

## Clustering

Steps (cont'd)

- Verify clustering tendency (Sec. 16.6)
- Clustering algorithm: Chapters 12–15
- <u>Cluster validation</u>: Verify that choices of alg. params. & cluster shape match data's clustering structure (Chapt. 16)
- Interpretation: The expert interprets results with other information
- Warning: <u>Each step is subjective</u> and depends on expert's biases!



- **Clustering** Applications
- <u>Data reduction</u> (compression): Represent each cluster with single item
- Suggest hypotheses about nature of data
- <u>Test hypotheses</u> about data, e.g. that certain feats. are correlated while others are independent
- Prediction based on groups: e.g. Slide 7.2

**Clustering** Types of Features

- <u>Nominal</u>: Name only, no quantitative comparisons possible, e.g. {male, female}
- <u>Ordinal</u>: Can be meaningfully ordered, but no quantitative meaning on the differences, e.g. {4, 3, 2, 1} to represent {excellent, very good, good, poor}
- <u>Interval-scaled</u>: Difference is meaningful, ratio is not, e.g. temperature measures on Celsius scale
- <u>Ratio-scaled</u>: Difference and ratio both meaningful, e.g. weight
- Each type possesses the properties of the preceding types

7

Clustering Cluster Types

- Start with  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and place into m clusters  $C_1, \dots, C_m$
- Type 1: <u>Hard</u> (crisp)

$$C_i \neq \emptyset, \ i = 1, \dots, m \qquad \bigcup_{i=1}^m C_i = X$$
$$C_i \cap C_j = \emptyset, \ i \neq j, \ i, j \in \{1, \dots, m\}$$

- F.v.'s in  $C_i$  "more similar" to others in  $C_i$  than those in  $C_j$ ,  $j \neq i$
- Type 2: <u>Fuzzy</u>:  $C_j$  has <u>membership function</u>  $\mu_j: X \to [0, 1]$  s.t.

$$\sum_{j=1}^{m} \mu_j(\mathbf{x}_i) = 1, \ i \in \{1, \dots, N\}$$
$$0 < \sum_{i=1}^{N} \mu_j(\mathbf{x}_i) < N, \ j \in \{1, \dots, m\}$$

8

6

Proximity Measures Definitions	Proximity Measures Definitions (cont'd)
• Dissimilarity measure is func. $d : X \times X \to \Re$ s.t. $\exists d_0 \in \Re : -\infty < d_0 \le d(\mathbf{x}, \mathbf{y}) < +\infty, \ \forall \mathbf{x}, \mathbf{y} \in X$ $d(\mathbf{x}, \mathbf{x}) = d_0 \ \forall \mathbf{x} \in X$	<ul> <li>Can also define proximity measures between sets of f.v.'s</li> </ul>
$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})  \forall \mathbf{x}, \mathbf{y} \in X$ • $d$ is a <u>metric DM</u> if $d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$ and $d(\mathbf{x}, \mathbf{z}) \le d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})  \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$	• Let $U = \{D_1, \dots, D_k\}, D_i \subset X$ , PM $\alpha : U \times U \to \Re$
- E.g. $d_2(\cdot, \cdot)$ = Euclidean distance, $d_0 = 0$	• E.g. $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}, U = \{\{x_1, x_2\}, \{x_1, x_4\}, \{x_3, x_4, x_5\}, \{x_1, x_2, x_3, x_4, x_5\}\},\$
• Similarity measure is func. $s: X \times X \to \Re$ s.t.	$d_{min}^{ss} \left( D_i, D_j \right) = \min_{\mathbf{x} \in D_i,  \mathbf{y} \in D_j} d_2(\mathbf{x}, \mathbf{y})$
$\exists s_0 \in \Re : -\infty < s(\mathbf{x}, \mathbf{y}) \le s_0 < +\infty, \ \forall \mathbf{x}, \mathbf{y} \in X$ $s(\mathbf{x}, \mathbf{x}) = s_0 \ \forall \mathbf{x} \in X$ $s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}) \ \forall \mathbf{x}, \mathbf{y} \in X$	• Min. value is $d_{min,0}^{ss} = 0$ , $d_{min}^{ss} (D_i, D_i) = d_{min,0}^{ss}$ , and $d_{min}^{ss} (D_i, D_j) = d_{min}^{ss} (D_j, D_i)$ , so $d_{min}^{ss} (\cdot, \cdot)$ is a DM
• s is a metric SM if $s(\mathbf{x}, \mathbf{y}) = s_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$ and $s(\mathbf{x}, \mathbf{y}) s(\mathbf{y}, \mathbf{z}) \leq [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})] s(\mathbf{x}, \mathbf{z})$ $\forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$	• However, $d_{min}^{ss} (\{\mathbf{x}_1, \mathbf{x}_2\}, \{\mathbf{x}_1, \mathbf{x}_4\}) = d_{min,0}^{ss}$ and $\{\mathbf{x}_1, \mathbf{x}_2\} \neq \{\mathbf{x}_1, \mathbf{x}_4\}$ , so not a metric DM
9	10
Proximity Measures Between Points	

Real-Valued Vectors Example Dissimilarity Measures (pp. 361–362)

• Common, general-purpose metric DM is weighted *L<sub>p</sub>* norm:

$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^{\ell} w_i |x_i - y_i|^p\right)^{1/p}$$

• Special cases include weighted Euclidian distance (p = 2), weighted Manhattan distance

