



UNIwersytet PRzyrodniczy we Wroclawiu

## Data Mining Wykład 8

Analiza skupień (grupowanie)  
Grupowanie hierarchiczne O-Cluster

---

---

---

---

---

---

---

---

### Plan wykładu

- Wprowadzanie
- Definicja problemu
- Klasyfikacja metod grupowania
- Grupowanie hierarchiczne

UNIwersytet PRzyrodniczy we Wroclawiu

---

---

---

---

---

---

---

---

### Sformułowanie problemu

Dany jest zbiór obiektów (rekordów). Znajdź naturalne pogrupowanie obiektów w klasy (klastry, skupienia) obiektów o podobnych cechach

- **Grupowanie:**

proces grupowania obiektów, rzeczywistych bądź abstrakcyjnych, w klasy, nazywane klastrami lub skupieniami, o podobnych cechach

UNIwersytet PRzyrodniczy we Wroclawiu

---

---

---

---

---

---

---

---

## Czym jest klaster?

- **Istnieje wiele definicji:**

Zbiór obiektów, które są "podobne"

Zbiór obiektów, takich, że odległość pomiędzy dwoma dowolnymi obiektami należącymi do klastra jest mniejsza aniżeli odległość pomiędzy dowolnym obiektem należącym do klastra i dowolnym obiektem nie należącym do tego klastra

Spójny obszar przestrzeni wielowymiarowej, charakteryzujący się dużą gęstością występowania obiektów

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

---

---

## Przykłady (1)

- **Zbiór dokumentów:**

Zbiór punktów w przestrzeni wielowymiarowej, w której pojedynczy wymiar odpowiada jednemu słowu z określonego słownika

Współrzędne dokumentu w przestrzeni są zdefiniowane względną częstością występowania słów ze słownika.

Klastry dokumentów odpowiadają grupom dokumentów dotyczących podobnej tematyki

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

---

---

## Przykłady (2)

- **Zbiór sekwencji stron WWW:**

Pojedyncza sekwencja opisuje sekwencję dostępow do stron WWW danego serwera realizowaną w ramach jednej sesji przez użytkownika

Klastry sekwencji odpowiadają grupom użytkowników danego serwera, którzy realizowali dostęp do tego serwera w podobny sposób

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

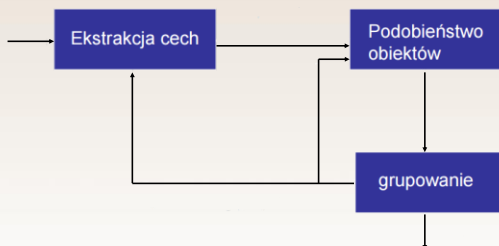
---

---

---

---

### Składowe procesu grupowania (1)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

### Składowe procesu grupowania (2)

- Proces grupowania:

Reprezentacja obiektów (zawiera ekstrakcję/selekcję cech obiektów)

Definicja miary podobieństwa pomiędzy obiektami (zależy od dziedziny zastosowań)

Grupowanie obiektów (klastry)

Znajdowanie charakterystyki klastrów

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

### Miary odległości (1)

- Dyskusja dotycząca podobieństwa, lub odległości, dwóch obiektów wymaga przyjęcia miary odległości pomiędzy dwoma obiektami  $x$  i  $y$  reprezentowanymi przez punkty w przestrzeni wielowymiarowej
- Klasyczne aksjomaty dla miary odległości

$$1. D(x, y) = 0 \Leftrightarrow x = y$$

$$2. D(x, y) = D(y, x)$$

$$3. D(x, y) \leq D(x, z) + D(z, y) \text{ (nierówność trójkąta)}$$

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

## Miary odległości (2)

- Dana jest k-wymiarowa przestrzeń euklidesowa, odległość pomiędzy dwoma punktami  $x=[x_1, x_2, \dots, x_k]$  oraz  $y=[y_1, y_2, \dots, y_k]$  można zdefiniować:

Odległość euklidesowa: 
$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

Odległość Manhattan: 
$$\sum_{i=1}^k |x_i - y_i|$$

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

## Miary odległości (3)

Odległość max z wymiarów: 
$$\max_{i=1}^k |x_i - y_i|$$

Odległość Minkowskiego: 
$$\left(\sum_{i=1}^k |x_i - y_i|^q\right)^{1/q}$$

W przypadku, gdy obiekty nie poddają się transformacji do przestrzeni euklidesowej, proces grupowania wymaga zdefiniowania innych miar odległości (podobieństwa): sekwencja dostępów do stron WWW, sekwencje DNA, sekwencje zbiorów, zbiory atrybutów kategorycznych, dokumenty tekstowe, XML, grafy, itp..

