# Clustering, Distance Methods, and Ordination

Shyh-Kang Jeng

Department of Electrical Engineering/
Graduate Institute of Communication/
Graduate Institute of Networking and
Multimedia

---

# Outlinet

* Introduction
* Similarity Measures
* Hierarchical Clustering Methods
* Nonhierarchical Clustering Methods
* Clustering Based on Statistical Models
* Multidimensional Scaling

---

# Outlinet

* Introduction
* Similarity Measures
* Hierarchical Clustering Methods
* Nonhierarchical Clustering Methods
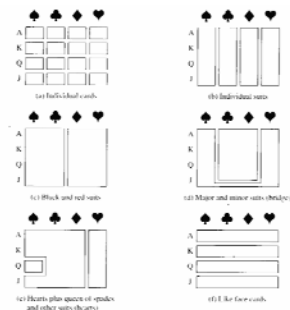* Clustering Based on Statistical Models
* Multidimensional Scaling

---

# Clustering

* Searching data for a structure of "natural" groupings
* An exploratory technique
* Provides means for
  - Assessing dimensionality
  - Identifying outliers
  - Suggesting interesting hypotheses concerning relationships

---

# Classification vs. Clustering

* Classification
  - Known number of groups
  - Assign new observations to one of these groups
* Cluster analysis
  - No assumptions on the number of groups or the group structure
  - Based on similarities or distances (dissimilarities)

---

# Difficulty in Natural Grouping

## Outlinet

- Introduction
- <u>Similarity Measures</u>
- Hierarchical Clustering Methods
- Nonhierarchical Clustering Methods
- Clustering Based on Statistical Models
- Multidimensional Scaling

## Choice of Similarity Measure

- Nature of variables
  - Discrete, continuous, binary
- Scale of measurement
  - Nominal, ordinal, interval, ratio
- Subject matter knowledge
- Items: proximity indicated by some sort of distance
- Variables: grouped by correlation coefficient or measures of association

## Some Well-known Distances

- Euclidean distance
$$d(\mathbf{x},\mathbf{y}) = \sqrt{(\mathbf{x}-\mathbf{y})'(\mathbf{x}-\mathbf{y})}$$
- Statistical distance
$$d(\mathbf{x},\mathbf{y}) = \sqrt{(\mathbf{x}-\mathbf{y})'\mathbf{A}(\mathbf{x}-\mathbf{y})}$$
- Minkowski metric
$$d(\mathbf{x},\mathbf{y}) = \left[\sum_{i=1}^{p}|x_i - y_i|^m\right]^{1/m}$$

## Two Popular Measures of Distance for Nonnegative Variables

- Canberra metric
$$d(\mathbf{x},\mathbf{y}) = \sum_{i=1}^{p}\frac{|x_i - y_i|}{x_i + y_i}$$
- Czekanowski coefficient
$$d(\mathbf{x},\mathbf{y}) = 1 - \frac{2\sum_{i=1}^{p}\min(x_i, y_i)}{\sum_{i=1}^{p}(x_i + y_i)}$$

## A Caveat

- Use "true" distances when possible
  - i.e., distances satisfying distance properties
- Most clustering algorithms will accept subjectively assigned distance numbers that may not satisfy, for example, the triangle inequality

## Example of Binary Variable

|  | Variable | | | | |
| --- | --- | --- | --- | --- | --- |
|  | 1 | 2 | 3 | 4 | 5 |
| Item $i$ | 1 | 0 | 0 | 1 | 1 |
| Item $j$ | 1 | 1 | 0 | 1 | 0 |

## Squared Euclidean Distance for Binary Variables

☛ Squared Euclidean distance

$$d(\mathbf{x}_i, \mathbf{x}_k) = \sum_{j=1}^{p} (x_{ij} - x_{kj})^2$$

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

☛ Suffers from weighting the 1-1 and 0-0 matches equally
– e.g., two people both read ancient Greek is stronger evidence of similarity than the absence of this capability

## Contingency Table

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

## Some Binary Similarity Coefficients

| 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).]

