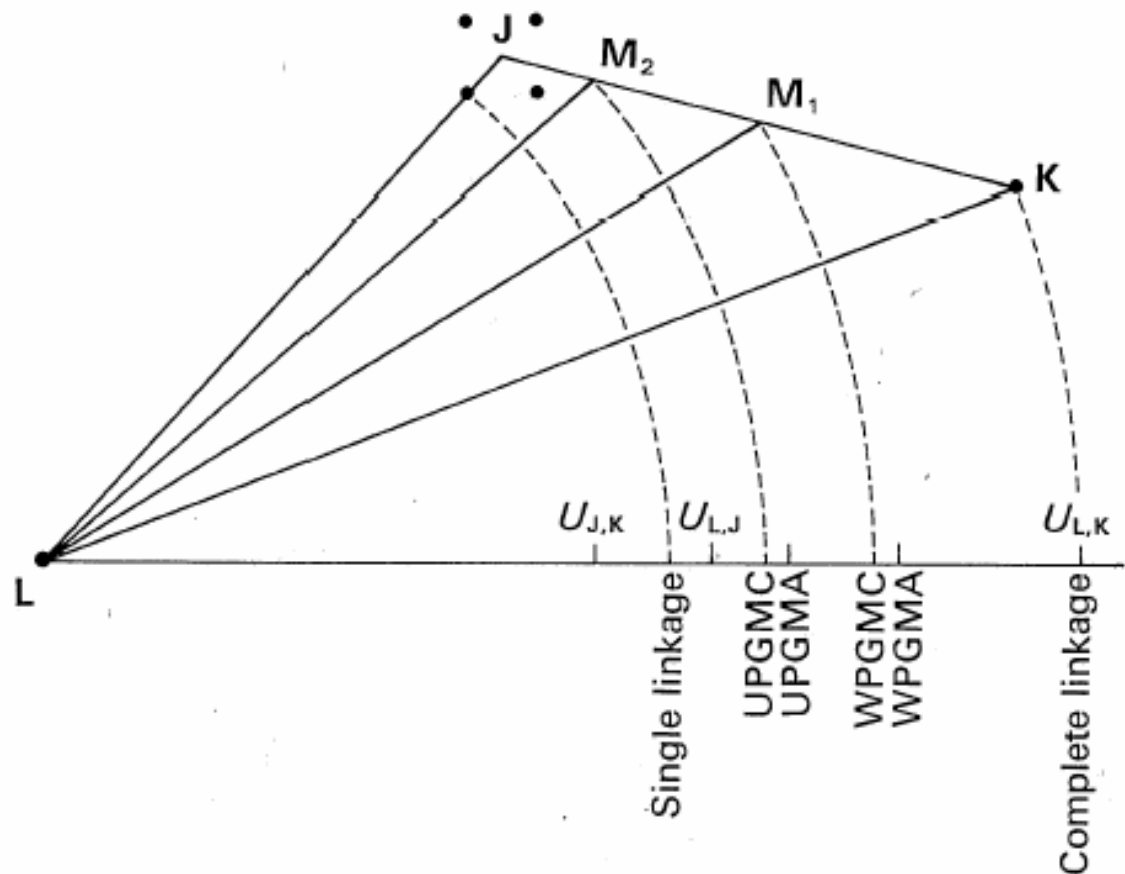


FIGURE 5-11

The effects of several clustering methods (for explanation see text) on the criterion for admitting L (containing one OTU) to the cluster formed of four OTU's in J plus one in K. OTU's are indicated by solid circles. Abscissa is distance Δ .