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

## Zmienne binarne (1)

- W jaki sposób obliczyć podobieństwo (lub niepodobieństwo) pomiędzy dwoma obiektami opisanymi zmiennymi binarnymi:
- Podejście: konstruujemy macierz niepodobieństwa

	obiekt j			
	1	0	Sum	
obiekt i	1	q	r	q+r
	0	s	t	s+t
Sum	q+s	r+t	p	

- q – liczba zmiennych przyjmujących wartość 1 dla obu obiektów
- r – ... 1 dla obiektu i, i wartość 0 dla j
- s – ... 0 dla obiektu i, i wartość 1 dla j
- t – ... 0 dla obu obiektów

$p=q+r+s+t$  – łączna liczba zmiennych

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

## Zmienne binarne (2)

- Zmienne binarne symetryczne:

Zmienną binarną nazywamy symetryczną jeżeli obie wartości tej zmiennej posiadają tą samą wagę (np. płeć)

Niepodobieństwo pomiędzy obiektami  $i$  oraz  $j$  jest zdefiniowane następująco:

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

---

---

## Zmienne binarne (3)

- Zmienne binarne asymetryczne:

zmienną binarną nazywamy asymetryczną jeżeli obie wartości tej zmiennej posiadają różne wagi (np. wynik badania EKG)

Niepodobieństwo pomiędzy obiektami  $i$  oraz  $j$  jest zdefiniowane następująco:

$$d(i, j) = \frac{r + s}{q + r + s}$$

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

---

---

## Zmienne binarne (4)

- Dana jest tablica zawierająca informacje o pacjentach

imię	ból	gorączka	katar	test1	test2	test3	test4
Jack	Y	Y	N	P	N	N	N
Mary	N	Y	N	P	N	P	N
Jim	Y	Y	Y	N	N	N	N
...	...	...	...	...	...	...	...

$$d_{\text{sym}}(\text{jack}, \text{mary}) = \frac{2}{4} = 0.5 \quad d_{\text{sym}}(\text{jack}, \text{mary}) = \frac{2}{7} = 0.29$$

$$d_{\text{sym}}(\text{jack}, \text{jim}) = \frac{2}{4} = 0.5 \quad d_{\text{sym}}(\text{jack}, \text{jim}) = \frac{2}{7} = 0.29$$

$$d_{\text{sym}}(\text{jim}, \text{mary}) = \frac{4}{5} = 0.8 \quad d_{\text{sym}}(\text{jim}, \text{mary}) = \frac{4}{7} = 0.57$$

UNIWERSYTET PRZYRODNICZY WE WROCŁAWIU

---

---

---

---

---

---

---

---

---

---

## Zmienne kategoryczne

- Zmienna kategoryczna jest generalizacją zmiennej binarnej: może przyjmować więcej niż dwie wartości (np. dochód: wysoki, średni, niski)
- Niepodobieństwo (podobieństwo) pomiędzy obiektami  $i, j$ , opisanymi zmiennymi kategorycznymi, można zdefiniować następująco:

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

- gdzie  $p$  oznacza łączną liczbę zmiennych,  $m$  oznacza liczbę zmiennych, których wartość jest identyczna dla obu obiektów,  $n$  oznacza liczbę zmiennych, których wartość jest różna dla obu obiektów

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

## Metody grupowania - Typy metod

- Istnieje wiele różnych metod i algorytmów grupowania:
  - Dla danych liczbowych i/lub danych symbolicznych
  - Deterministyczne i probabilistyczne
  - Rozłączne i przecinające się
  - Hierarchiczne i płaskie
  - Monoteiczny i politeiczny
  - Przyrostowe i nieprzyrostowe