15

## Example 12.1

|  | 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 |

## Example 12.1

$$X_1 = \begin{cases} 1 & \text{height} \geq 72 \text{ in.} \\ 0 & \text{height} < 72 \text{ in.} \end{cases}$$

$$X_4 = \begin{cases} 1 & \text{blond hair} \\ 0 & \text{not blond hair} \end{cases}$$

$$X_2 = \begin{cases} 1 & \text{weight} \geq 150 \text{ lb} \\ 0 & \text{weight} < 150 \text{ lb} \end{cases}$$

$$X_5 = \begin{cases} 1 & \text{right handed} \\ 0 & \text{left handed} \end{cases}$$

$$X_3 = \begin{cases} 1 & \text{brown eyes} \\ 0 & \text{otherwise} \end{cases}$$

$$X_6 = \begin{cases} 1 & \text{female} \\ 0 & \text{male} \end{cases}$$

## Example 12.1

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

Individual 2

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

18

3

# Example 12.1: Similarity Matrix with Coefficient 1

$$
\begin{array}{c}
\phantom{1}\\1\\2\\3\\4\\5
\end{array}
\begin{array}{ccccc}
1 & 2 & 3 & 4 & 5 \\
\left[\begin{array}{c}1 \\ 1/6 \\ 4/6 \\ 4/6 \\ 0\end{array}\right. & \begin{array}{c} \\ 1 \\ 3/6 \\ 3/6 \\ 5/6\end{array} & \begin{array}{c} \\ \\ 1 \\ 2/6 \\ 2/6\end{array} & \begin{array}{c} \\ \\ \\ 1 \\ 2/6\end{array} & \left.\begin{array}{c} \\ \\ \\ \\ 1\end{array}\right]
\end{array}
$$

19

---

# Conversion of Similarities and Distances

- Similarities from distances
  - e.g., $\widetilde{s}_{ik} = 1/(1 + d_{ik})$
- "True" distances from similarities
  - Matrix of similarities must be nonnegative definite
  - e.g., $d_{ik} = \sqrt{2(1 - \widetilde{s}_{ik})}, \quad \widetilde{s}_{ii} = 1$

20

---

# Contingency Table

|  |  | Variable $k$ | | Totals |
|---|---|---|---|---|
|  |  | 1 | 0 |  |
| Variable $i$ | 1 | $a$ | $b$ | $a + b$ |
|  | 0 | $c$ | $d$ | $c + d$ |
| Totals |  | $a+c$ | $b+d$ | $n = a + b + c + d$ |

---

# Product Moment Correlation as a Measure of Similarity

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

- Related to the chi-square statistic ($r^2 = \chi^2/n$) for testing independence
  - For $n$ fixed, large similarity is consistent with presence of dependence

22

---

# Example 12.2 Similarities of 11 Languages

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 | piec | negy | nelja | |
| 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 | |

23

---

# Example 12.2 Similarities of 11 Languages

Concordant First Letters for Numbers in 11 Languages

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

24

4

## Agglomerative Methods

- Initially a many clusters as objects
- The most similar objects are first grouped
- Initial groups are merged according to their similarities
- Eventually, all subgroups are fused into a single cluster
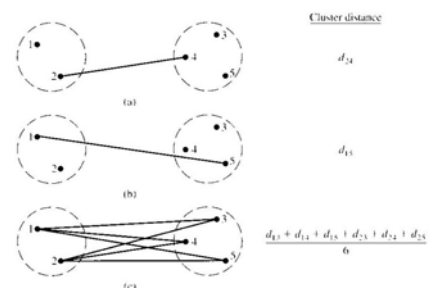
25

## Outlinet

- Introduction
- Similarity Measures
- <u>Hierarchical Clustering Methods</u>
- Nonhierarchical Clustering Methods
- Clustering Based on Statistical Models
- Multidimensional Scaling

## Divisive Methods

- Initial single group is divided into two subgroups such that objects in one subgroup are "far from" objects in the other
- These subgroups are then further divided into dissimilar subgroups
- Continues until there are as many subgroups as objects

