

# Klasyfikacja i grupowanie

# Pojęcia

Podobieństwo a odległość:

- odległość obiektów
- rozproszenie obiektów

Definicje miar podobieństwa

Cechy ilościowe

Cechy jakościowe

Porównanie obiektów do zbioru obiektów

Porównanie zbiorów obiektów

# Określenie podobieństwa na podstawie odległości cech obiektów

$$simi(\bullet) = \frac{1}{1 + \alpha dist(\bullet)^\beta}$$

$$simi(\bullet) = 1 - \frac{dist(\bullet)}{\max(dist(\bullet))}$$

$$simi(\bullet) = e^{-\alpha dist(\bullet)}$$

# Problemy określania podobieństwa

- Skalowanie i normalizacja
- Cechy ilościowe
- Cechy jakościowe
- Cechy ilościowo-jakościowe

# Cechy ilościowe

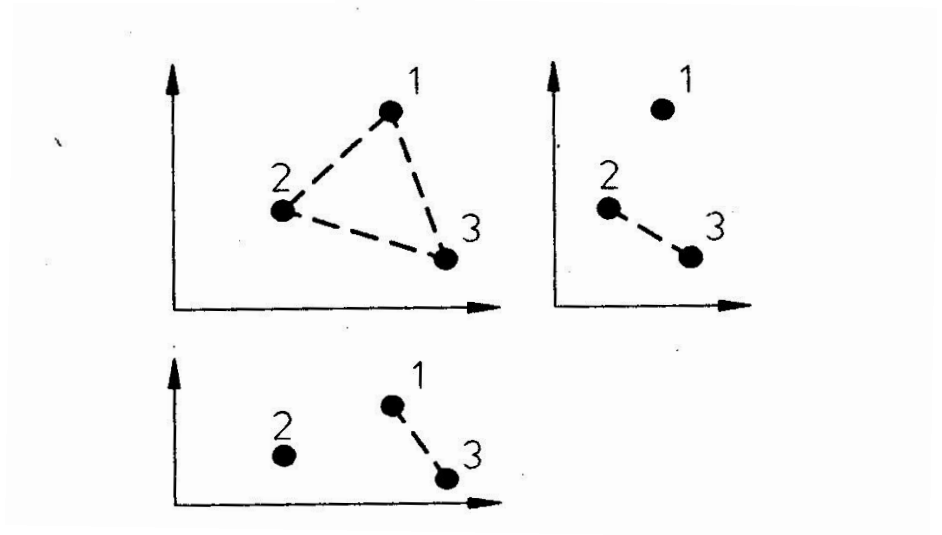
- Miara odległości Minkowskiego:

$$\text{dist}_{\mathbf{M}}(\mathbf{v}_i, \mathbf{v}_j) = \left( \sum_{k=1}^{\dim(\mathbf{v})} |v_{i,k} - v_{j,k}|^s \right)^{1/s}$$

- Miara odległości euklidesowej:

$$\text{dist}_{\mathbf{E}}(\mathbf{v}_i, \mathbf{v}_j) = \sqrt{(\mathbf{v}_i - \mathbf{v}_j)^T (\mathbf{v}_i - \mathbf{v}_j)}$$

- Wpływ skal współrzędnych na wynik grupowania



# Cechy ilościowe

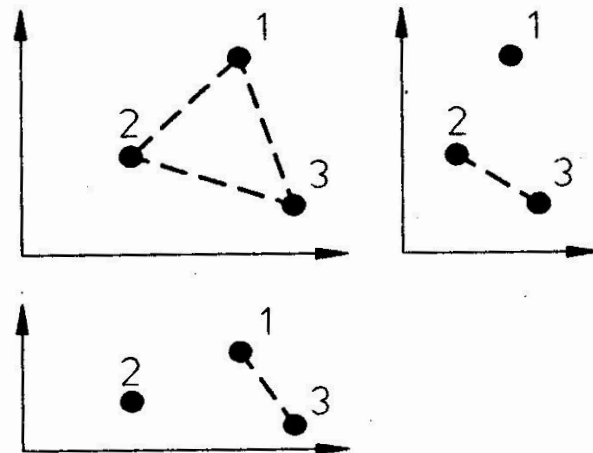
- Miara odległości Minkowskiego:

$$\text{dist}_{\mathbf{M}}(\mathbf{v}_i, \mathbf{v}_j) = \left( \sum_{k=1}^{\dim(\mathbf{v})} |v_{i,k} - v_{j,k}|^s \right)^{1/s}$$

