
Analiza Skupień - Grupowanie

Zaawansowana Eksploracja Danych



JERZY STEFANOWSKI
Inst. Informatyki PP
Wersja dla TPD 2013

Część II

Organizacja wykładu

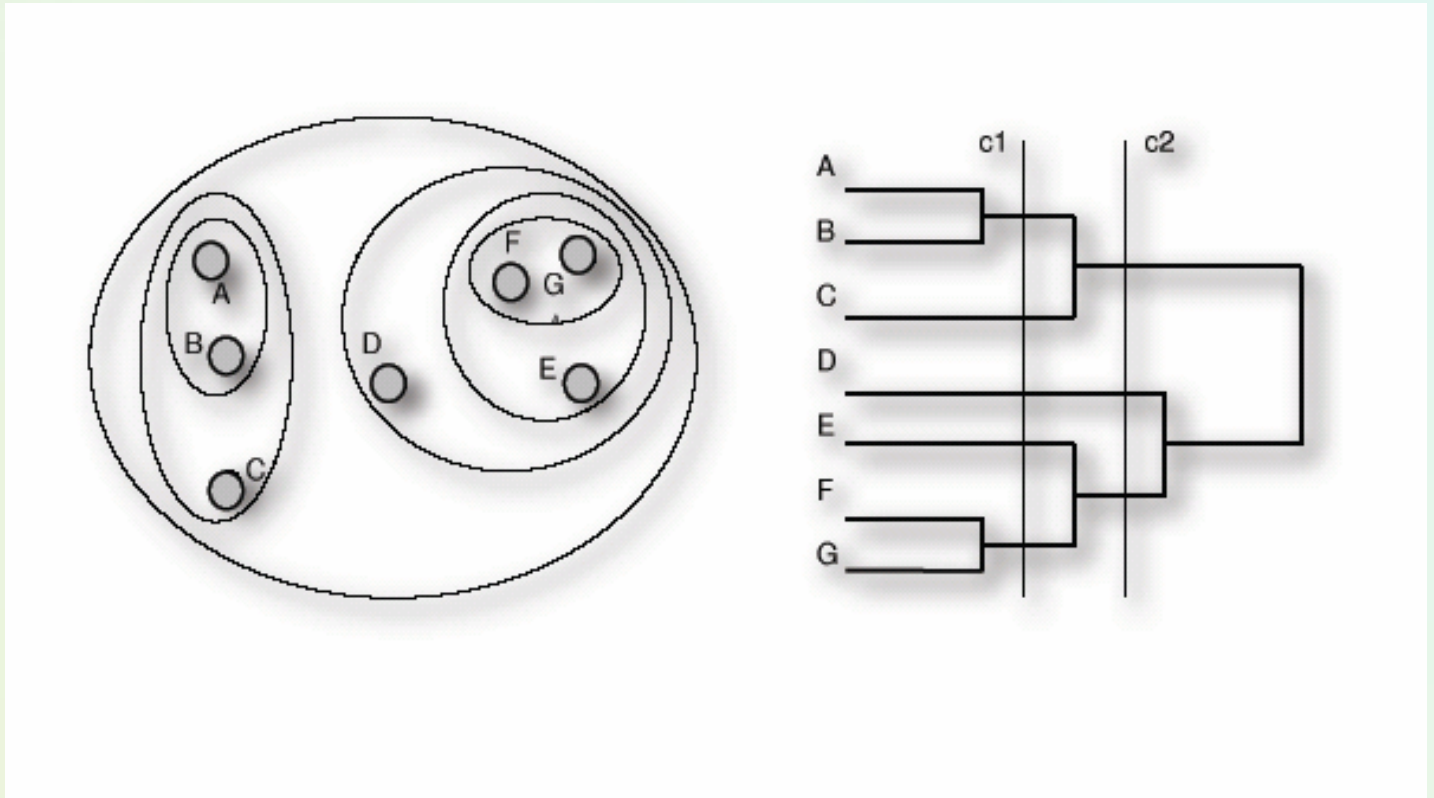
- Przypomnienie wyboru liczby skupień
- Studium przypadku użycia
- Rozszerzenia dla analizy danych o większych rozmiarach.
- Podsumowanie



Referencje do literatury (przykładowe)

- Koronacki J. Statystyczne systemy uczące się, WNT 2005.
- Pocięcha J., Podolec B., Sokołowski A., Zając K. „Metody taksonomiczne w badaniach społeczno-ekonomicznych”. PWN, Warszawa 1988,
- Stapor K. „Automatyczna klasyfikacja obiektów” Akademska Oficyna Wydawnicza EXIT, Warszawa 2005.
- Hand, Mannila, Smyth, „Eksploracja danych”, WNT 2005.
- Larose D: „Odkrywania wiedzy z danych”, PWN 2006.
- T.Morzy: „Eksploracja danych”, PWN 2013.
- Kucharczyk J. „Algorytmy analizy skupień w języku ALGOL 60” PWN Warszawa, 1982,
- Materiały szkoleniowe firmy Statsoft.
- Inne polecane przez wykładownicę

Jak wybrać liczbę skupień?



Analiza skupień w Statsoft -Statistica

Analiza Skupień – Statistica; więcej na www.statsoft.com. Przykład analizy danych o parametrach samochodów

STATISTICA: Analiza skupień - [Dane: CARS.STA 5v * 22c]

Plik Edycja Widok Analiza Wykresy Opcje Okno Pomoc

-.521072362755425

Zmienne Przypadki

LICZBOWE WARTOŚCI

Cena, wydajność, trzymanie się drogi różny

	1	2	3	4	5
	CENA	PRZYSP	HAMOWAN	WSK_TRZY	ZUŻYCIE
Acura	-,521	,477	-,007	,382	2,079
Audi	,866	,208	,319	-,091	-,677
BMW	,496	-,802	,192	-,091	-,154
Buick	-,614	1,689	,933	-,210	-,154
Corvette	1,235	-1,811	-,494	,973	-,677
Chrysler	-,614	,073	,427	-,210	-,154
Dodge	-,706	-,196	,481	,145	-,154
Eagle	-,614	1,218	-4,199	-,210	-,677
Ford	-,706	-1,542	,987	,145	-1,724
Honda	-,429	,410	-,007	,027	,369
Isuzu	-,798	,410	-,061	-4,230	1,067
Mazda	,126	,679	-,133	,500	-1,724
Mercedes	1,051	,006	,120	-,091	-,154
Mitsub.	-,614	-1,003	,084	,382	,718
Nissan	-,429	,073	-,007	,263	,997
Olds	-,614	-,734	,409	,382	2,114
Pontiac	-,614	,679	,536	,145	,195
Porsche	3,454	-2,215	-,296	,618	-1,026
Saab	,588	,679	,246	,263	,021
Toyota	-,059	1,218	,228	,736	-,851
VW	-,706	-,128	,102	,382	,195
Volvo	,219	,612	,138	-,210	,369

Metoda grupowania

Aglomeracja

Grupowanie metodą k-średnich

Grupowanie obiektów i cech

OK

Anuluj

Otwórz dane

SELECT CASES

W

Analiza skupień

Analiza Wykresy Opcje Okno Pomoc

Kolumny Wiersze

S.STA 5v * 22c

Cena, wydajność, trzymanie się drogi różny

1	2	3	4	5
CENA	PRZYSP	HAMOWAN	WSK TRZY	ZUŻYCIE
-521	,477	-,007	,382	2,079
,866	,208	,319	-,091	-,677
,496	-,802	,192	-,091	-,154

Przebieg aglomeracji (cars.sta)

Dalej... Pojedyncze wiązanie
Odległości euklidesowe

połącz. odległ.	Obj. Nr 1	Obj. Nr 2	Obj. Nr 3	Obj. Nr 4	Obj. Nr 5	Obj. Nr 6	Obj. Nr 7	Obj. Nr 8	Obj. Nr 9
4580484	Chrysler	Dodge							
5710964	Chrysler	Dodge	VW						
6231085	Audi	Mercedes							
6670490	Honda	Pontiac							
7060042	Saab	Volvo							
7313396	Chrysler	Dodge	VW	Honda	Pontiac				
7323840	Chrysler	Dodge	VW	Honda	Pontiac	Saab	Volvo		
7506309	Chrysler	Dodge	VW	Honda	Pontiac	Saab	Volvo	Niss	
9159300	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pontiac	S	
9824548	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pontiac	S	
1,023831	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pontiac	S	
1,127473	Mazda	Toyota							
1,164055	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
1,193655	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
1,284603	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
1,301269	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
1,855838	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
2,128886	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
2,317976	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
4,214866	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	
4,355048	Acura	Audi	Mercedes	Chrysler	Dodge	VW	Honda	Pont	

Analiza skupień: Aglomeracja

Zmienne: **WSZYSTKIE** OK

Wejście: **Dane surowe**

Grupowanie: **Przypadki (obiekty)**

Metoda aglomeracji (wiązania): **Pojedynczego wiązania**

Miara odległości: **Odległość euklidesowa**

D: 2 I: 2

Braki danych: **Usuwane przypadkami**

Przetwarzanie wsadowe i drukowanie

SELECT CASES W

Wyniki aglomeracji

Liczba zmiennych: 5

Liczba przyp.: 22

Łączenie przyp.

Braki danych były usuwane przypad.