27

## Inter-cluster Distance for Linkage Methods



8

## Example 12.3: Single Linkage

$$\mathbf{D} = \{d_{ik}\} = \begin{array}{c} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array}\begin{array}{ccccc} 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{ccccc} 0 & & & & \\ 9 & 0 & & & \\ 3 & 7 & 0 & & \\ 6 & 5 & 9 & 0 & \\ 11 & 10 & ⨀2 & 8 & 0 \end{array}\right] \end{array}$$

29

## Example 12.3: Single Linkage

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

30

## Example 12.3: Single Linkage

$$
\begin{array}{c c}
 & (135) \quad 2 \quad\; 4 \\
\begin{array}{c}(135)\\ 2\\ 4\end{array} &
\begin{bmatrix}
0 & & \\
7 & 0 & \\
6 & \textcircled{5} & 0
\end{bmatrix}
\end{array}
$$

---

## Example 12.3: Single Linkage

$$
\begin{array}{c c}
 & (135) \quad (24) \\
\begin{array}{c}(135)\\ (24)\end{array} &
\begin{bmatrix}
0 & \\
\textcircled{6} & 0
\end{bmatrix}
\end{array}
$$

---

## Example 12.3: Single Linkage
## Resultant Dendrogram

---

## Example 12.4
## Single Linkage of 11 Languages

|     | 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  |

---

## Example 12.4
## Single Linkage of 11 Languages

---

## Pros and Cons of
## Single Linkage



(a) Single linkage confused by near overlap    (b) Chaining effect

| Example 12.5: Complete Linkage |
|---|

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

37

| Example 12.5: Complete Linkage |
|---|

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

38

| Example 12.5: Complete Linkage |
|---|

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

39

| Example 12.5: Complete Linkage |
|---|



40

| Example 12.6 Complete Linkage of 11 Languages |
|---|



41

| Example 12.7 Clustering Variables |
|---|



42

7

## Example 12.7
### Correlations of Variables

| | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ | $X_7$ | $X_8$ |
|---|---|---|---|---|---|---|---|---|
| | 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 |

## Example 12.7:
### Complete Linkage Dendrogram



## Average Linkage

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

## Example 12.8
### Average Linkage of 11 Languages



## Example 12.9
### Average Linkage of Public Utilities



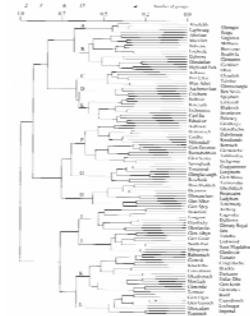## Example 12.9
### Average Linkage of Public Utilities

## Ward's Hierarchical Clustering Method

- For a given cluster $k$, let $\mathrm{ESS}_k$ be the sum of the squared deviation of every item in the cluster from the cluster mean
- At each step, the union of every possible pair of clusters is considered
- The two clusters whose combination results in the smallest increase in the sum of $\mathrm{Ess}_k$ are joined

49

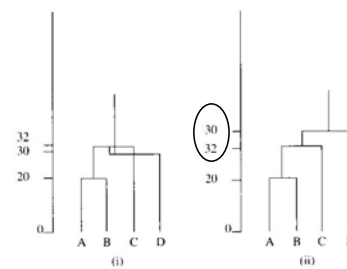## Example 12.10
### Ward's Clustering Pure Malt ScotchWhiskies



50

## Final Comments

- Sensitive to outliers, or "noise points"
- No reallocation of objects that may have been "incorrectly" grouped at an early stage
- Good idea to try several methods and check if the results are roughly consistent
- Check stability by perturbation

51

## Inversion



52

## Outlinet

- Introduction
- Similarity Measures
- Hierarchical Clustering Methods
- <u>Nonhierarchical Clustering Methods</u>
- Clustering Based on Statistical Models
- Multidimensional Scaling