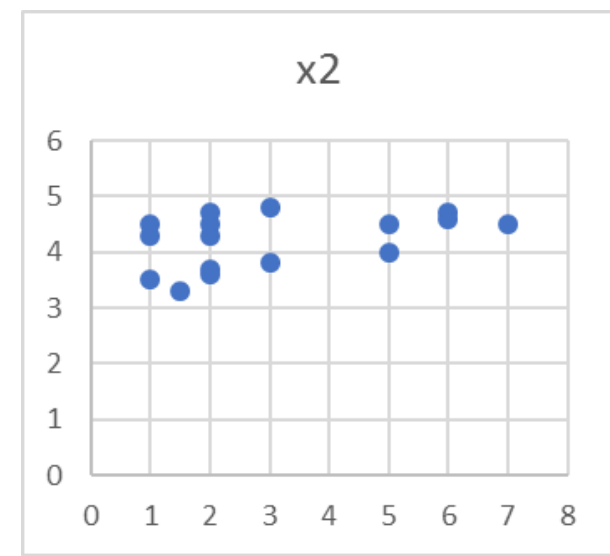
- Miara odległości euklidesowej:

$$\text{dist}_{\mathbf{E}}(\mathbf{v}_i, \mathbf{v}_j) = \sqrt{(\mathbf{v}_i - \mathbf{v}_j)^T (\mathbf{v}_i - \mathbf{v}_j)}$$

- Wpływ skal współrzędnych na wynik grupowania



# Cechy ilościowe



- Normalizacja

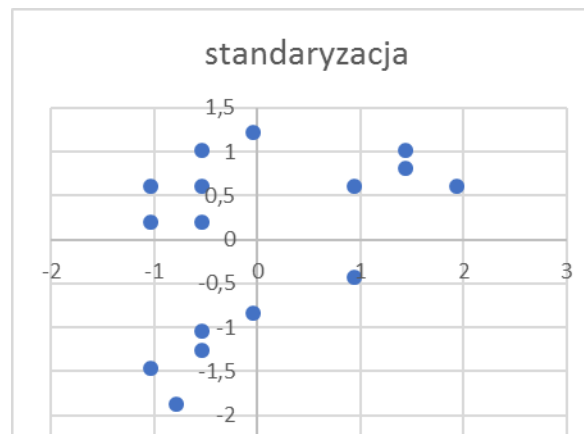
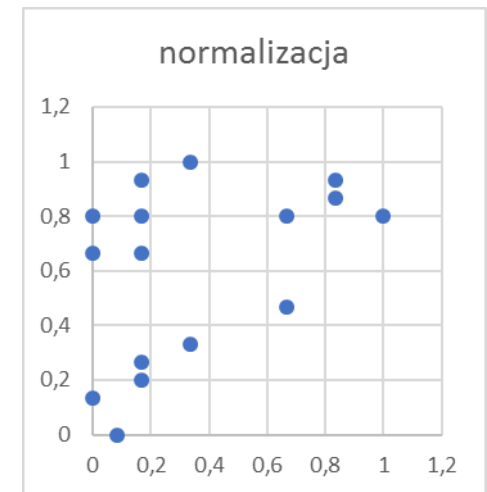
$$x'_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

- Standaryzacja

$$x'_i = \frac{x_i - \bar{x}}{\sigma}$$

gdzie:

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$



# Cechy ilościowe – wagi cech

- Ważona miara odległości

$$\text{dist}_{\mathbf{W}}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = \sqrt{(\underline{\mathbf{v}}_i - \underline{\mathbf{v}}_j)^T \underline{\mathbf{W}} (\underline{\mathbf{v}}_i - \underline{\mathbf{v}}_j)}$$

$$\underline{\mathbf{W}} = \begin{bmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & w_n \end{bmatrix}$$