From Sneath and Sokal, 1973, Numerical taxonomy

Ordinary clustering

- Obtain the data matrix
- Transform or standardize the data matrix
- Select the best resemblance or distance measure
- Compute the resemblance matrix
- Execute the clustering method (often UPGMA = average linkage)
- Rearrange the data and resemblance matrices
- Compute the cophenetic correlation coefficient

Binary similarity coefficients

(between two objects i and j)

	j	1	0
i			
1		a	b
0		c	d

Matches and mismatches

- $m = a + b$ (number of matches)
- $u = c + d$ (number of mismatches)
- $n = m + u = a + b + c + d$ (total sample size)

- Similarity (often 0 to 1)
- Dissimilarity (distance) (often 0 to 1)
- Correlation (-1 to 1)

Simple matching coefficient

- $SM = (a + d) / (a + b + c + d) = m / n$
- Euclidean distance for binary data:
- $D = 1 - SM = (b + c) / (a + b + c + d) = u / n$

Avoiding zero zero comparisons

- Jaccard = $J = a / (a + b + c)$
- Sørensen or Dice: $DICE = 2a / (2a + b + c)$

Correlation coefficients

$$\text{Yule: } (ad - bc) / (ad + bc)$$

$$PHI = (ad - bc) / \sqrt{(a + b)(c + d)(a + c)(b + d)}$$

Other binary coefficients

- Hamann = $H = (a + d - b - c) / (a + b + c + d)$
- Rogers and Tanimoto = $RT = (a + d) / (a + 2b + 2c + d)$
- Russel and Rao = $RR = a / (a + b + c + d)$
- Kulzynski 1 = $K1 = a / (b + c)$
- UN1 = $(2a + 2d) / (2a + b + c + 2d)$
- UN2 = $a / (a + 2b + 2c)$
- UN3 = $(a + d) / (b + c)$

Distances for quantitative (interval) data

Euclidean and taxonomic distance

$$EUCLID = E_{ij} = \sqrt{\sum_k (x_{ki} + x_{kj})^2}$$

$$DIST = d_{ij} = \sqrt{\frac{1}{n} \sum_k (x_{ki} + x_{kj})^2}$$

Bray-Curtis and Canberra distance

$$BRAYCURT = d_{ij} = \sum_k |x_{ki} - x_{kj}| / \sum_k (x_{ki} + x_{kj})$$

$$CANBERRA = \frac{1}{n} \sum_k |x_{ki} - x_{kj}| / \sum_k (x_{ki} + x_{kj})$$

Average Manhattan distance (city block)

$$MANHAT = M_{ij} = \frac{1}{n} \sum_k |x_{ki} - x_{kj}|$$

Chi-squared distance

$$CHISQ = d_{ij} = \sqrt{\sum_k \frac{\left(\frac{x_{ki}}{x_{\cdot i}} - \frac{x_{kj}}{x_{\cdot j}} \right)^2}{x_k}}$$

Cosine coefficient

$$COSINE = c_{ij} = \sum_k x_{ki} x_{kj} / \sqrt{\sum_k x_{ki}^2 \sum_k x_{kj}^2}$$