Metoda aglomeracji: **Pojedyncze wiązanie**

Miara odległości: **Odległości euklidesowe standaryzowane**

Poziomy hierarchiczny wykres drzewkowy OK

Pionowy wykres soplekowy Anuluj

Prostokątne gałęzie

Skaluj drzewo do odl_wiąz./odl_maks*100

Przebieg aglomeracji

Wykres przebiegu aglomeracji

Macierz odległości

Statystyki opisowe

Zapisz macierz odległości

Wrażliwość na dobór metod łączenia skupień

Diagram dla 22 przyp.
Pojedyncze wiązanie
Odległości euklidesowe

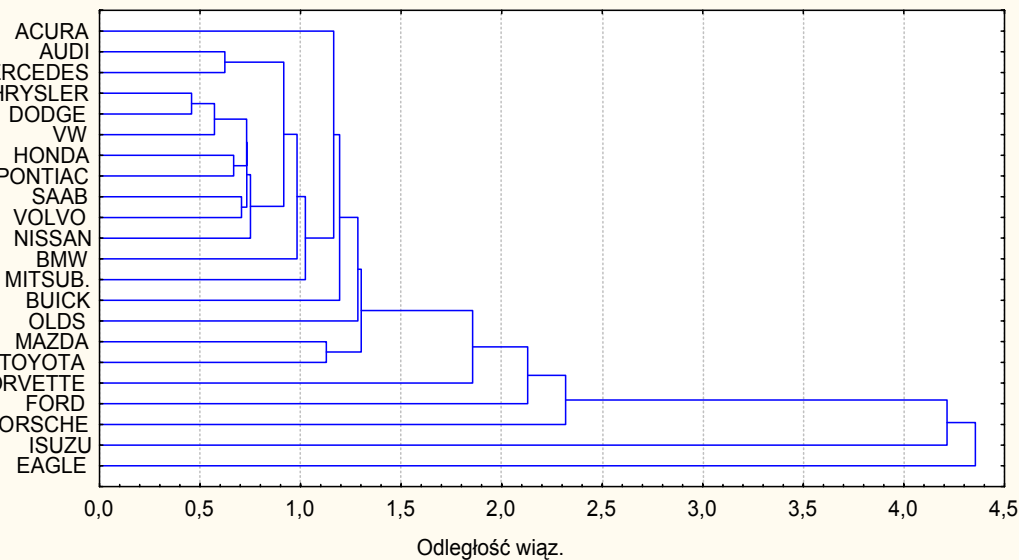
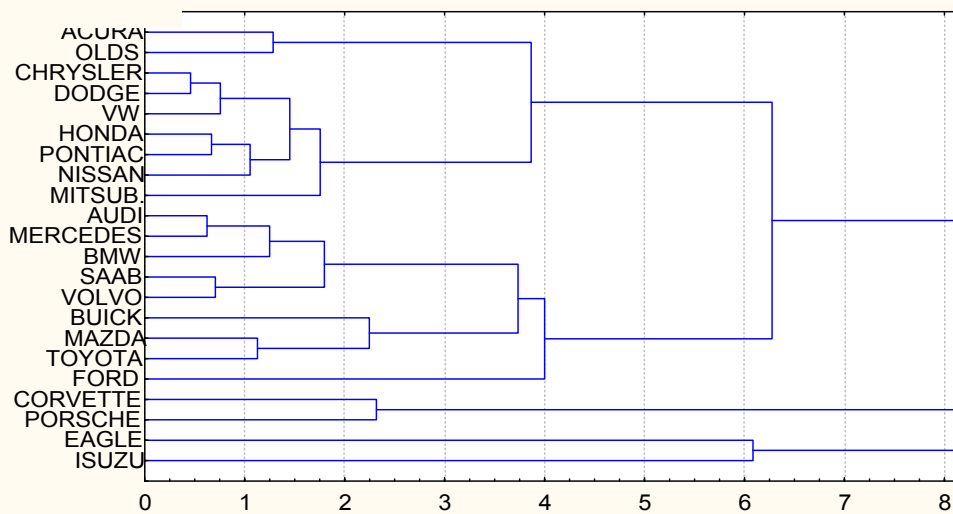
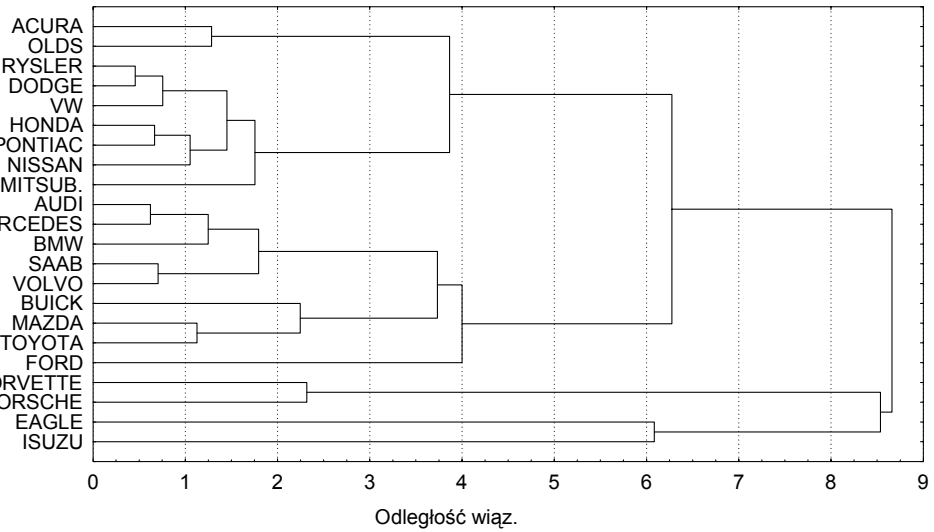


Diagram dla 22 przyp.
Metoda Warda
Odległości euklidesowe



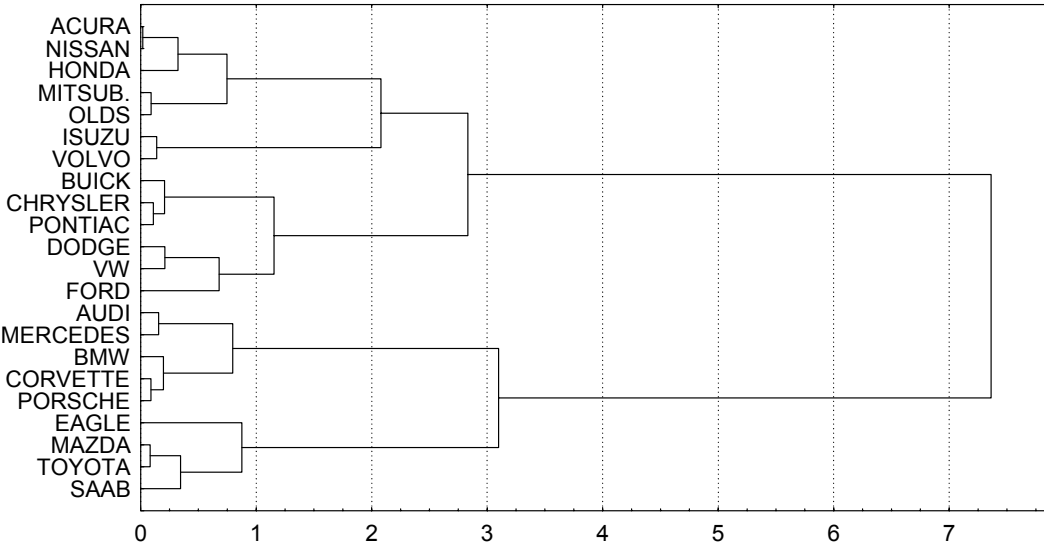
Przykłady użycia metody Warda

Diagram dla 22 przyp.
Metoda Warda
Odległości euklidesowe



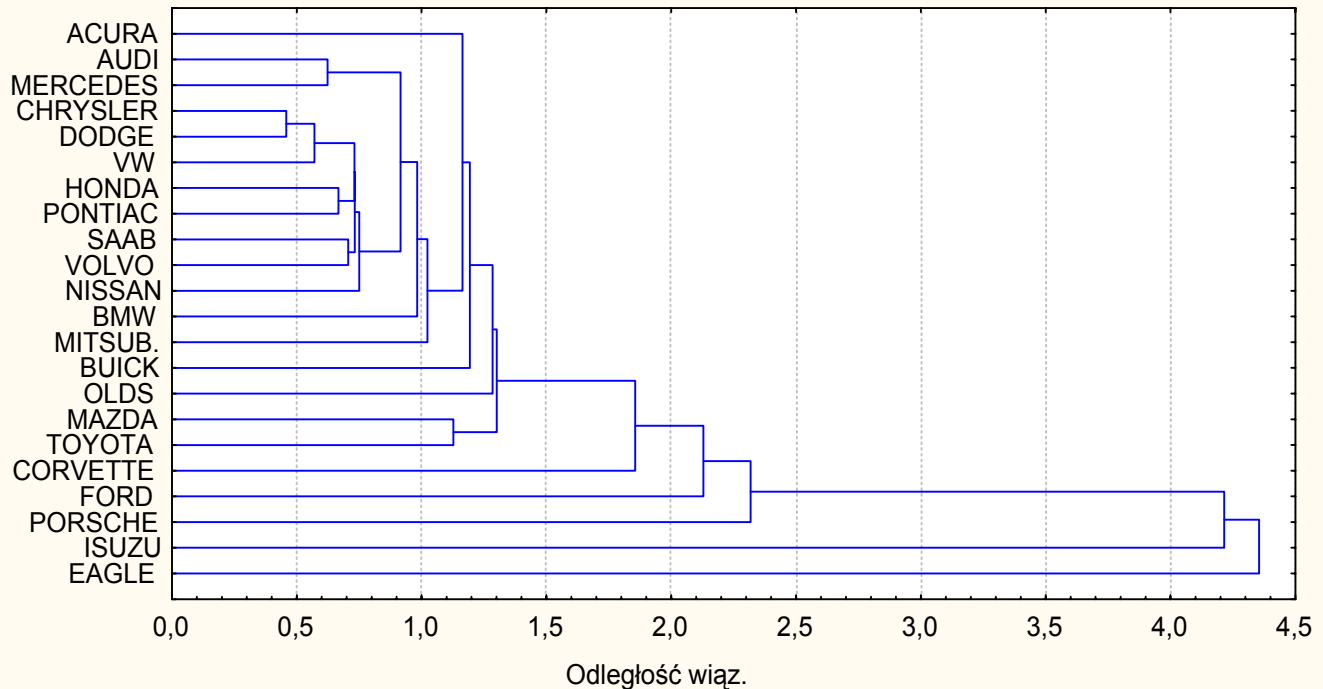
Cars data

Diagram dla 22 przyp.
Metoda Warda
1-r Pearsona



Dendrogram for „Single Linkage”