## K-means Method

- Partition the items into $K$ initial clusters
- Proceed through the list of items, assigning an item to the cluster whose centroid is nearest
- Recalculate the centroid for the cluster receiving the new item and for the cluster losing the item
- Repeat until no more reassignment

54

## Example 12.11
### *K*-means Method

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

## Example 12.11
### *K*-means Method

| | Coordinates of Centroid | |
|---|---|---|
| Cluster | x1 | x2 |
| (AB) | (5+(-1))/2 = 2 | (3+1)/2 = 2 |
| (CD) | (1+(-3))/2=-1 | (-2+(-2))/2=-2 |

## Example 12.11
### *K*-means Method

$$\overline{x}_{i,new} = \frac{n\overline{x}_i + x_{ji}}{n+1}$$

if the *j*th item is added to a group

$$\overline{x}_{i,new} = \frac{n\overline{x}_i - x_{ji}}{n-1}$$

if the *j*th item is removed from a group

## Example 12.11
### Final Clusters

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

## F Score

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

## Outlinet

- Introduction
- Similarity Measures
- Hierarchical Clustering Methods
- Nonhierarchical Clustering Methods
- Clustering Based on Statistical Models
- Multidimensional Scaling

## Normal Mixture Model

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

$$p_k \geq 0, \quad \sum_{k=1}^{K} p_k = 1, \quad f_k(\mathbf{x}) : N_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$f_{Mix}(\mathbf{x} \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \cdots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$$

$$= \sum_{k=1}^{K} \frac{p_k}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \right)$$

## Likelihood

$$L(p_1, \cdots, p_k, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \cdots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$$

$$= \prod_{j=1}^{N} f_{Mix}(\mathbf{x}_j \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \cdots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K)$$

$$= \prod_{j=1}^{N} \sum_{k=1}^{K} \frac{p_k}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left( -\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_j - \boldsymbol{\mu}_k) \right)$$

## Statistical Approach

Obtain the maximum likelihood estimates and

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

Determine $K$ via maximizing

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

or

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

## BIC for Special Structures

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

## Software Package *MCLUST*

- Combines hierarchical clustering, EM algorithm, and BIC
- In the E step of EM, a matrix is created whose $j$th row contains the estimates of the conditional probabilities that observation $\mathbf{x}_j$ belongs to cluster 1, 2, . . ., $K$
- At convergence $\mathbf{x}_j$ is assigned to cluster $k$ for which the conditional probability of membership is largest

## Example 12.13
## Clustering of Iris Data

$$K = 3$$
$$p = 4, \quad \boldsymbol{\Sigma}_k = \eta_k \mathbf{I}, \quad k = 1,2,3$$
$$\hat{\eta}_1 = 0.076, \quad \hat{\eta}_2 = 0.163, \quad \hat{\eta}_3 = 0.163$$
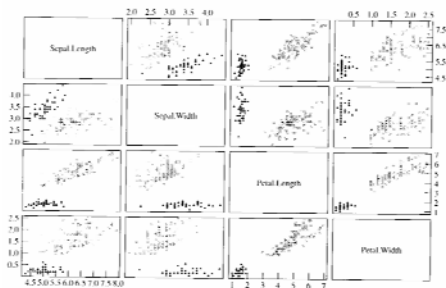$$\hat{p}_1 = 0.3333, \quad \hat{p}_2 = 0.4133, \quad \hat{p}_3 = 0.2534$$
$$\text{BIC} = -853.8$$

$$\boldsymbol{\mu}_1 = \begin{bmatrix} 5.01 \\ 3.43 \\ 1.46 \\ 0.25 \end{bmatrix}, \quad \boldsymbol{\mu}_2 = \begin{bmatrix} 5.90 \\ 2.75 \\ 4.40 \\ 1.43 \end{bmatrix}, \quad \boldsymbol{\mu}_3 = \begin{bmatrix} 6.85 \\ 3.07 \\ 5.73 \\ 2.07 \end{bmatrix},$$
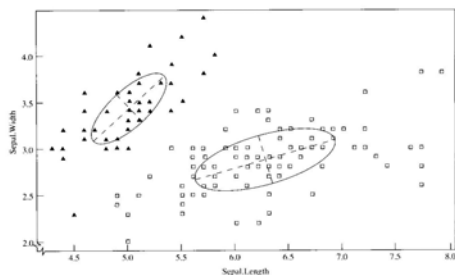
## Example 12.13
## Clustering of Iris Data