Step 1. Obtain the data matrix

Object

1 2 3 4 5

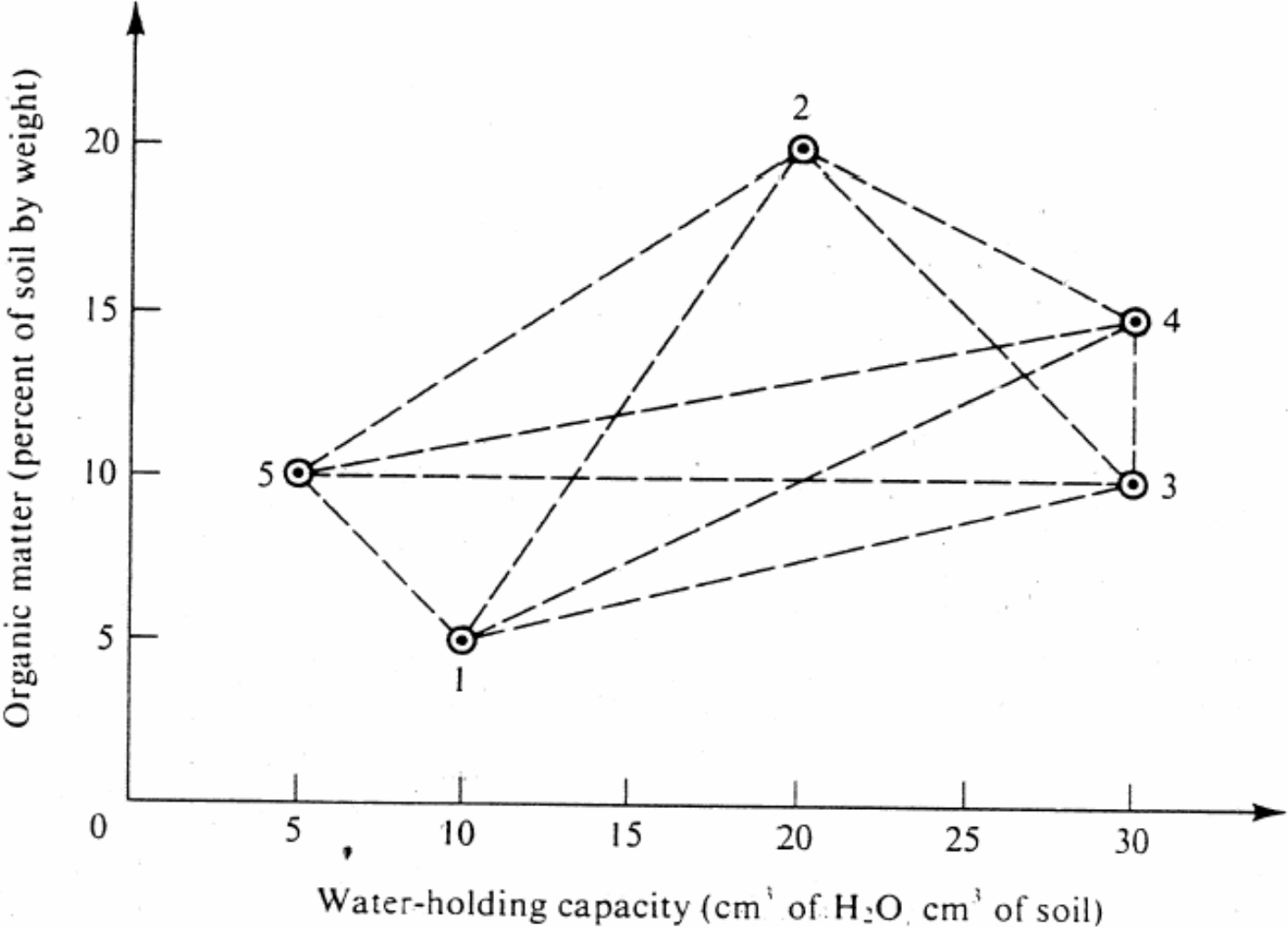
1	10	20	30	30	5
2	5	20	10	15	10

Feature

Objects and features

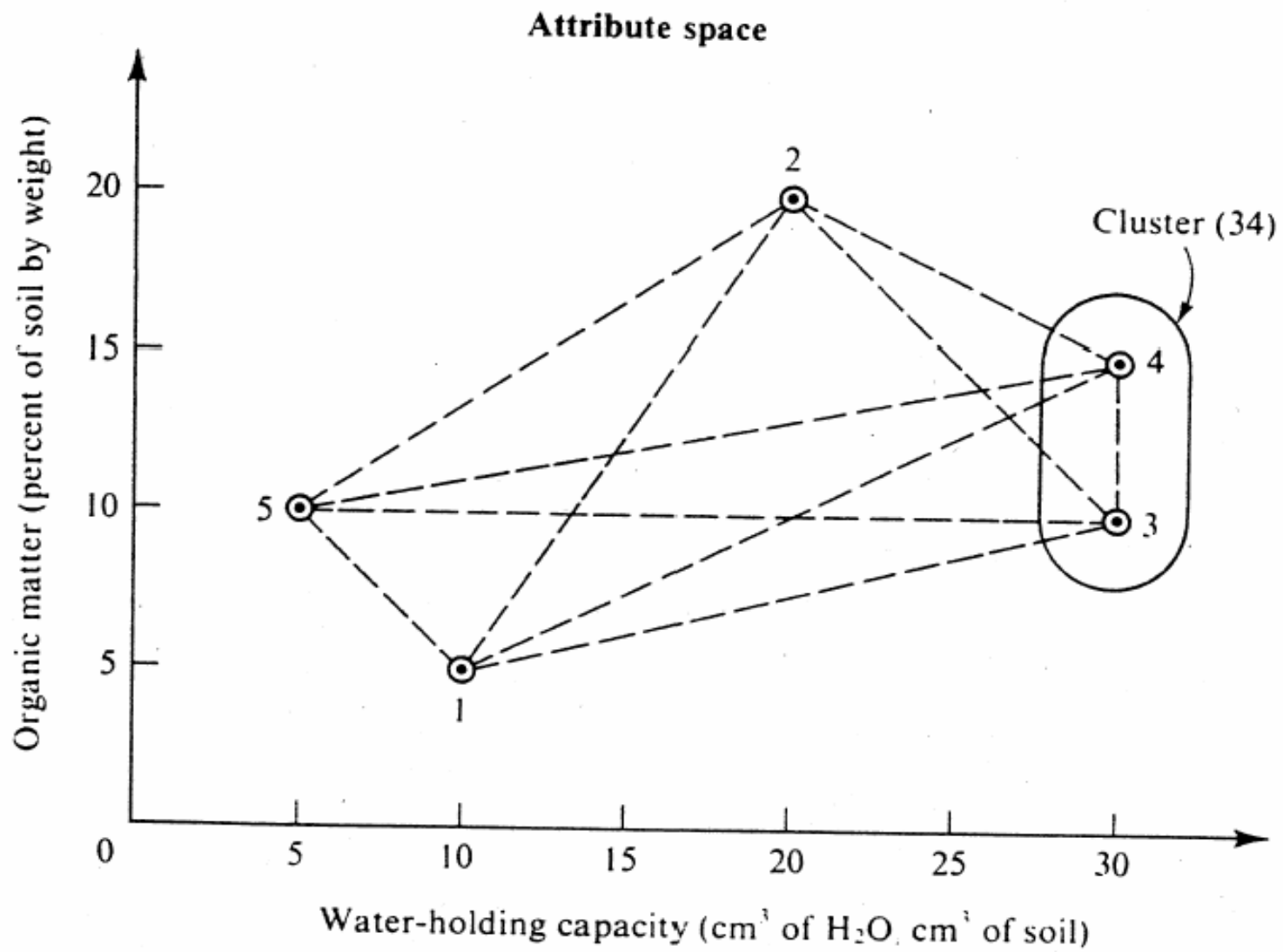
- The five objects are plots of farm land
- The features are
 - 1. Water-holding capacity (%)
 - 2. Weight % soil organic matter
- Objective: find the two most similar plots

Attribute space



Resemblance matrix

	1	2	3	4	5
1	-	-	-	-	-
2	18.0	-	-	-	-
3	20.6	14.1	-	-	-
4	22.4	11.2	5.00	-	-
5	7.07	18.0	25.0	25.5	-