Diagram dla 22 przyp.
Pojedyncze wiązanie
Odległości euklidesowe



Opis tworzenia dendrogramu

- Łączenie obiektów w kolejnych krokach

STATISTICA: Analiza skupień - [Przebieg aglomeracji [cars.staj]]

Plik Edycja Widok Analiza Wykresy Opcje Okno Pomoc

Dodge Kolumny Wiersze

Pojedyncze wiązanie
Odlegności euklidesowe

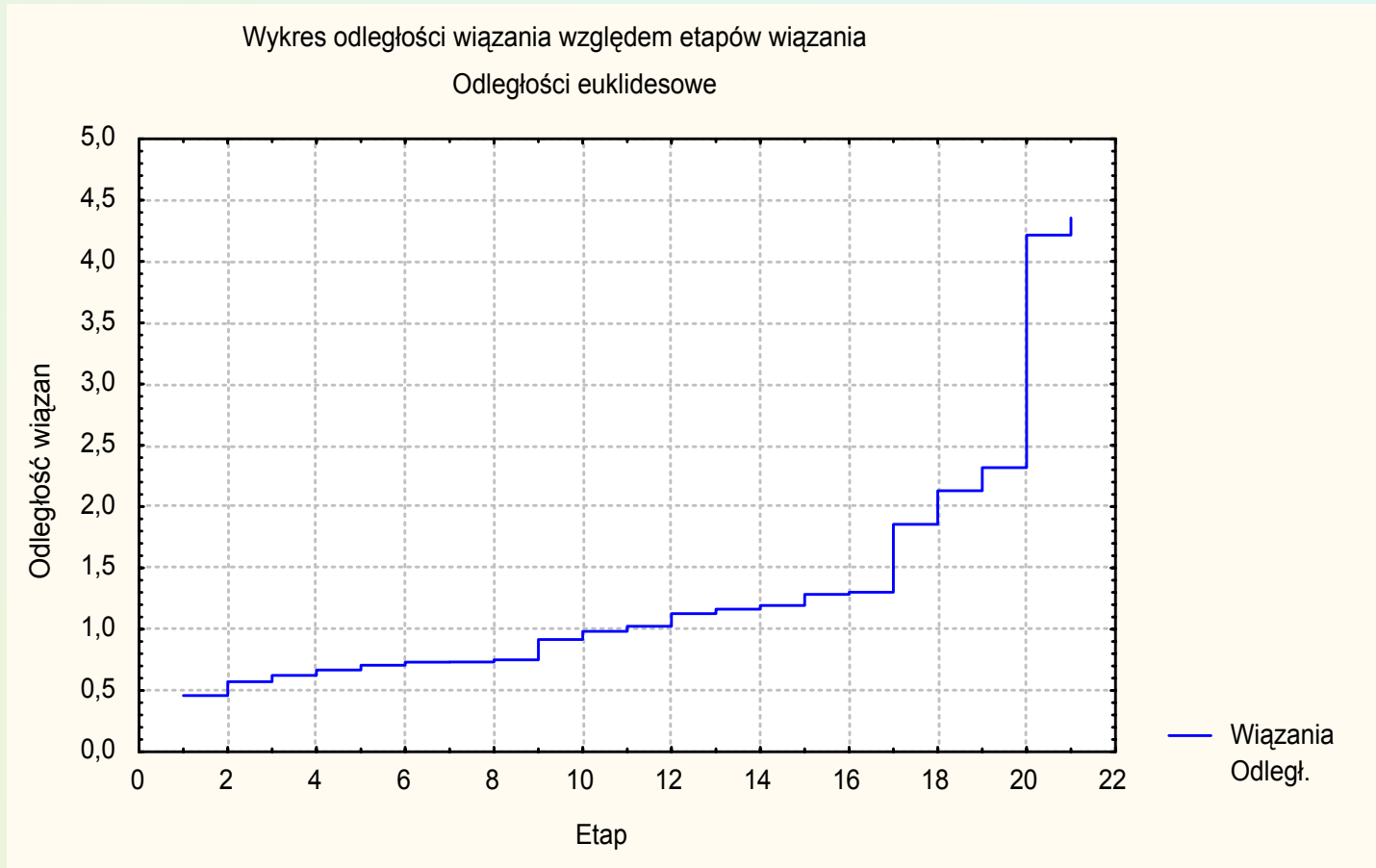
po ³ cz. odleg ³ .	Obj. Nr 1	Obj. Nr 2	Obj. Nr 3	Obj. Nr 4	Obj. Nr 5
,4580484	Chrysler	Dodge			
,5710964	Chrysler	Dodge	VW		
,6231085	Audi	Mercedes			
,6670490	Honda	Pontiac			
,7060042	Saab	Volvo			
,7313396	Chrysler	Dodge	VW	Honda	Por
,7323840	Chrysler	Dodge	VW	Honda	Por
,7506309	Chrysler	Dodge	VW	Honda	Por
,9159300	Audi	Mercedes	Chrysler	Dodge	
,9824548	Audi	Mercedes	Chrysler	Dodge	
1,023831	Audi	Mercedes	Chrysler	Dodge	
1,127473	Mazda	Toyota			
1,164055	Acura	Audi	Mercedes	Chrysler	Dc
1,193655	Acura	Audi	Mercedes	Chrysler	Dc
1,284603	Acura	Audi	Mercedes	Chrysler	Dc
1,301269	Acura	Audi	Mercedes	Chrysler	Dc
1,855838	Acura	Audi	Mercedes	Chrysler	Dc
2,128886	Acura	Audi	Mercedes	Chrysler	Dc
2,317976	Acura	Audi	Mercedes	Chrysler	Dc
4,	Acura	Audi	Mercedes	Chrysler	Dc
4,	Acura	Audi	Mercedes	Chrysler	Dc

Gotowy Wyjście: WYŁĄCZONE Sel: NIE Waga: WYŁĄCZONA

Start Windows Commander 4.0... STATISTICA: Analiza... Document3 - Microsoft W... 16:40

Analiza procesu łączenia

- Wykres kolankowy – a cut point („kolanko” / knee)



czba zmiennych: 5
czba przyp.: 22
ązanie przypadków met.k-ś
aki danych usuwano przypadkami
czba skupień: 4
związanie odnaleziono po 1 iteracjach

Analiza wariacji [Anuluj]

Średnie skupień i odległości euklidesowe

Wykres średnich

Statystyki opisowe każdego skupienia

Elementy każdego skupienia i odległości

Zapisz klasyfikacje i odległości

Analiza skupień: Grupowanie metodą k-średnich

Zmienne: WSZYSTKIE [OK]

Grupowanie: Przypadki (obiekty) [Anuluj]

Liczba skupień: 4

Liczba iteracji: 10

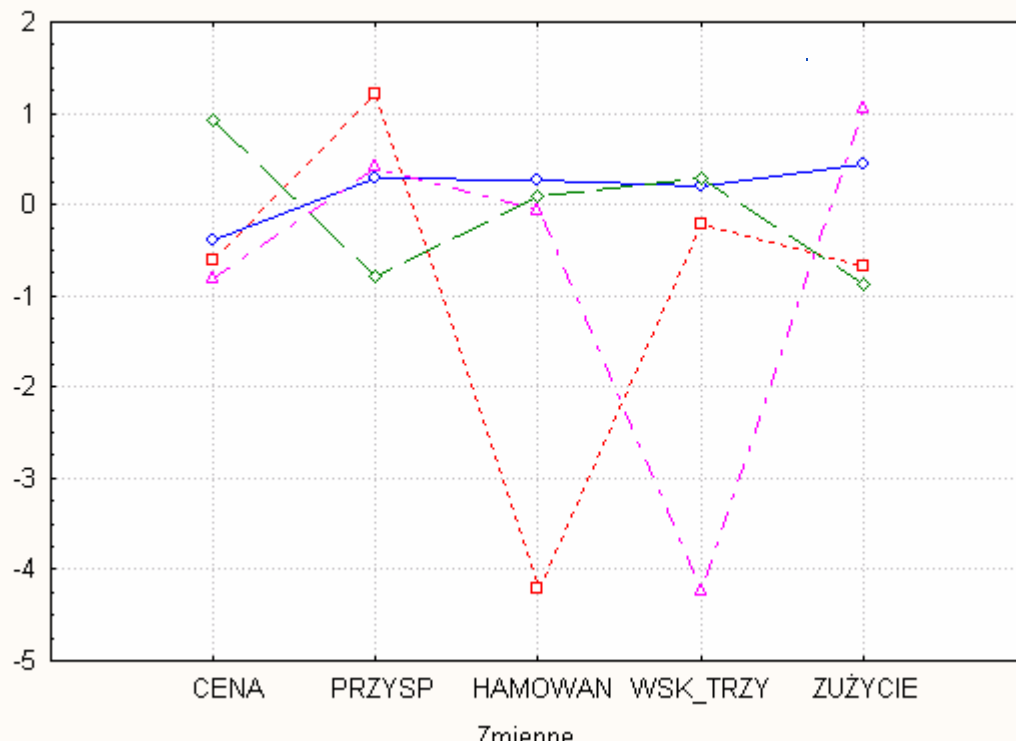
Braki danych: Usuwane przypadkami

Wstępne centra skupień

- Wybierz obserwacje tak, aby zmaksymalizować odległości skupień
- Sortuj odległości i weź obserwacje przy stałym interwale
- Wybierz pierwszych N (liczba skupień) obserwacji

Przetwarzanie wsadowe i drukowanie [SELECT CASES] [10] [W]