- Kowariancja

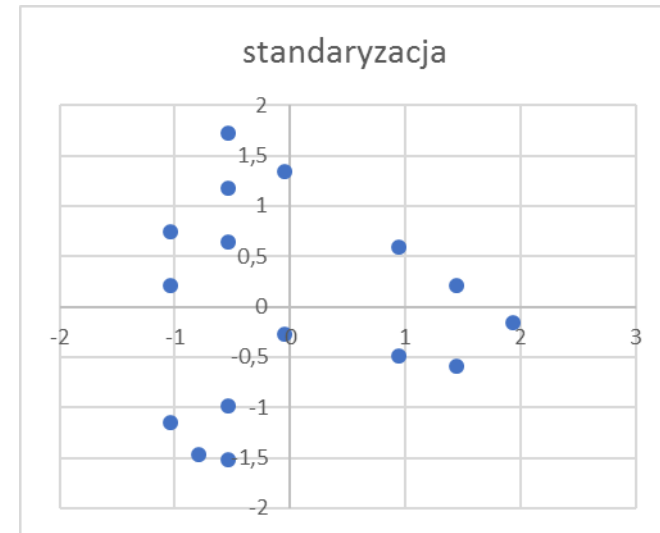
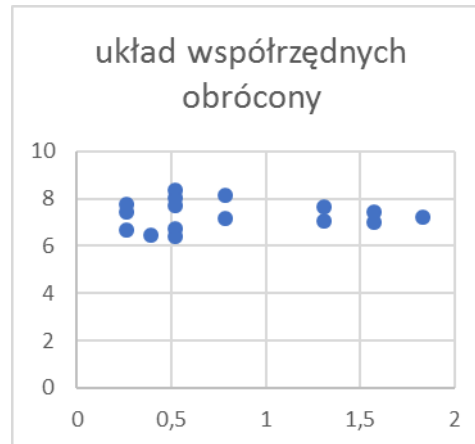
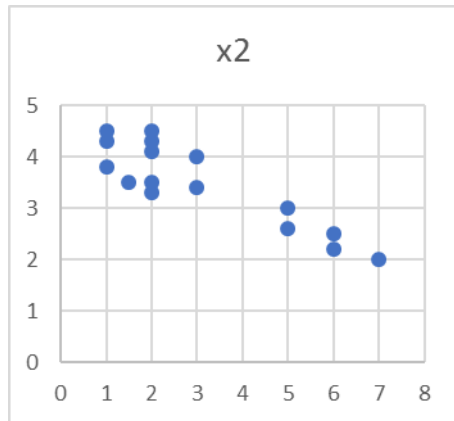
$$\mathbf{W} = \mathbf{C}^{-1}$$

$$\mathbf{C} = \frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$



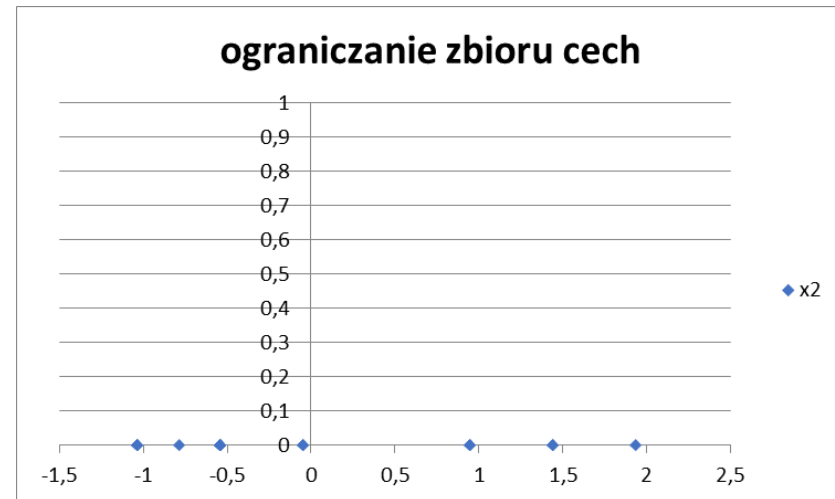
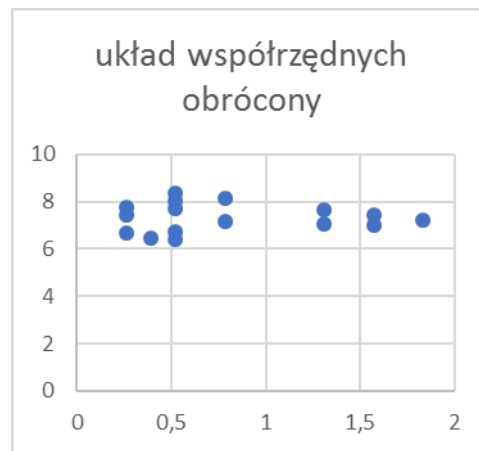
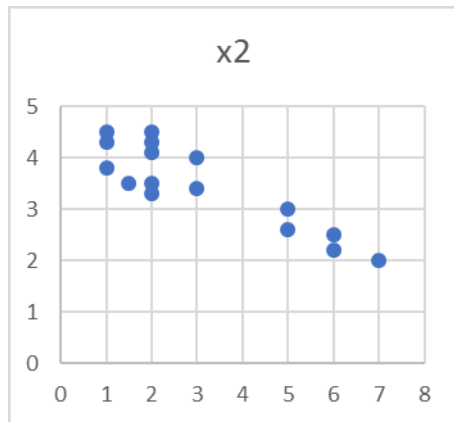
# Cechy ilościowe – wagi cech