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

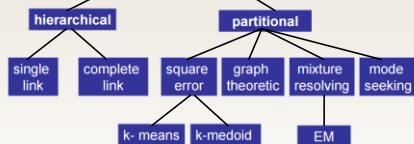
---

---

---

## Metody grupowania (1)

### Klasyfikacja metod



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

## Metody grupowania (4)

- Dwa podstawowe podejścia do procesu grupowania obiektów:
- **Metody hierarchiczne**

generują zagnieżdżoną sekwencję podziałów zbiorów obiektów w procesie grupowania

- **Metody z iteracyjno- optymalizacyjne**

generują tylko jeden podział (partycję) zbioru obiektów w dowolnym momencie procesu grupowania

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

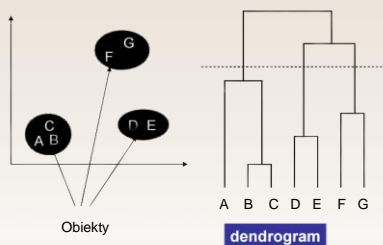
---

---

---

---

## Metody grupowania hierarchicznego (1)



Metoda grupowania hierarchicznego polega na sekwencyjnym grupowaniu obiektów – drzewo klastrow (tzw. dendrogram)

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

## Metody grupowania hierarchicznego (2)

- **Podejście podziałowe**

(top-down): początkowo, wszystkie obiekty przypisujemy do jednego klastra; następnie, w kolejnych iteracjach, klaster jest dzielony na mniejsze klastry, które, z kolei, dzielone są na kolejne mniejsze klastry

- **Podejście aglomeracyjne**

(bottom-up): początkowo, każdy obiekt stanowi osobny klaster, następnie, w kolejnych iteracjach, klastry są łączone w większe klastry

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

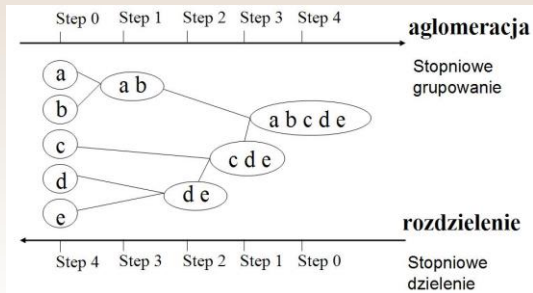
---

---

---

---

## Metody grupowania hierarchicznego (2)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

## Miary odległości (1)

- W obu podejściach, aglomeracyjnym i podziałowym, liczba klastrów jest ustalona z góry przez użytkownika i stanowi warunek stopu procesu grupowania:

4 podstawowe (najczęściej stosowane) miary odległości pomiędzy klastrami są zdefiniowane następująco,

- gdzie  $|p - p'|$  oznacza odległość pomiędzy dwoma obiektami (lub punktami),
- $p$  i  $p'$  mi oznacza średnią wartość klastra  $C_i$ ,
- $n_i$  oznacza liczbę obiektów należących do klastra  $C_i$

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

## Miary odległości (2)

Minimalna odległość:

$$d_{\min}(C_i, C_j) = \min_{p \in C_i, p' \in C_j} \|p - p'\|$$

Maksymalna odległość:

$$d_{\max}(C_i, C_j) = \max_{p \in C_i, p' \in C_j} \|p - p'\|$$

Odległość średnich:

$$d_{\text{mean}}(C_i, C_j) = \|m_i - m_j\|$$

Średnia odległość:

$$d_{\text{ave}}(C_i, C_j) = 1/(n_i n_j) \sum_{p \in C_i} \sum_{p' \in C_j} \|p - p'\|$$

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---



### Ogólny hierarchiczny aglomeracyjny algorytm grupowania

Dane wejściowe: baza danych D obiektów (n- obiektów)

Dane wyjściowe: dendrogram reprezentujący grupowanie obiektów

- 1) umieść każdy obiekt w osobnym klastrze;
- 2) skonstruuj macierz odległości pomiędzy klastrami;
- 3) Dla każdej wartości niepodobieństwa  $dk$  (dk może się zmieniać w kolejnych iteracjach) powtarzaj:
  - Utwórz graf klastrow, w którym każda para klastrow, której wzajemna odległość jest mniejsza niż  $dk$ , jest połączona lukiem aż wszystkie klastry utworzą graf spójny