Wykres średnich każdego skupienia



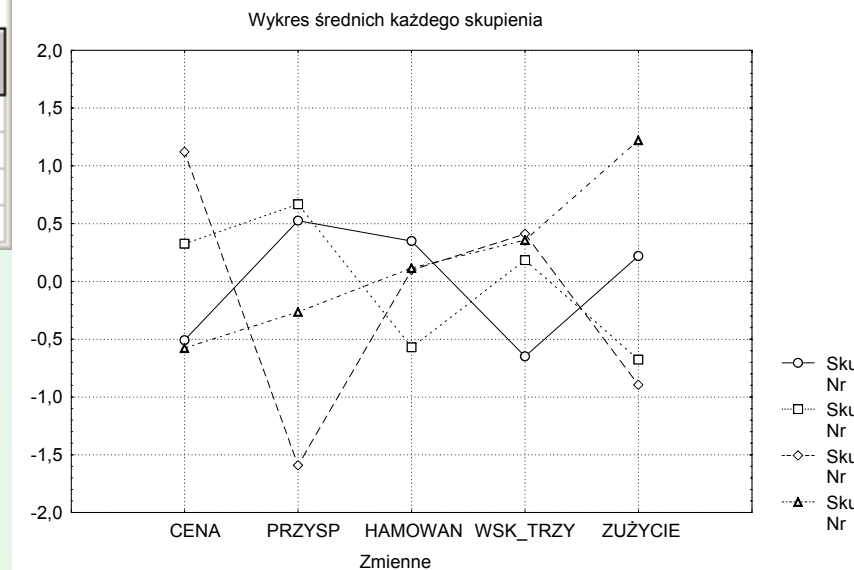
Analiza Skupień – optymalizacja k-średnich

- Skupien. Nr 1
- Skupien. Nr 2
- Skupien. Nr 3
- Skupien. Nr 4

Wsparcie do charakterystyki skupień

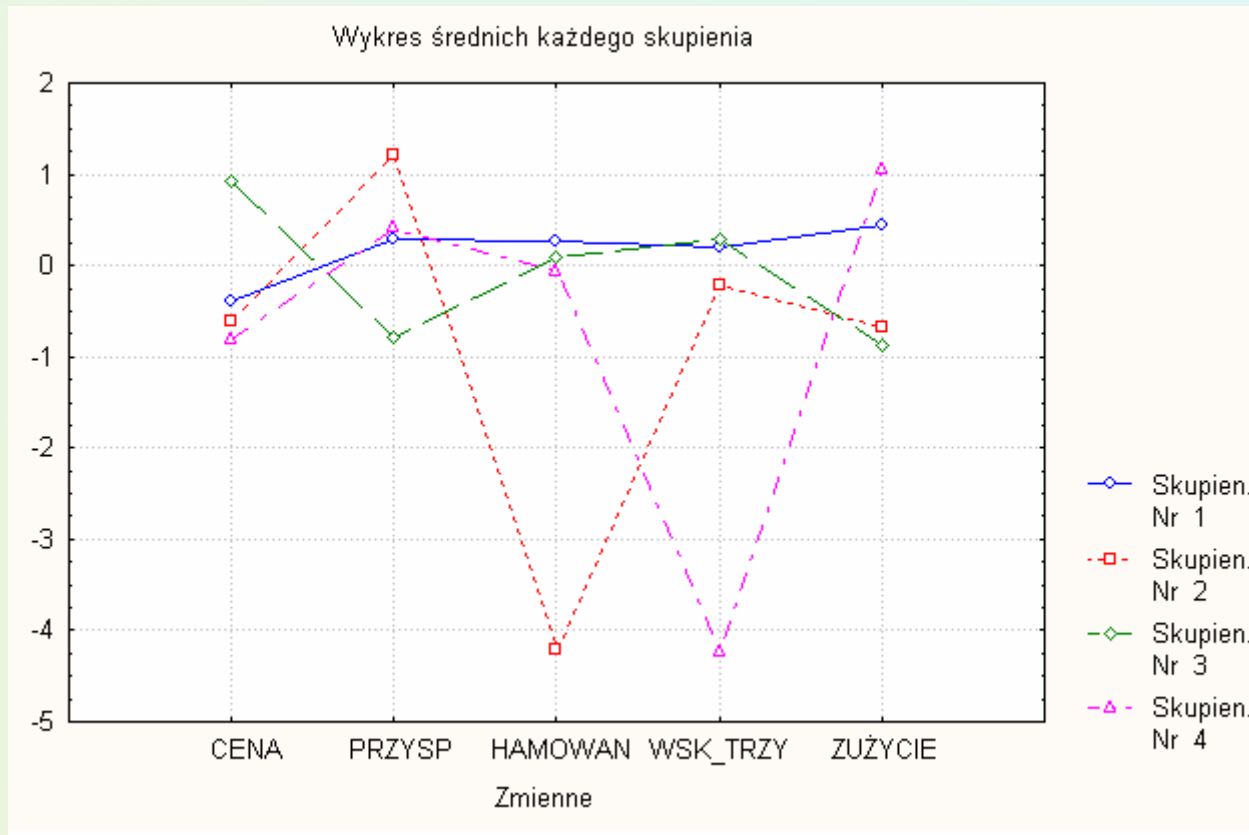
Odległości euklidesowe skupień (cars.sta)				
Dalej...	Odległości pod przekątną Kwadr. odległości nad przekątną			
Skupien. Numer	Nr 1	Nr 2	Nr 3	Nr 4
Nr 1	0,000000	,612157	1,912964	,539154
Nr 2	,782405	0,000000	1,256528	1,156810
Nr 3	1,383100	1,120950	0,000000	1,824839
Nr 4	,734271	1,075551	1,350866	0,000000

Średnie skup. (cars.sta)				
ANALIZA SKUPIEŃ	Skupien. Nr 1	Skupien. Nr 2	Skupien. Nr 3	Skupien. Nr 4
CENA		,326371	1,11989	-,576541
PRZYSP	,525328	,667950	-1,59237	-,263101
HAMOWAN	,349673	-,569764	,09733	,116307
WSK_TRZY	-,648829	,184539	,41118	,357969
ZUŻYCIE	,219949	-,677062	-,89508	1,220614



Elementy skupienia numer 2 (cars.sta)						
ANALIZA SKUPIEŃ	i odległości od środka właściwego skupienia W skupieniu jest 6 przyp					
	Audi	Eagle	Mazda	Mercedes	Saab	Toyota
Odległ.	,523036	1,703612	,533962	,598004	,495550	,533369

Wizualizacja centroidów



Analiza skupień w WEKA

WEKA zakładka Clustering

- Stopniowy przyrost implementacji
 - *k*-Means
 - EM
 - Cobweb
 - X-means
 - FarthestFirst...
 - DbScann
 - Oraz nowe
- Możliwości prostej wizualizacji i ew. porównania przydziałów do „wzorcowej klasyfikacji” – jeśli jest dostępna w pliku arff

Przykład 1. Bank data in WEKA

- Prosty przykład – sample of customers of the bank
 - Bank data (bank-data.csv -> bank.arff)
 - Pre-processing wykonany na wersji csv
 - 600 klientów opisanych przez 11 atrybutó

```
id,age,sex,region,income,married,children,car,save_act,current_act,mortgage,pep
ID12101,48,FEMALE,INNER_CITY,17546.0,NO,1,NO,NO,NO,NO,YES
ID12102,40,MALE,TOWN,30085.1,YES,3,YES,NO,YES,YES,NO
ID12103,51,FEMALE,INNER_CITY,16575.4,YES,0,YES,YES,YES,NO,NO
ID12104,23,FEMALE,TOWN,20375.4,YES,3,NO,NO,YES,NO,NO
ID12105,57,FEMALE,RURAL,50576.3,YES,0,NO,YES,NO,NO,NO
.....
.....
```

- Skorzystaj z k-means

Charakterystyka danych

Weka Explorer [Window Controls]

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Undo | Save...

Filter: Choose **None** [Apply]

Current relation: Relation: bank, Instances: 600, Attributes: 11

Attributes:

No.	Name
1	age
2	sex
3	region
4	income
5	married
6	children
7	car
8	save_act
9	current_act
10	mortgage
11	pep

Selected attribute:

Name: age, Type: Numeric, Missing: 0 (0%), Distinct: 50, Unique: 0 (0%)

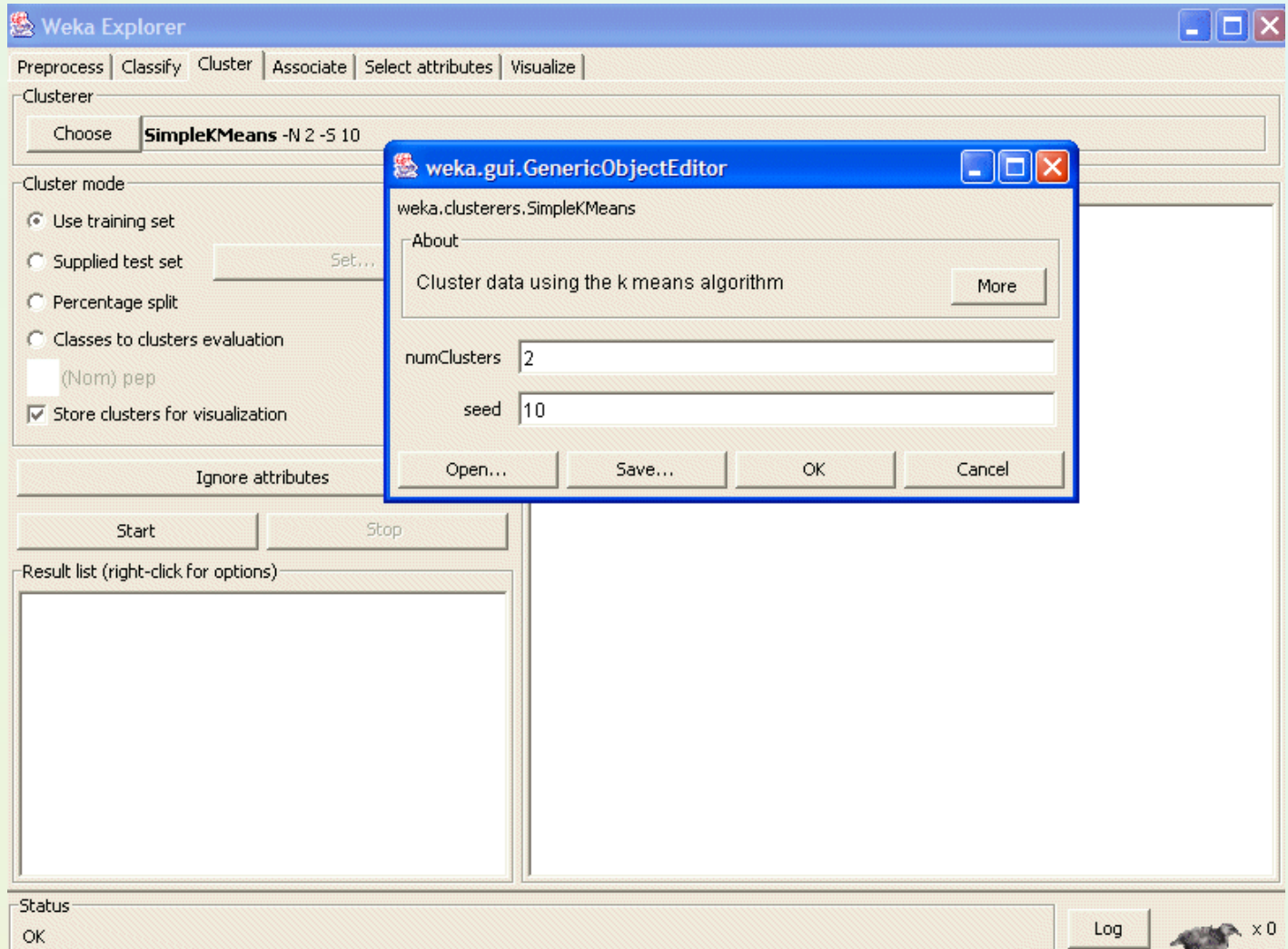
Statistic	Value
Minimum	18
Maximum	67
Mean	42.395
StdDev	14.425