Revised resemblance matrix

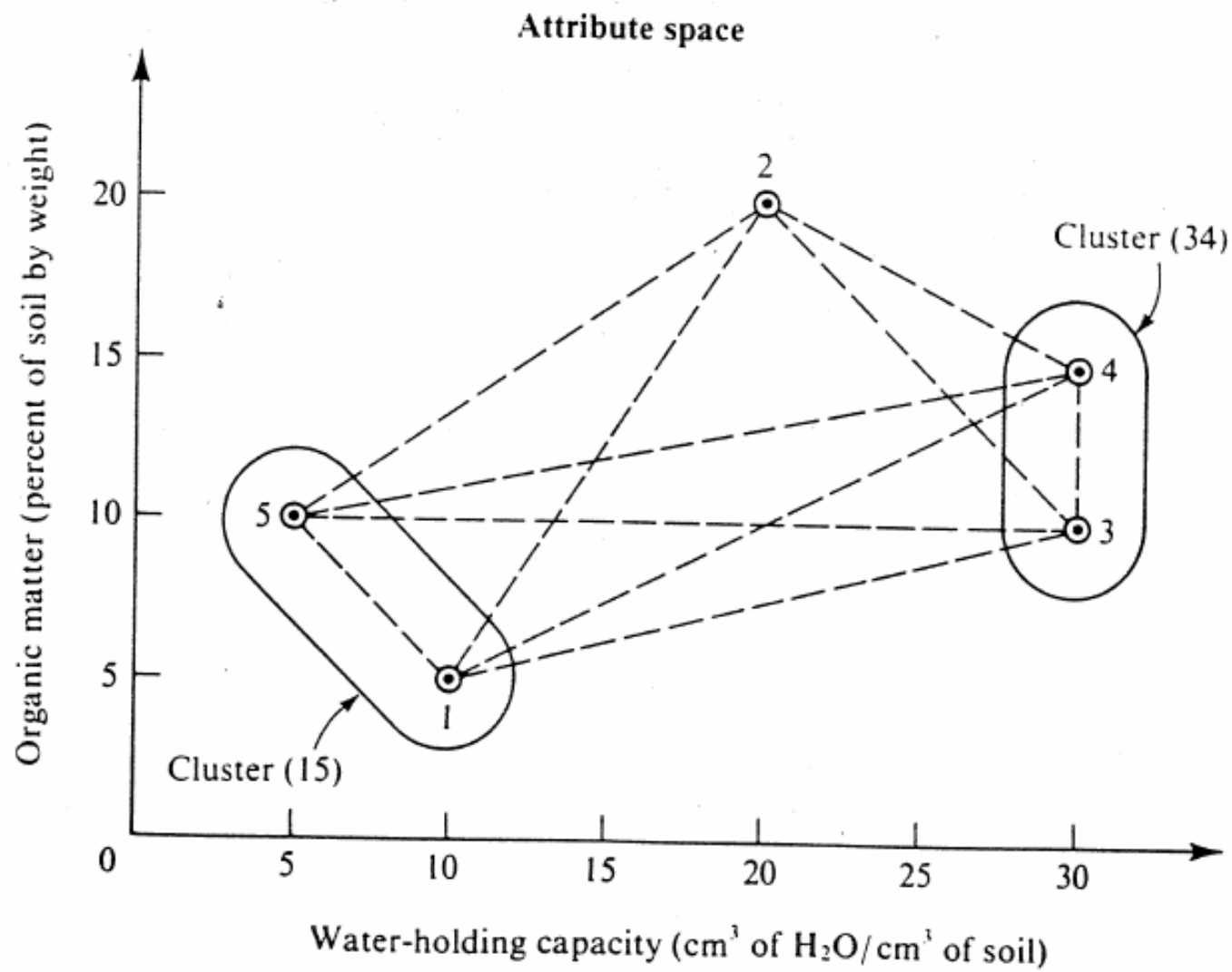