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

---

---

### Hierarchiczny aglomeracyjny algorytm grupowania

- Umieść każdy obiekt w osobnym klastrze. Skonstruuj macierz przyległości zawierającą odległości pomiędzy każdą parą klastrow
- Korzystając z macierzy przyległości znajdź najbliższą parę klastrow. Połącz znalezione klastry tworząc nowy klastr. Uaktualnij macierz przyległości po operacji połączenia
- Jeżeli wszystkie obiekty należą do jednego klastra, zakończ procedurę grupowania, w przeciwnym razie przejdź do kroku 2

UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

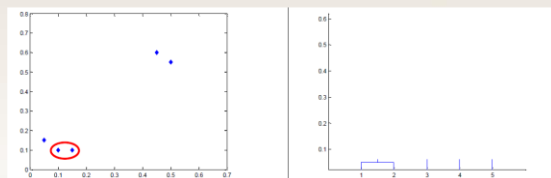
---

---

---

---

### Przykład (1)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

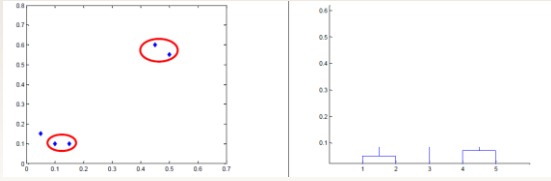
---

---

---

---

Przykład (2)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

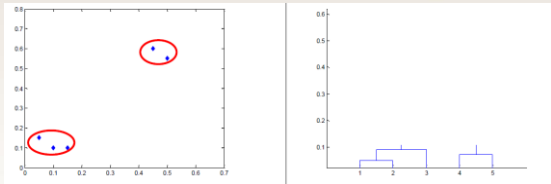
---

---

---

---

Przykład (3)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

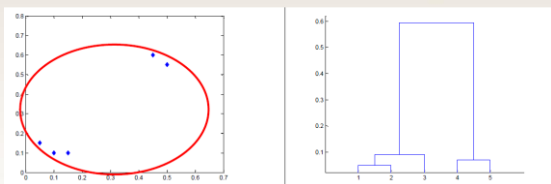
---

---

---

---

Przykład (4)



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

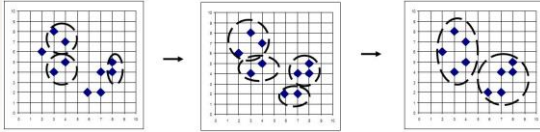
---

---

## Aglomeracyjny algorytm grupowania danych

**Algorithm 2:** Aglomeracyjny algorytm grupowania danych

**węzicie:**  $P$  - zbiór przykładów  $P$ ;  
 $D(C_i, C_j)$  - funkcja do mierzenia odległości między dwoma skupieniami  $C_i$  i  $C_j$ ;  
**wyjście:** Dendrogram skupień.  
**begin**  
 while pozostało więcej niż jedno skupienie do  
   Niech  $C_i$  i  $C_j$  będą skupieniami minimalizującymi odległość  $D(C_i, C_j)$  między dwoma skupieniami;  
    $C_k = C_i \cup C_j$ ;  
   usuń skupienie  $C_j$ ;  
**end**



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

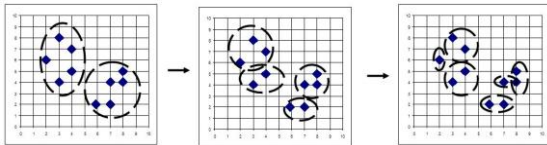
---

---

## Rozdzielający algorytm grupowania danych

**Algorithm 3:** Rozdzielający algorytm grupowania danych

**węzicie:**  $P$  - zbiór przykładów  $P$ ;  
 $D(C_i, C_j)$  - funkcja do mierzenia odległości między dwoma skupieniami  $C_i$  i  $C_j$ ;  
**wyjście:** Dendrogram skupień.  
**begin**  
 while istnieje skupienie składające z więcej niż dwóch elementów do  
   Niech  $C_i$  będzie skupieniem, którego po dzieleniu tworzą się dwa nowe skupienia  $C_j$  i  $C_k$  o maksymalnej odległości  $D(C_j, C_k)$ ;  
   Podziel  $C_i$  na dwa skupienia  $C_j$  i  $C_k$ ;  
   Usuń skupienie  $C_i$ ;  
**end**



UNIWERSYTET PRZYRODNICZY WE WROCLAWIU

---

---

---

---

---

---

---

---

---

---

---

---