- Ważona miara odległości – transformacje przestrzeni wartości cech; osie główne



# Cechy ilościowe – wagi cech

- Ważona miara odległości – ograniczanie zbioru uwzględnianych cech



# Cechy ilościowe – wagi cech

- miara odległości – modyfikacja Hamminga:

$$\text{dist}_{\mathbf{H}}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = \frac{1}{\dim(\underline{\mathbf{v}})} \sum_{k=1}^{\dim(\underline{\mathbf{v}})} |v_{ik} - v_{jk}|$$

gdzie:  $v_{i,k}$  – jest wartością odpowiedniej  $i$ -tej cechy (obiektu) i  $k$ -tej współrzędnej.

Normalizacja tej modyfikacji:

$$\text{dist}_{\mathbf{Hn}}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = \frac{1}{\dim(\underline{\mathbf{v}})} \sum_{k=1}^{\dim(\underline{\mathbf{v}})} \frac{|v_{ik} - v_{jk}|}{\max_m(v_{mk}) - \min_m(v_{mk})}$$

# Cechy jakościowe – przejście na cechy ilościowe

1. Metryka nominalna (*overlap metric*):

$$\text{dist}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = \begin{cases} 0 & \Leftrightarrow \underline{v}_i = \underline{v}_j \\ \infty & \Leftrightarrow \underline{v}_i \neq \underline{v}_j \end{cases}$$

2. Rangi dla wartości jakościowych:

$$\text{dist}_{\underline{\mathbf{v}}_i \underline{\mathbf{v}}_j} = \frac{1}{\text{dim}(\underline{\mathbf{v}})} \sum_{k=1}^{\text{dim}(\underline{\mathbf{v}})} |\text{rng}(v_{ik}) - \text{rng}(v_{jk})|$$

# Cechy mieszane

- Zamiana wartości cech jakościowych na ilościowe i dopiero wtedy wyznaczenie podobieństwa
- Wyznaczenie osobnych podobieństw dla poszczególnych cech:

$$simi(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = \alpha simi_{quan}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) + (1 - \alpha) simi_{qual}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j)$$

$$simi(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j) = simi_{quan}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j)^\alpha \cdot simi_{qual}(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}_j)^{(1-\alpha)}$$

# Porównanie obiektów ze zbiorami obiektów

- Podobieństwo średnie:

$$simi(\underline{\mathbf{v}}_i, \{\underline{\mathbf{v}}_j\}) = \frac{1}{|\{\underline{\mathbf{v}}_j\}|} \sum_{\underline{\mathbf{v}} \in \{\underline{\mathbf{v}}_j\}} simi(\underline{\mathbf{v}}_i, \underline{\mathbf{v}})$$

- Podobieństwo między danym elementem a reprezentantem zbioru:

$$simi(\underline{\mathbf{v}}_i, \{\underline{\mathbf{v}}_j\}) = simi\left(\underline{\mathbf{v}}_i, \frac{1}{|\{\underline{\mathbf{v}}_j\}|} \sum_{\underline{\mathbf{v}} \in \{\underline{\mathbf{v}}_j\}} \underline{\mathbf{v}}\right)$$

$$simi(\underline{\mathbf{v}}_i, \{\underline{\mathbf{v}}_j\}) = \max_{\underline{\mathbf{v}} \in \{\underline{\mathbf{v}}_j\}} (simi(\underline{\mathbf{v}}_i, \underline{\mathbf{v}}))$$