1

2

5

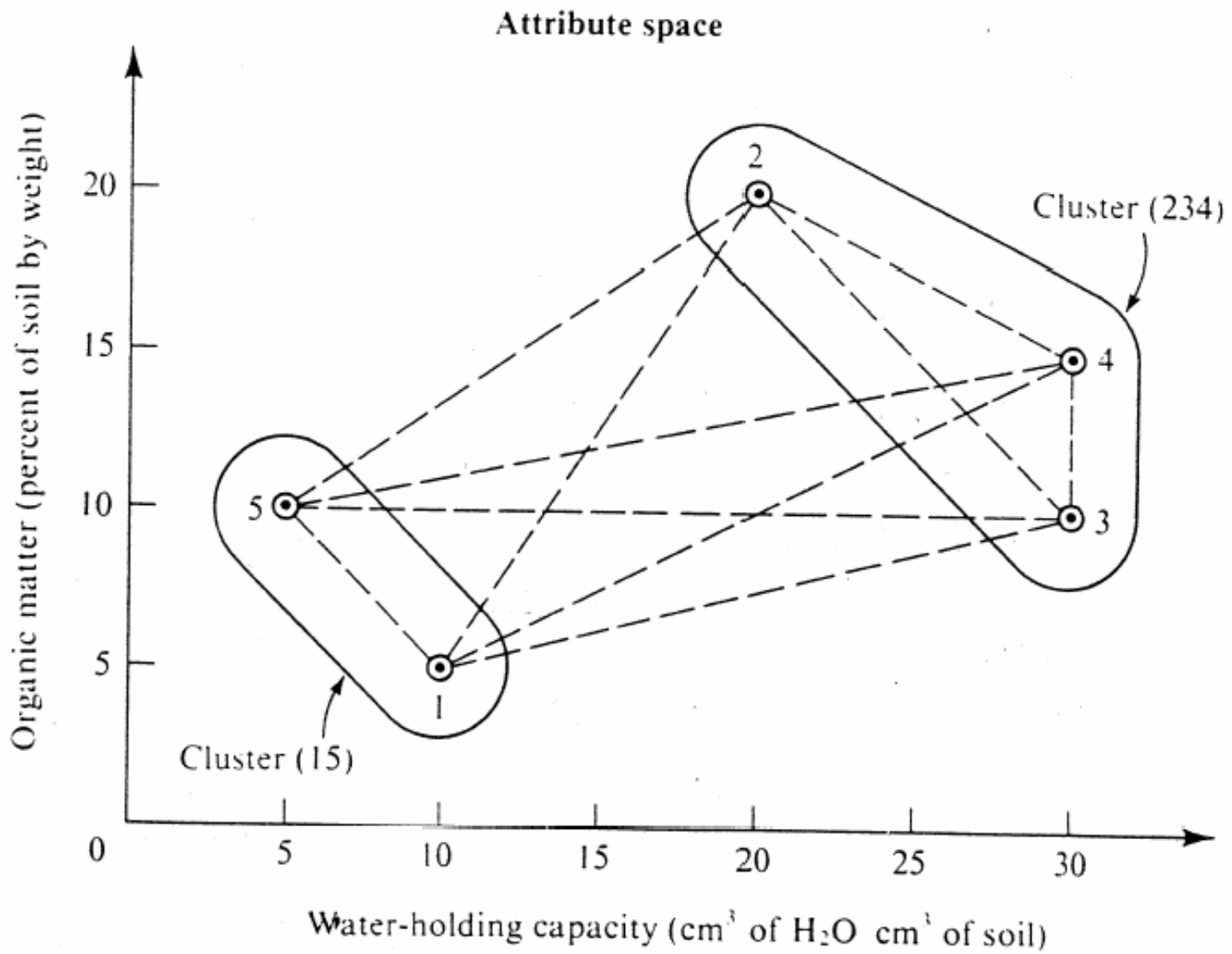
(34)