Colour: pep (Nom) [Visualize All]

Status: OK [Log] x 0

Simple k-means

- Dobór parametrów



Clustering results

Clusterer
Choose **SimpleKMeans -N 6 -S 10**

Cluster mode
 Use training set
 Supplied test set (Set...)
 Percentage split (% 66)
 Classes to clusters evaluation (Nom) pep
 Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)
16:47:12 - SimpleKMeans

- View in main window
- View in separate window**
- Save result buffer
- Load model
- Save model
- Re-evaluate model on current test set
- Visualize cluster assignments
- Visualize tree

Clusterer output

Cluster 2
Mean/Mode: 44.0479 MALE INNER_CITY 28547.224 YES
Std Devs: 14.2211 N/A N/A 12696.446

Cluster 3
Mean/Mode: 40.5068 MALE TOWN 25975.293 YES 0 YES
Std Devs: 13.6353 N/A N/A 11111.66

Cluster 4
Mean/Mode: 49.7843 FEMALE INNER_CITY 33917.4538 NO
Std Devs: 13.6872 N/A N/A 14195.168

Cluster 5
Mean/Mode: 41.5234 FEMALE TOWN 26191.8366 YES 0 NO
Std Devs: 13.5728 N/A N/A 11737.313

Clustered Instances

0	66	(11%)
1	85	(14%)
2	146	(24%)
3	73	(12%)
4	102	(17%)
5	128	(21%)

Status: OK Log x 0

- Jak analizować okno z wynikami

Okno wynikowe

- Co można odnaleźć
- Jak opisano centra skupień?

```
16:47:12 - SimpleKMeans
kMeans
-----
Number of iterations: 9

Cluster centroids:

Cluster 0
  Mean/Mode: 36.6061 FEMALE RURAL 23215.9002 NO 3 NO YES YES NO NO
  Std Devs: 14.4317 N/A N/A 12378.3336 N/A N/A N/A N/A N/A N/A
Cluster 1
  Mean/Mode: 38.1176 FEMALE INNER_CITY 24775.7982 YES 1 NO YES YES YES YES
  Std Devs: 13.793 N/A N/A 12444.5713 N/A N/A N/A N/A N/A N/A
Cluster 2
  Mean/Mode: 44.0479 MALE INNER_CITY 28547.224 YES 0 YES YES YES NO NO
  Std Devs: 14.2211 N/A N/A 12696.4468 N/A N/A N/A N/A N/A N/A
Cluster 3
  Mean/Mode: 40.5068 MALE TOWN 25975.293 YES 0 YES NO YES YES YES
  Std Devs: 13.6353 N/A N/A 11111.66 N/A N/A N/A N/A N/A N/A
Cluster 4
  Mean/Mode: 49.7843 FEMALE INNER_CITY 33917.4538 NO 0 YES YES YES NO YES
  Std Devs: 13.6872 N/A N/A 14195.1688 N/A N/A N/A N/A N/A N/A
Cluster 5
  Mean/Mode: 41.5234 FEMALE TOWN 26191.8366 YES 0 NO YES YES NO NO
  Std Devs: 13.5728 N/A N/A 11737.3135 N/A N/A N/A N/A N/A N/A

Clustered Instances
0      66 ( 11%)
1      85 ( 14%)
2     146 ( 24%)
3      73 ( 12%)
4     102 ( 17%)
5     128 ( 21%)
```

Gdzie jest przydział do skupień?

```
TextPad - [D:\Bamshad\CLASS\ECT584\WEKA\Cluster\bank-kmeans.arff]
File Edit Search View Tools Macros Configure Window Help

1 @relation bank_clustered
2
3 @attribute Instance_number numeric
4 @attribute age numeric
5 @attribute sex {FEMALE,MALE}
6 @attribute region {INNER_CITY,TOWN,RURAL,SUBURBAN}
7 @attribute income numeric
8 @attribute married {NO,YES}
9 @attribute children {0,1,2,3}
10 @attribute car {NO,YES}
11 @attribute save_act {NO,YES}
12 @attribute current_act {NO,YES}
13 @attribute mortgage {NO,YES}
14 @attribute pep {YES,NO}
15 @attribute Cluster {cluster0,cluster1,cluster2,cluster3,cluster4,cluster5}
16
17 @data
18 0,48,FEMALE,INNER_CITY,17546,NO,1,NO,NO,NO,NO,YES,cluster1
19 1,40,MALE,TOWN,30085,1,YES,3,YES,NO,YES,YES,NO,cluster3
20 2,51,FEMALE,INNER_CITY,16575,4,YES,0,YES,YES,YES,YES,NO,NO,cluster2
21 3,23,FEMALE,TOWN,20375,4,YES,3,NO,NO,YES,NO,NO,cluster5
22 4,57,FEMALE,RURAL,50576,3,YES,0,NO,YES,NO,NO,NO,cluster5
23 5,57,FEMALE,TOWN,37869,6,YES,2,NO,YES,YES,NO,YES,cluster5
24 6,22,MALE,RURAL,8877,07,NO,0,NO,NO,YES,NO,YES,cluster0
25 7,58,MALE,TOWN,24946,6,YES,0,YES,YES,YES,NO,NO,cluster2
26 8,37,FEMALE,SUBURBAN,25304,3,YES,2,YES,NO,NO,NO,NO,cluster5
27 9,54,MALE,TOWN,24212,1,YES,2,YES,YES,YES,NO,NO,cluster2
28 10,66,FEMALE,TOWN,59803,9,YES,0,NO,YES,YES,NO,NO,cluster5
29 11,52,FEMALE,INNER_CITY,26658,8,NO,0,YES,YES,YES,YES,NO,cluster4
30 12,44,FEMALE,TOWN,15735,8,YES,1,NO,YES,YES,YES,YES,cluster1
31 13,66,FEMALE,TOWN,55204,7,YES,1,YES,YES,YES,YES,YES,cluster1
32 14,36,MALE,RURAL,19474,6,YES,0,NO,YES,YES,YES,NO,cluster5
33 15,38,FEMALE,INNER_CITY,22342,1,YES,0,YES,YES,YES,YES,YES,NO,cluster2
34 16,37,FEMALE,TOWN,17729,8,YES,2,NO,NO,NO,YES,NO,cluster5
35 17,46,FEMALE,SUBURBAN,41016,YES,0,NO,YES,NO,YES,NO,cluster5
36 18,62,FEMALE,INNER_CITY,26909,2,YES,0,NO,YES,NO,NO,YES,cluster4
37 19,31,MALE,TOWN,22522,8,YES,0,YES,YES,YES,NO,NO,cluster2
38 20,61,MALE,INNER_CITY,57880,7,YES,2,NO,YES,NO,NO,YES,cluster2
39 21,50,MALE,TOWN,16497,3,YES,2,NO,YES,YES,NO,NO,cluster5
```

Ocena jakości skupień



Dwa różne spojrzenia

- Wewnętrzne (ocena tylko charakterystyki skupień)
 - Brak dodatkowych źródeł informacji, np. zbioru odniesienia etykiet
 - Miary oceny oparte na danych (internal measures)
- Zewnętrzne
 - „Benchmarking on existing labels”
 - Porównanie skupień z tzw. ground-truth categories / partitions
- Ocena ekspercka

Czy można poszukiwać pojedynczej miary

- Pewne „trudne” rady

“The problem of how to judge the quality of a clustering is difficult and there seems to be no universal answer to it.”