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{\ell} w_i |x_i - y_i|,$$

and weighted  $L_{\infty}$  norm

$$d_{\infty}(\mathbf{x}, \mathbf{y}) = \max_{1 \le i \le \ell} \{ w_i | x_i - y_i | \}$$

• Generalization of weighted  $L_2$  norm is

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T B(\mathbf{x} - \mathbf{y})}$$
,

e.g. Mahalanobis distance

Proximity Measures Between Points Real-Valued Vectors Example Similarity Measures (pp. 362–363)

• Inner product:

$$s_{inner}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^{\ell} x_i y_i$$

- If  $\|\mathbf{x}\|_2, \|\mathbf{y}\|_2 \leq a$ , then  $-a^2 \leq s_{inner}(\mathbf{x}, \mathbf{y}) \leq a^2$
- Tanimoto distance:

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2^2 + \|\mathbf{y}\|_2^2 - \mathbf{x}^T \mathbf{y}} = \frac{1}{1 + \frac{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}{\mathbf{x}^T \mathbf{y}}},$$

which is inversely prop. to (squared Euclid. dist.)/(correlation measure)

11

	Proximity Measures Between Points Fuzzy Measures
Proximity Measures Between Points Discrete-Valued Vectors	<ul> <li>Let x<sub>i</sub> ∈ [0, 1] be measure of how much x possesses ith feature</li> </ul>
<ul> <li>If the coordinates of f.v.'s come from {0,,k-1}, can use SMs and DMs defined for real-valued f.v.'s, (e.g. weighted Lp norm) plus:</li> <li><u>Hamming distance</u>: DM measuring number of places where x and y differ</li> <li><u>Tanimoto measure</u>: SM measuring number of places where x and y are same, divided by total number of places</li> <li>* Ignore places i where x<sub>i</sub> = y<sub>i</sub> = 0</li> <li>Useful for ordinal features where x<sub>i</sub> is degree to which x possesses ith feature</li> </ul>	• If $x_i, y_i \in \{0, 1\}$ , then $(x_i \equiv y_i) = ((\neg x_i \land \neg y_i) \lor (x_i \land y_i))$ • Generalize to fuzzy values: $s(x_i, y_i) = \max \{\min \{1 - x_i, 1 - y_i\}, \min \{x_i, y_i\}\}$ • To measure similarity between vectors: $s_F^p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^{\ell} s(x_i, y_i)^p\right)^{1/p}$ $(\ell^{1/p})/2 \leq s_F^q(\cdot, \cdot) \leq \ell^{1/p}$ • So $s_F^{\infty} = \max_{1 \leq i \leq \ell} s(x_i, y_i)$ and $s_F^1 = \sum_{i=1}^{\ell} s(x_i, y_i) = \text{generalization of Hamming distance}$
13	14
Prox. Measures Between a Point and a Set	<b>Prox. Measures Between a Point and a Set</b> Representatives
<ul> <li>Might want to measure proximity of point x to existing cluster C</li> <li>Can measure proximity α by using <u>all points</u> of C or by using a <u>representative</u> of C</li> </ul>	<ul> <li>Alternative: Measure distance between point x and a <u>representative</u> of the set C</li> <li>Appropriate choice of representative depends</li> </ul>
• If all points of $C$ used, common choices: $\alpha_{max}^{ps}(\mathbf{x}, C) = \max_{\mathbf{y} \in C} \{\alpha(\mathbf{x}, \mathbf{y})\}$ $\alpha_{min}^{ps}(\mathbf{x}, C) = \min_{\mathbf{y} \in C} \{\alpha(\mathbf{x}, \mathbf{y})\}$ $\alpha_{avg}^{ps}(\mathbf{x}, C) = \frac{1}{ C } \sum_{\mathbf{y} \in C} \alpha(\mathbf{x}, \mathbf{y}) ,$ where $\alpha(\mathbf{x}, \mathbf{y})$ is any measure between $\mathbf{x}$ and $\mathbf{y}$	on type of cluster Compact Elongated <u>Point</u> <u>Hyperplane</u> · · · · · · + · · (a) (b) (c)
15	16



## Overview of Clustering Algorithms

Exhaustive Search

- Want to find set of clusters that maximizes SM or minimizes DM
- Option 1: Try all possible clusters of size *m* for various values of *m*
- Number of ways to partition N items into m nonempty subsets is exactly given by the <u>Stirling</u> <u>numbers of the second kind</u>, which are:

$$\gg {N \choose m} \geq \left(\frac{N}{m}\right)^m$$

• Thus brute-force approach infeasible

## Overview of Clustering Algorithms Categories of Algorithms

- <u>Sequential algorithms</u> (Chapt. 12) produce a single clustering, use straightforward greedy approaches, and output depends on the order the f.v.'s are presented to the algorithm
- <u>Hierarchical algorithms</u> (Chapt. 13) produce a sequence (hierarchy) of clusterings, and are of two types:
  - <u>Agglomerative</u>: Repeatedly merge two clusters into one
  - <u>Divisive</u>: Repeatedly divide one cluster into two
- Algorithms based on cost function optimization (Chapt. 14) evaluate the goodness of a clustering with a cost function, typically *m* is fixed
  - <u>Crisp</u> (hard): Each f.v. belongs to only one cluster, e.g. Isodata algorithm
  - <u>Fuzzy</u>: Each f.v. can belong to a cluster up to a certain degree, as indicated by a <u>membership function</u>
  - Many more
- Other various methods (Chapt. 15)

22

21