# Proximity of data in clustering analysis

- Interval-scaled variables
- Binary variables
- Nominal, ordinal, and ratio variables
- Variables of mixed types

# Interval-scaled variables

- Standardize data
  - Calculate the mean absolute deviation:
    $$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + ... + |x_{nf} - m_f|)$$
    where
    $$m_f = \frac{1}{n}(x_{1f} + x_{2f} + ... + x_{nf})$$
  - Calculate the standardized measurement (z-score)
    $$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- Using mean absolute deviation is more robust than using standard deviation

# Object similarity measures

- Distances are normally used to measure the similarity or dissimilarity between two data objects
- Some popular ones include Minkowski distance:

$$d(i,j) = \sqrt[q]{(|x_{i_1} - x_{j_1}|^q + |x_{i_2} - x_{j_2}|^q + ... + |x_{i_p} - x_{j_p}|^q)}$$

  where $i = (x_{i1}, x_{i2}, ..., x_{ip})$ and $j = (x_{j1}, x_{j2}, ..., x_{jp})$ are two p-dimensional data objects, and q is a positive integer, and

- If q = 1, d is Manhattan distance

$$d(i,j) = |x_{i_1} - x_{j_1}| + |x_{i_2} - x_{j_2}| + ... + |x_{i_p} - x_{j_p}|$$

---

# Object similarity measures

- If $q = 2$, $d$ is Euclidean distance:

$$d(i,j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + ... + |x_{i_p} - x_{j_p}|^2)}$$

- Metric properties
  - $d(i,j) \geq 0$
  - $d(i,i) = 0$
  - $d(i,j) = d(j,i)$
  - $d(i,j) \leq d(i,k) + d(k,j)$
- Also one can use other dissimilarity measures, such as weighted distance

# Binary variables

■ A contingency table for binary data

|  | | Object $j$ | | |
|---|---|---|---|---|
|  | | 1 | 0 | sum |
|  | 1 | $a$ | $b$ | $a+b$ |
| Object $i$ | 0 | $c$ | $d$ | $c+d$ |
|  | sum | $a+c$ | $b+d$ | $p$ |

■ Simple matching coefficient (invariant, for symmetric binary variables):

$$d(i, j) = \frac{b + c}{a + b + c + d}$$

■ Jaccard coefficient (noninvariant, for asymmetric binary variables):

$$d(i, j) = \frac{b + c}{a + b + c}$$

---

# Dissimilarity between binary variables

■ Example

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|---|---|---|---|---|---|---|---|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- gender is a symmetric attribute
- the remaining attributes are asymmetric binary
- let the values Y and P be set to 1, and the value N be set to 0

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

# Nominal variables

- A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- Method 1: Simple matching
  - $m$: # of matches, $p$: total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- Method 2: use a large number of binary variables
  - creating a new binary variable for each of the $M$ nominal states

# Ordinal variables

- An ordinal variable can be discrete or continuous
- Order is important, e.g., rank $r_{if} \in \{1,..., M_f\}$
- Can be treated like interval-scaled
  - replacing $x_{if}$ by their rank
  - map the range of each variable onto [0, 1] by replacing $i$-th object in the $f$-th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

  - compute the dissimilarity using methods for interval-scaled variables

# Ratio-scaled variables

- Positive measurements on nonlinear scales, approximately at exponential scale, such as $Ae^{Bt}$ or $Ae^{-Bt}$
- Methods:
  - treat them like interval-scaled variables — not a good choice! (why?)
  - apply logarithmic transformation
    $$y_{if} = log(x_{if})$$
  - treat them as continuous ordinal data and treat their rank as interval-scaled.

---

# Variables of mixed types

- A database may contain all the six types of variables
  - symmetric binary, asymmetric binary, nominal, ordinal, interval and ratio.
- One may use a weighted formula to combine their effects.
  $$d(i,j) = \frac{\Sigma_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\Sigma_{f=1}^{p} \delta_{ij}^{(f)}}$$
  - f is binary or nominal:
    - $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, or $d_{ij}^{(f)} = 1$ otherwise
  - f is interval-based: use the normalized distance
  - f is ordinal or ratio-scaled
    - compute ranks $r_{if}$ and
    - and treat $z_{if}$ as interval-scaled     $z_{if} = \dfrac{r_{if} - 1}{M_f - 1}$