67

## Example 12.13
## Clustering of Iris Data

$K = 2$ is the best solution
$\hat{p}_1 = 0.3333, \quad \hat{p}_2 = 0.6667$

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

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

68

## Example 12.13
## Clustering of Iris Data



69

## Outlinet

- Introduction
- Similarity Measures
- Hierarchical Clustering Methods
- Nonhierarchical Clustering Methods
- Clustering Based on Statistical Models
- Multidimensional Scaling

## Multidimensional Scaling (MDS)

- Displays (transformed) multivariate data in low-dimensional space
- Different from plots based on PC
  - Primary objective is to "fit" the original data into low-dimensional system
  - Distortion caused by reduction of dimensionality is minimized
- Distortion
  - Similarities or dissimilarities among data

71

## Multidimensional Scaling

- Given a set of similarities (or distances) between every pair of $N$ items
- Find a representation of the items in few dimensions
- Inter-item proximities "nearly" match the original similarities (or distances)

72

12

## Non-metric and Metric MDS

- Non-metric MDS
  - Uses only the rank orders of the *N(N-1)/2* original similarities and not their magnitudes
- Metric MDS
  - Actual magnitudes of original similarities are used
  - Also known as principal coordinate analysis

## Objective

$N$ items, $M = N(N-1)/2$ similarities

Assume no ties, and arrange

$$s_{i_1 k_1} < s_{i_2 k_2} < \cdots < s_{i_M k_M}$$

Find a $q$ - dimensional configuration, such that

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

## Kruskal's Stress

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

$\hat{d}_{ik}^{(q)}$ are numbers known to satisfy the ordering

They are not distances, and merely reference numbers

## Takane's Stress

$$\text{SStress} = \left\{ \frac{\sum_k \sum_{i<k} \left( d_{ik}^2 - \hat{d}_{ik}^2 \right)^2}{\sum_k \sum_{i<k} d_{ik}^4} \right\}^{1/2}$$

## Basic Algorithm

- Obtain and order the *M* pairs of similarities
- Try a configuration in $q$ dimensions
  - Determine inter-item distances and reference numbers
  - Minimize Kruskal's or Takane's stress
- Move the points around to obtain an improved configuration
- Repeat until minimum stress is obtained

## Example 12.14
## MDS of U.S. Cities
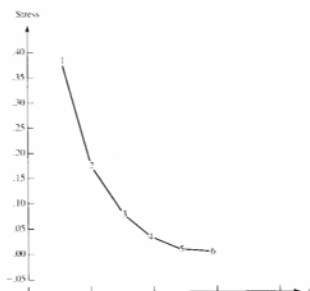
## Example 12.14
## MDS of U.S. Cities



## Example 12.14
## MDS of U.S. Cities



## Example 12.15
## MDS of Public Utilities



## Example 12.15
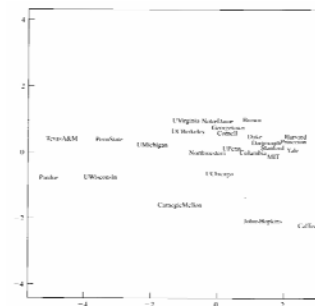## MDS of Public Utilities
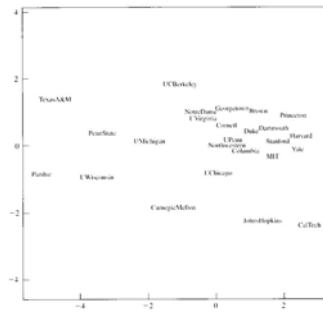


## Example 12.16
## MDS of Universities



## Example 12.16
## Metric MDS of Universities

# Example 12.16
## Non-metric MDS of Universities