1	-	-	-	-
2	18.0	-	-	-
5	7.07	18.0	-	-
(34)	21.5	12.7	25.3	-



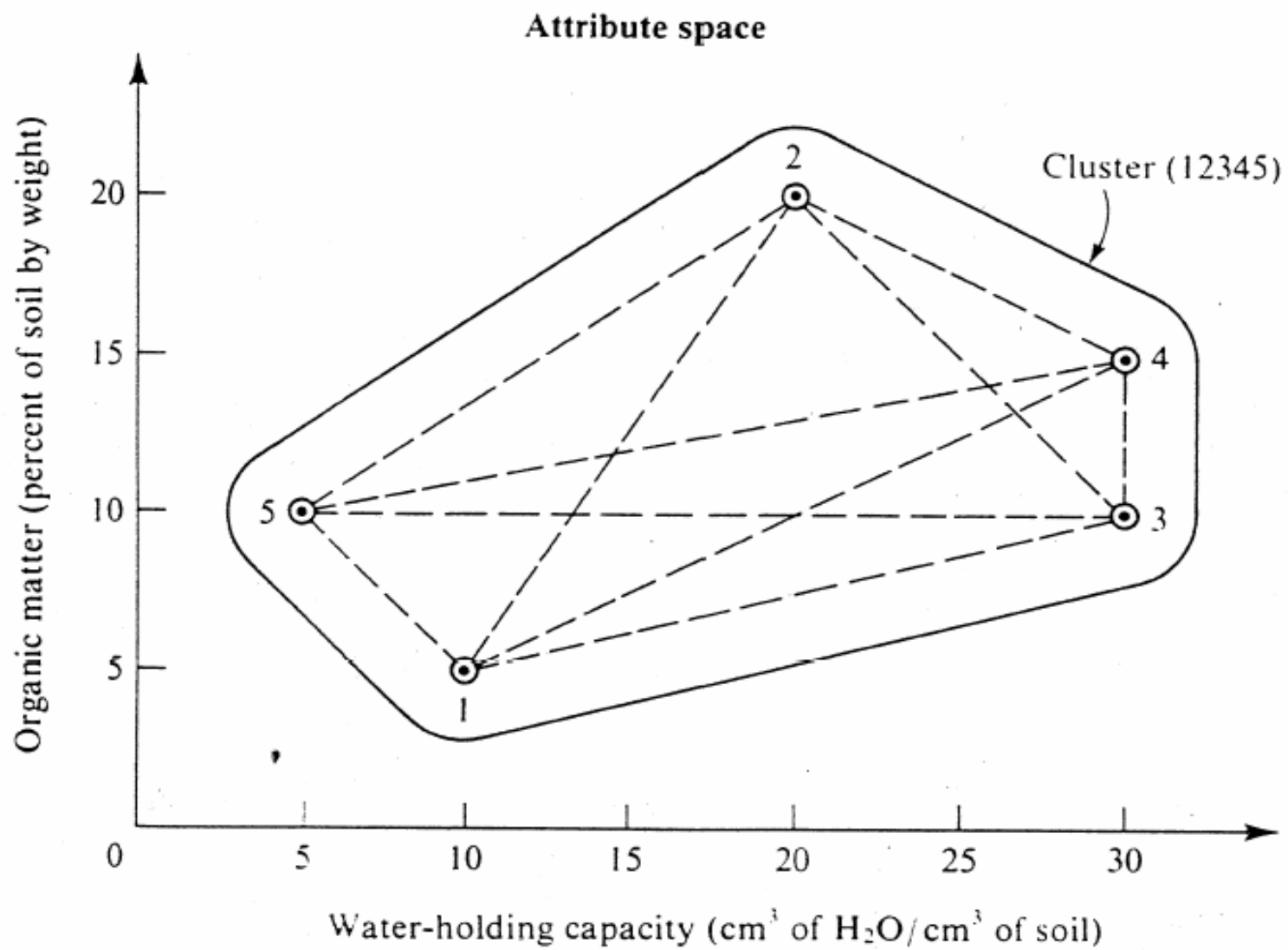
Revised resemblance matrix

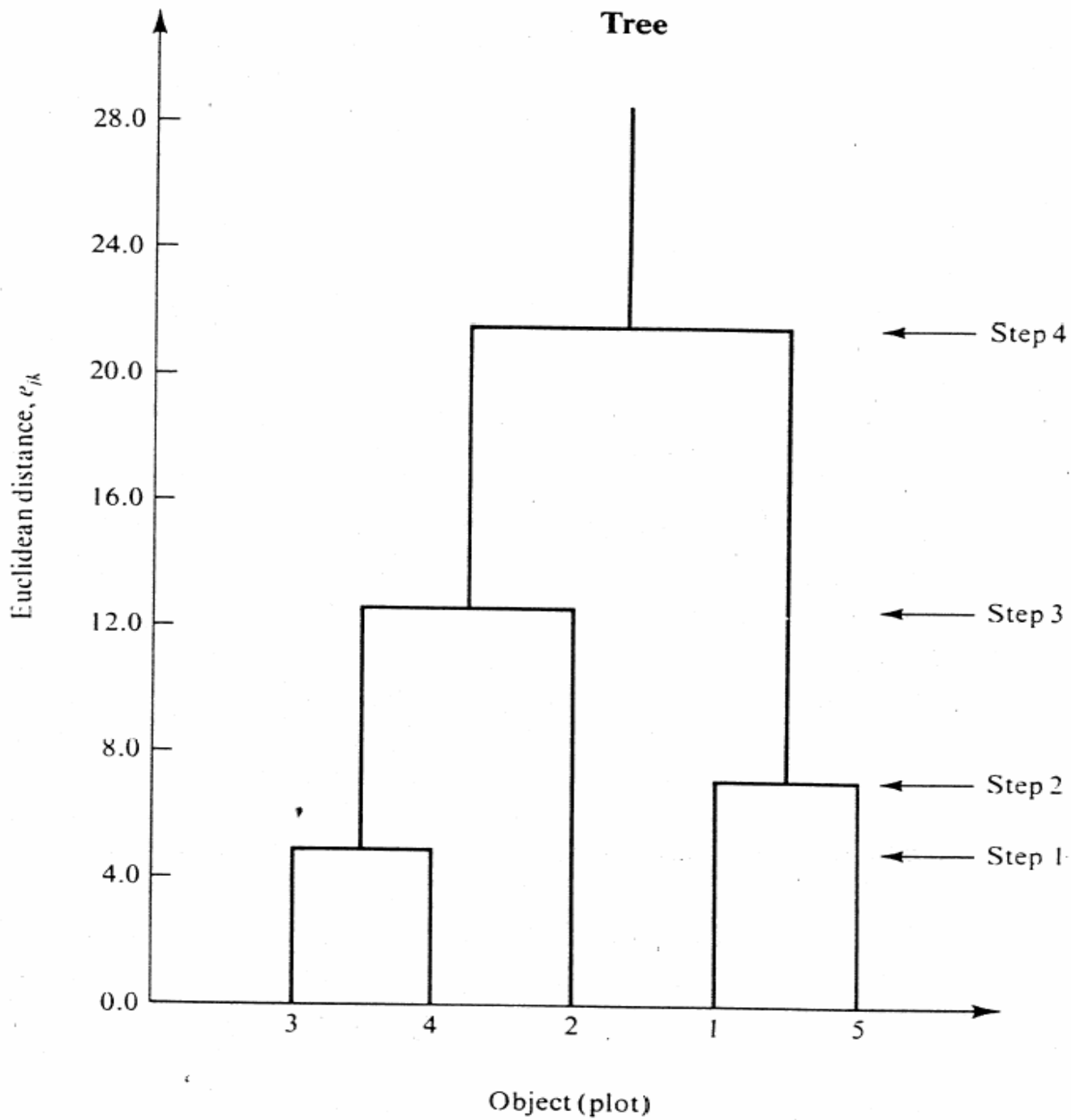
	2	(34)	(15)
2	-	-	-
(34)	12.7	-	-
(15)	18.0	23.4	-



Rvised resemblance matrix

	(15)	(234)
(15)	-	-
(234)	21.6	-





Rearranged data matrix

		Object				
		3	4	2	1	5
Attribute	1	30	30	20	10	5
	2	10	15	20	5	10

Rearranged resemblance matrix

		Object				
		3	4	2	1	5
Object	3	—	—	—	—	—
	4	5.00	—	—	—	—
	2	14.1	11.2	—	—	—
	1	20.6	22.4	18.0	—	—
	5	25.0	25.5	18.0	7.07	—

Cophenetic correlation coefficient (Pearson product-moment correlation coefficient)

- A comparison of the similarities according to the similarity matrix and the similarities according to the dendrogram

$$r_{X,Y} = \frac{\sum xy - (1/n)(\sum x)(\sum y)}{\sqrt{(\sum x^2 - (1/n)(\sum x)^2)(\sum y^2 - (1/n)(\sum y)^2)}}$$

NTSYS

- Import matrix
- Transpose matrix if objects are rows (they are supposed to be columns in NTSYS) (transp in transformation / **general**)
- Consider log1 or autoscaling (standardization)
- Select similarity or distance measure (**similarity**)
- Produce similarity matrix

NTSYS (continued)

- Select clustering procedure (often UPGMA) (**clustering**)
- Calculate cophenetic matrix (**clustering**)
- Compare similarity matrix with cophenetic matrix (made from the dendrogram) and write down the cophenetic correlation (**graphics, matrix comparison**)
- Write dendrogram (**graphics, treeplot**)