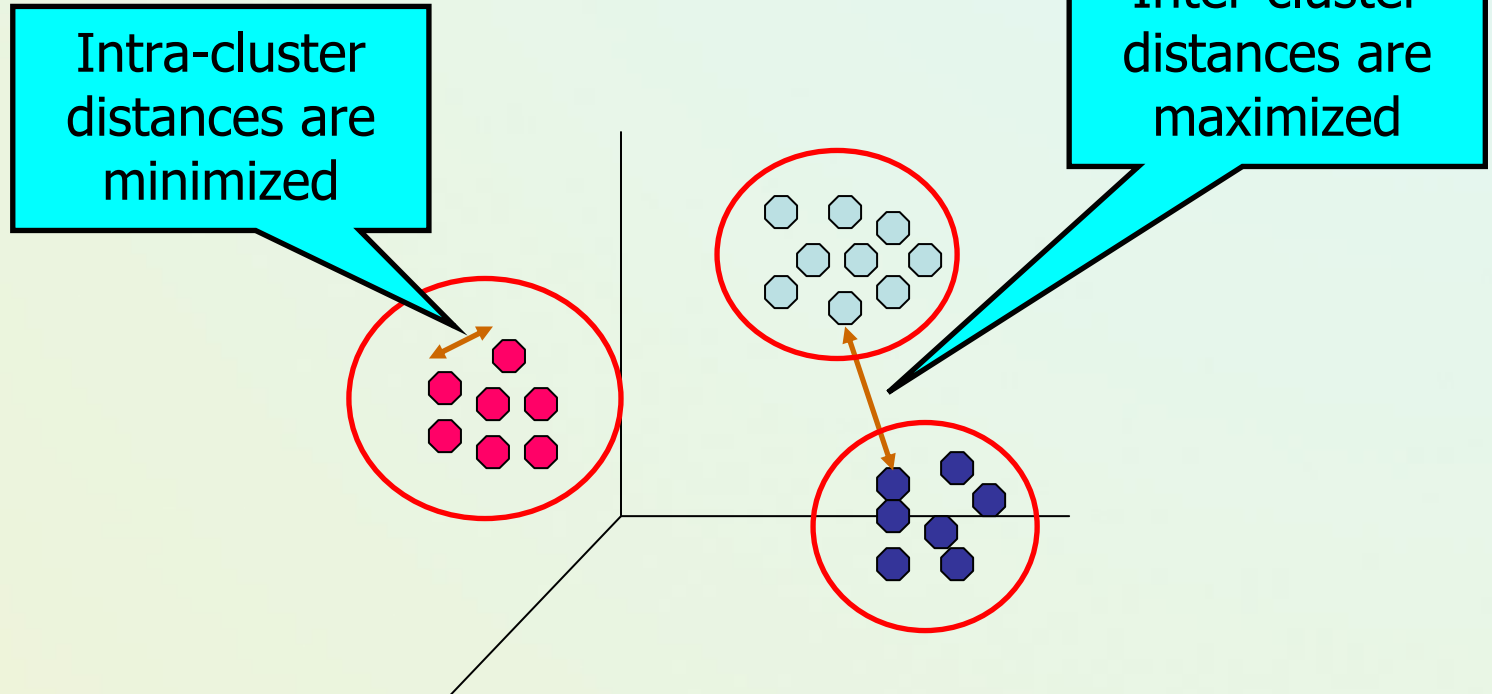
“The nature of processes leading to useful classifications remains little understood, despite considerable effort in this direction.”
— R. Michalski, R. Stepp [MS83]

“How do you know the resulting classifications are any good?”
— D. Fisher [Fis87]

Ocena jakości skupień

Miary oceny oparte na danych (internal measures)

- Oparte na odległościach lub ...
- Duże podobieństwo obiektów wewnątrz skupienia (*Compactness*)
- Same skupienia dość odległe (*Isolation*)



Typowe miary zmienności skupień

- Intuicja → „zmienność wewnątrz-skupieniowa” $wc(C)$ i „zmienność między-skupieniowa” $bc(C)$
 - Można definiować różnymi sposobami
 - Wykorzystaj średni obiekt w skupieniu \mathbf{r}_k (centroids)
 - Wtedy, np.
$$wc(C) = \sum_{k=1}^K \sum_{\mathbf{x} \in C_k} d(\mathbf{x}, \mathbf{r}_k) \quad \mathbf{r}_k = \frac{1}{n_k} \sum_{\mathbf{x} \in C_k} \mathbf{x}$$
$$bc(C) = \sum_{1 \leq j < k \leq K} d(\mathbf{r}_j, \mathbf{r}_k)$$
 - Zamiast bc odległość od globalnego centrum danych (inter-class distance)
$$id = \sum_{C_j} d(\mathbf{r}_j, \mathbf{r}_{glob})$$
- Preferencja dla zwartych, jednorodnych skupień dość odległych od centrum danych

Inne kryteria wewnętrznej jakości skupień

- Compactness → determining the weakest connection within the cluster, i.e., the largest distance between two objects R_i and R_k within the cluster.
- **Isolation** → determining the strongest connection of a cluster to another cluster, i.e., the smallest distance between a cluster centroid and another cluster centroid.

$$\left(\sum_{C_j} \left(\frac{\max(D(R_i, R_k)) \text{ where } (R_i, R_k) \in C_j}{\min(D(C_j, C_m)) \text{ where } C_m \neq C_j} \right) \right)^{-1}$$

- Object positioning → the quality of clustering is determined by the extent to which each object R_j has been correctly positioned in given clusters

$$\sum_{R_j} (\max(D((R_i, R_k))) - \min(D(R_i, R_m)))$$

where $(R_i, R_k) \in C_j$ and $R_m \notin C_j$.

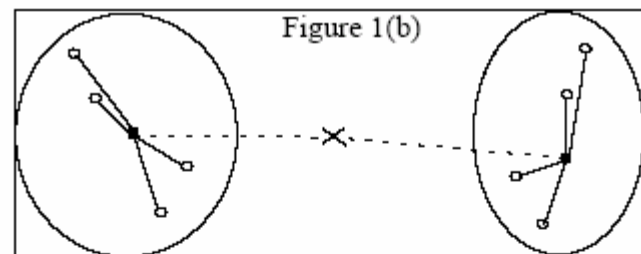
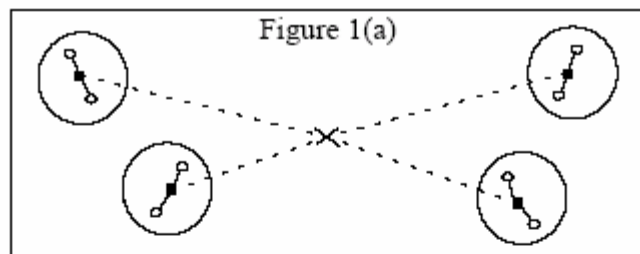


Figure 1: Minimum Total Distance Criterion

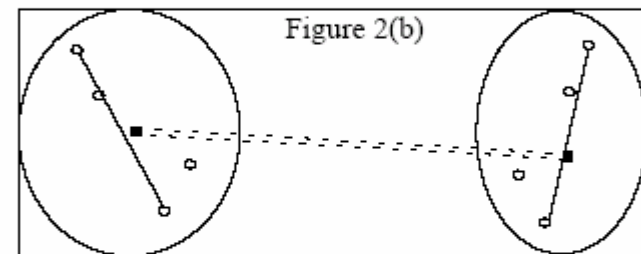
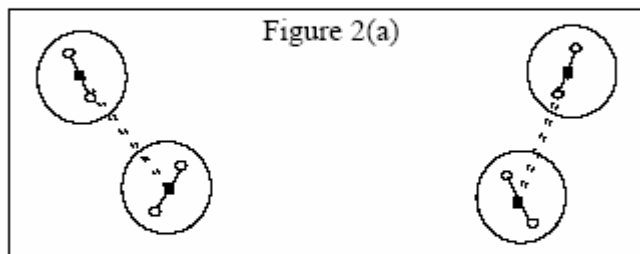


Figure 2: Separated Clusters Criterion

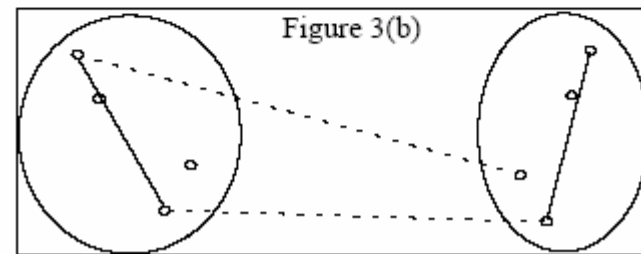
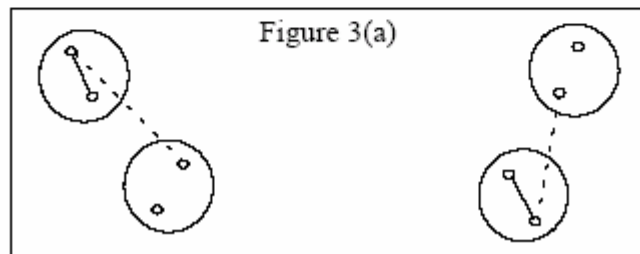
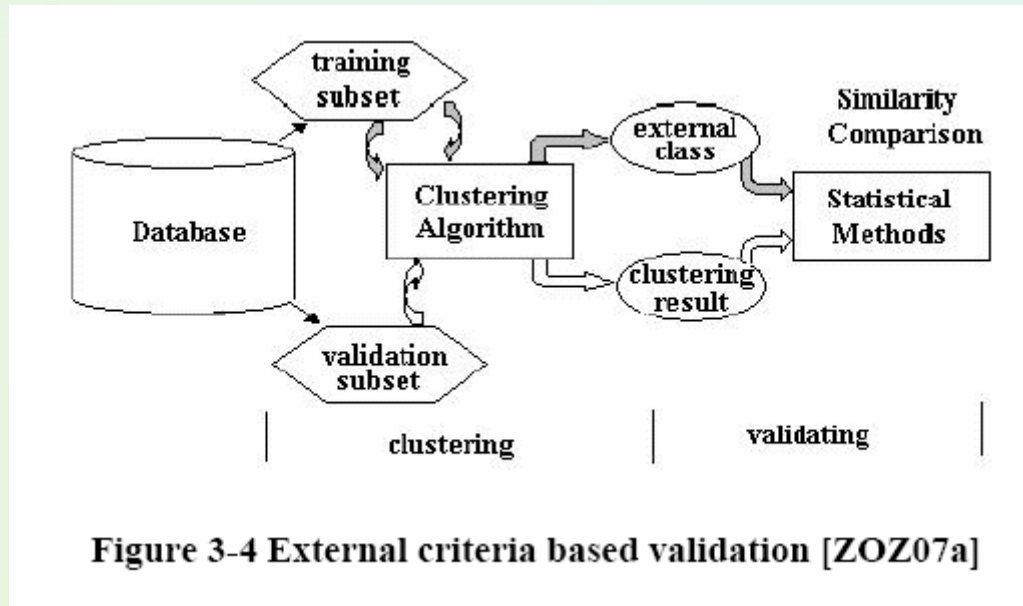


Figure 3: Object Positioning Criterion

Dodatkowe zewnętrzne informacje

- Dysponujemy referencyjnymi etykietami



Odniesienie do zewnętrznego podziału

- Dostępne referencyjne etykiety (manually labeled data)
 - Ekspert etykietuje w zależności od własności danych
 - Istnieją benchmarki TREC, Reuterds, itp..
 - Tryb sem-supervised
- Różne podejścia:
 - „Accuracy of clustering: Percentage of pairs of tuples in the same cluster that share common label” (etykiety w skupieniach)
 - Faworyzujemy małe „czyste” skupienia
 - Czy powinniśmy mieć zgodność liczby skupień i etykiet

Ogólne zasady

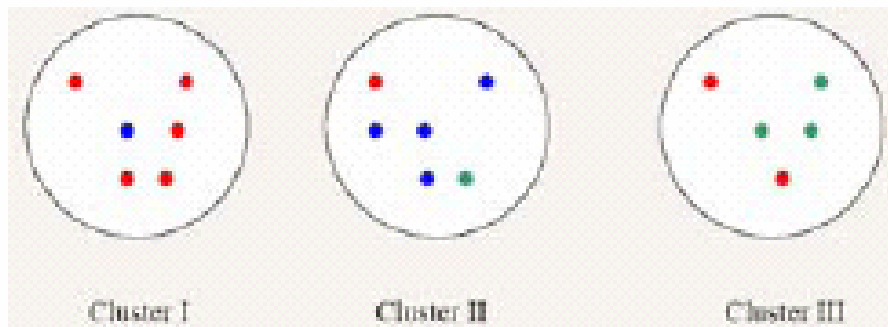
- Homogeneity - Jednorodności
 - Każde skupienie zawiera przykłady z jak najmniejszej liczby etykiet klas
 - Idealnie – tylko jedna klasa
- „Completeness”
 - Każda klasy reprezentowana w możliwie najmniejszej liczbie skupień
- Typowe miary
 - Purity
 - F-miara

Purity – najprostszza miara

- Simple measure: **purity**, the ratio between the dominant class in the cluster and the size of cluster
- Assume documents with C gold standard classes, while our clustering algorithms produce K clusters, $\omega_1, \omega_2, \dots, \omega_K$ with n_i members.

$$Purity(\omega_i) = \frac{1}{n_i} \max_j (n_{ij}) \quad j \in C$$

- Example



Cluster I: Purity = $1/6$ ($\max(5, 1, 0)$) = $5/6$

Cluster II: Purity = $1/6$ ($\max(1, 4, 1)$) = $4/6$

Cluster III: Purity = $1/5$ ($\max(2, 0, 3)$) = $3/5$

Subiektywizm eksperta

- Możliwe różne punkty widzenia



Clustering is subjective



Simpson's Family



School Employees



Females



Males

Ocena grupowania

- Inna niż w przypadku uczenia nadzorowanego (predykcji wartości)
- Poprawność grupowania zależna od oceny obserwatora / analityka
- Różne metody AS są skuteczne przy różnych rodzajach skupień i założeniach, co do danych:
 - Co rozumie się przez skupienie, jaki ma kształt, dobór miary odległości → sferyczne vs. inne
- Dla pewnych metod i zastosowań:
 - Miary zmienności wewnątrz i między – skupieniowych
 - Idea zbiorów kategorii odniesienia (np. TREC)



Grupowanie Dużych Repozytoriów

- Skalowalność (przykłady lecz także „High dimensionality”)
- Uwzględnianie złożonych typów danych
- Konstruowanie dowolnych kształtów skupień
- Wspomaganie parametryzacji (jak dobrać k , parametry DBSCAN, itp.)
- Odporność na szum i obserwacje nietypowe
- Grupowanie przyrostowe
- Przetwarzanie strumieni danych
- „Concept drift”

Problemy i wyzwania

- Znaczący postęp w zakresie skalowalnych algorytmów:
 - Partitioning: *k*-means, *k*-medoids, PAM, CLARANS
 - Hierarchical: BIRCH, CURE
 - Density-based: DBSCAN, CLIQUE, OPTICS
 - Grid-based: STING, WaveCluster.
 - Model-based: Autoclass, Denclue, Cobweb.
- Obecne techniki ciągle nie spełniają wystarczająco dobrze stawianych wymagań.
- Otwarte problemy i wyzwania badawcze; zwłaszcza dla nietypowych i złożonych danych.

Metody hierarchiczne dla dużych zbiorów danych

- Niektóre z ograniczeń metod aglomeracyjnych:
 - słaba skalowalność: złożoność czasowa przynajmniej $O(n^2)$, gdzie n jest liczbą obiektów,
 - „krytyczne” znaczenie decyzji o wyborze punktu połączenia kolejnych skupień w trakcie budowania drzewa hierarchii,
 - algorytmy nie zmieniają, ani nie poprawiają, wcześniej podjętych decyzji.
- Rozwinięcia algorytmów hierarchicznych oraz ich integracja z metodami gęstościowymi:
 - BIRCH (1996): użycie drzew „CF-tree”, uczenie przyrostowe i stopniowa poprawa jakości pod-skupień.
 - CURE (1998): wybór losowy odpowiednio rozproszonych punktów, wstępne grupowanie z określeniem ich punktów reprezentatywnych, łączenie grup w nowe skupienia wraz z przesuwaniem punktów reprezentatywnych w stronę środków tworzonego skupienia zgodnie z „shrinking factor α ”; eliminacja wpływu „outliners”.

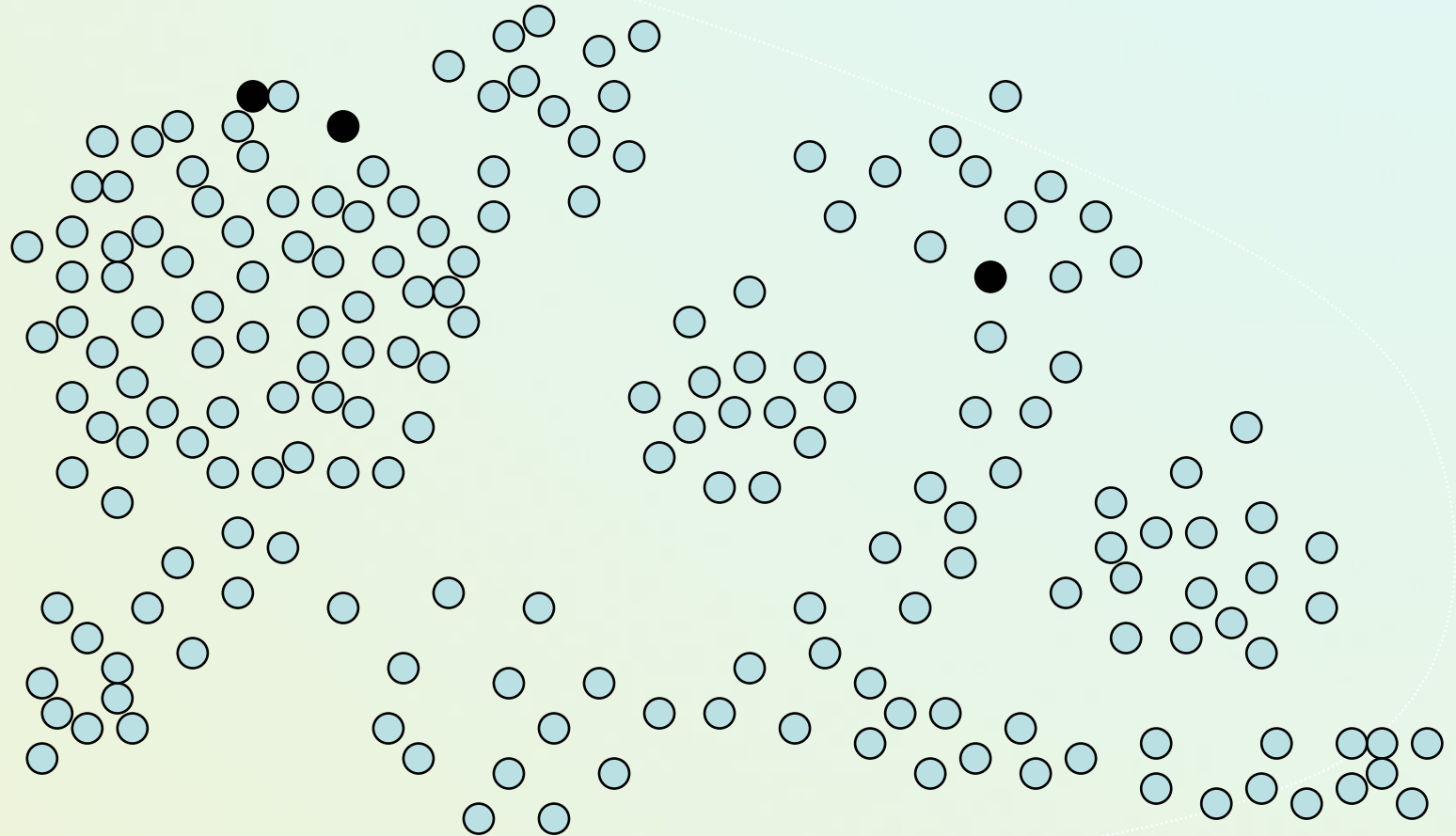
BIRCH – ang. Balanced Iterative Reducing and Clustering using Hierarchies – Zhang et al. (1996)

- Wykorzystuje hierarchiczne drzewo CF (Clustering Feature)
- Działanie algorytmu:
 - **Faza 1**: przyrostowo przeczytaj raz DB w celu zbudowania w pamięci początkowej struktury drzewa CF (rodzaj wielopoziomowej kompresji danych zachowującej wewnętrzną strukturę zgrupowań danych).
 - **Faza 2**: zastosuj wybrany (inny) algorytm skupień dla lepszego pogrupowania obiektów w liściach drzewa CF.
- *Dobra skalowalność*: znajduje zadawalające grupowanie po jednokrotnym przeczytaniu bazy danych i ulepsza je wykorzystując niedużo dodatkowych operacji odczytu DB.
- *Ograniczenia*: zaproponowany dla danych liczbowych, wrażliwość wyników na kolejność prezentacji przykładów.

Metody gęstościowe

- Podstawowe metody wykorzystują miary odległości między obiektami
- Inne metody wykorzystują pojęcie gęstości (ang. density) – lokalne sąsiedztwo punktu/skupienia, a także „gęsto” połączonych punktów
- Właściwości metod gęstościowych:
 - Wykrywanie skupień o dowolnych kształtach (niesferycznych)
 - Odporność na „szum informacyjny”
 - Jednokrotne przeglądanie DB
 - Potrzebna parametryzacja oceny gęstości i warunków zatrzymania
- Interesujące algorytmy:
 - DBSCAN: Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - DENCLUE: Hinneburg & D. Keim (KDD'98)
 - CLIQUE: Agrawal, et al. (SIGMOD'98)

Modelowanie dowolnych kształtów



DBSCAN: Algorytm gęstościowy

- DBSCAN: Density Based Spatial Clustering of Applications with Noise.
 - Wykorzystuje pojęcie „*density-based cluster*”: Skupienie będące maksymalnym zbiorem punktów gęsto połączonych „*density-connected points*”.
 - Poszukuje się zgrupowań odpowiednio gęsto (blisko siebie) położonych obiektów (*dense regions/clusters*) oddzielonych od siebie obszarami o niskiej gęstości („noise”)
 - Możliwość wykrywania skupień o dowolnym kształcie w obecności szumu informacyjnego (noise)

DBSCAN: Podstawowe pojęcia

Parametry:

- ***Eps***: Maksymalny promień sąsiedztwa
- ***MinPts***: minimalna liczba punktów (obiektów) w *Eps*-sąsiedztwie badanego punktu

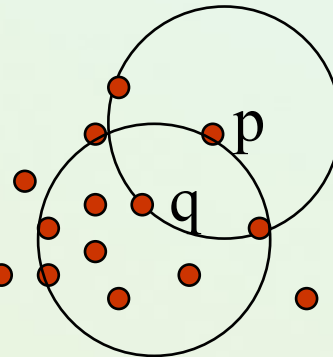
$N_{Eps}(p)$: {punkt q należy do D | $dist(p,q) \leq Eps$ }

Directly density-reachable: A point p is directly density-reachable from a point q wrt. ***Eps***, ***MinPts*** if

1) p belongs to $N_{Eps}(q)$

2) core point condition:

$$|N_{Eps}(q)| \geq MinPts$$



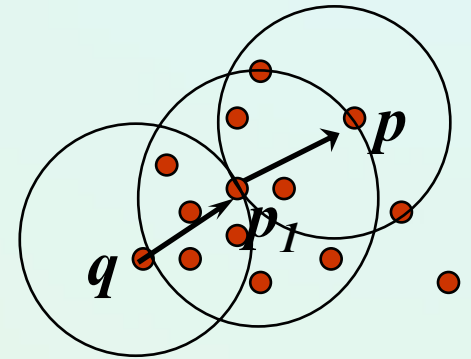
MinPts = 5

Eps = 1 cm

DBSCAN: Podstawowe pojęcia (II)

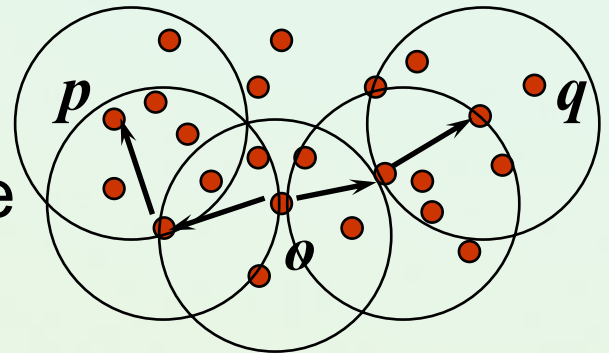
■ Density-reachable:

- A point p is density-reachable from a point q wrt. Eps , $MinPts$ if there is a chain of points p_1, \dots, p_n , $p_1 = q$, $p_n = p$ such that p_{i+1} is directly density-reachable from p_i

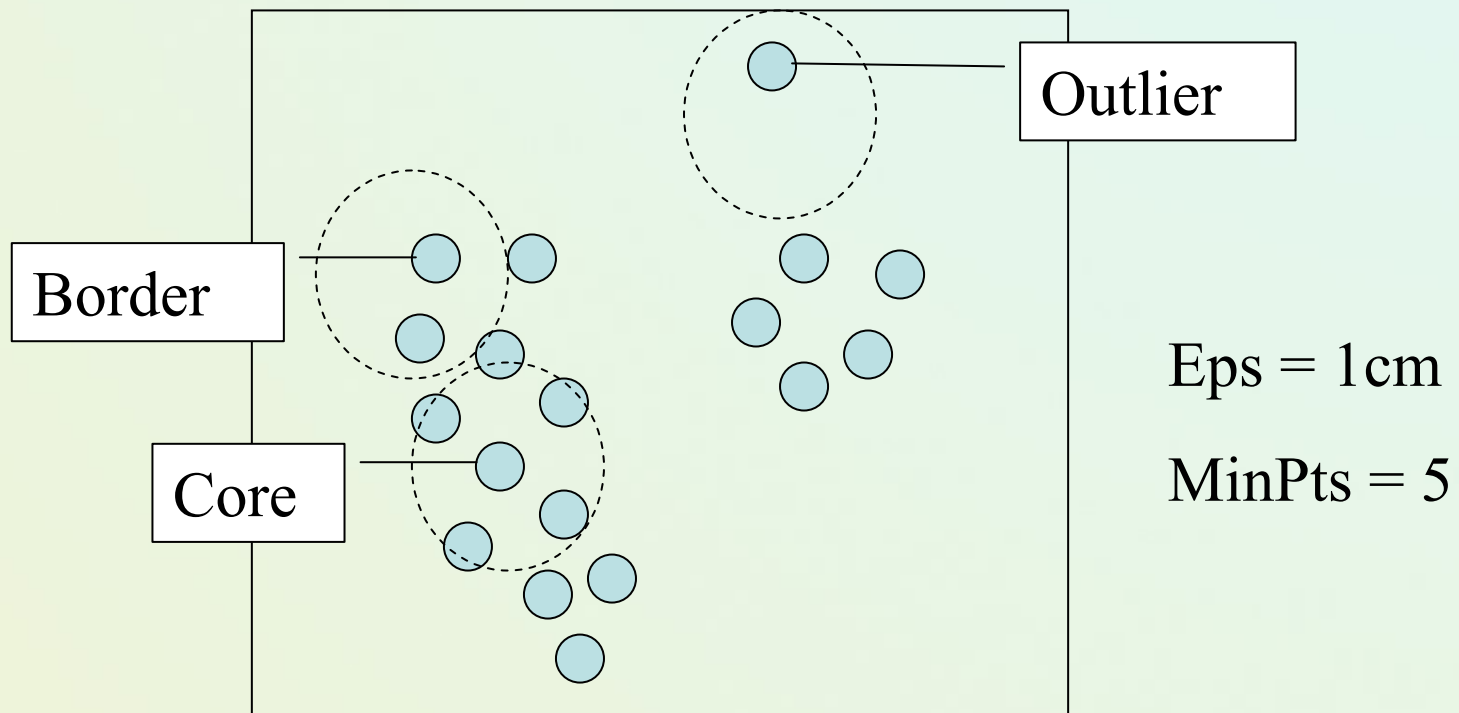


■ Density-connected

- A point p is density-connected to a point q wrt. Eps , $MinPts$ if there is a point o such that both, p and q are density-reachable from o wrt. Eps and $MinPts$.



DBSCAN: General Ideas



DBSCAN: Zarys algorytmu

- Wybierz punkt startowy p
- Odnajdź wszystkie punkty do gęstościowego osiągnięcia z p (density-reachable from p wrt ***Eps*** and ***MinPts***).
- Jeśli p jest rdzeniem (*core point*), utwórz skupienie.
- Jeśli p jest punktem granicznym (border point) i żadne punkty nie są z niego gęstościowo osiągalne, DBSCAN wybiera następny punkt z bazy danych
- Proces jest konytuowany dopóki żaden nowy punkt nie może być dodany to dowolnego skupienia.
- Złożoność: $O(n \log n)$ w przypadku użycia specjalnego „spatial index”, w przeciwnym razie $O(n^2)$.

Clustering in Data Mining – szukaj więcej!

Data Clustering: A Review

A.K. JAIN

Michigan State University

M.N. MURTY

Indian Institute of Science

AND

P.J. FLYNN

The Ohio State University

Clustering is the unsupervised classification of patterns (observations, data items, or feature vectors) into groups (clusters). The clustering problem has been addressed in many contexts and by researchers in many disciplines; this reflects its broad appeal and usefulness as one of the steps in exploratory data analysis. However, clustering is a difficult problem combinatorially, and differences in assumptions and contexts in different communities has made the transfer of useful generic concepts and methodologies slow to occur. This paper presents an overview of pattern clustering methods from a statistical pattern recognition perspective, with a goal of providing useful advice and references to fundamental concepts accessible to the broad community of clustering practitioners. We present a taxonomy of clustering techniques, and identify cross-cutting themes and recent advances. We also describe some important applications of clustering algorithms such as image segmentation, object recognition, and information retrieval.

Categories and Subject Descriptors: I.5.1 [Pattern Recognition]: Models; I.5.3 [Pattern Recognition]: Clustering; I.5.4 [Pattern Recognition]: Applications—*Computer vision*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Clustering*; I.2.6 [Artificial Intelligence]: Learning—*Knowledge acquisition*

Analiza skupień - podsumowanie

- Liczne i ważne zastosowanie praktyczne analizy skupień (AS).
- AS używana „samodzielnie” w zgłębianiu danych, lub jako jedno z narzędzi podczas wstępnego przetwarzania w procesie KDD.
- Jakość skupień i działanie wielu algorytmów związane są określeniem miary odległości obiektów.
- Podstawowe klasy metod:
 - hierarchiczne,
 - podziałowo/optymalizacyjne,
 - gęstościowe,
 - „grid-based”,
 - wykorzystujące modele matematyczne (np. probabilistyczne lub neuronowe)
- Ważne zagadnienie to także wykrywanie obiektów nietypowych (outliers discovery).

Może pytanie lub komentarze?

