

MAREK WALESIAK

ANDRZEJ DUDEK

Akademia Ekonomiczna

Wrocław

**SYMULACYJNA OPTYMALIZACJA WYBORU
PROCEDURY KLASYFIKACYJNEJ DLA DANEGO TYPU DANYCH
– CHARAKTERYSTYKA PROBLEMU**

1. Wprowadzenie

W literaturze przedmiotu w typowej procedurze klasyfikacyjnej wyodrębnia się osiem etapów¹: 1) wybór obiektów do klasyfikacji; 2) wybór zmiennych charakteryzujących obiekty; 3) wybór formuły normalizacji wartości zmiennych; 4) wybór miary odległości²; 5) wybór metody klasyfikacji; 6) ustalenie liczby klas; 7) walidacja wyników klasyfikacji; 8) opis (interpretacja) i profilowanie klas. Do newralgicznych zalicza się etapy dotyczące wyboru formuły normalizacji wartości zmiennych, miary odległości, metody klasyfikacji i ustalenia liczby klas, które mają w znacznej mierze arbitralny charakter.

W artykule zaprezentowano szczegółową charakterystykę dziewięciu ścieżek w symulacyjnej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych, a wyodrębnionych w zależności od typu skali pomiaru zmiennych. Liczba rozpatrywanych wariantów procedury klasyfikacyjnej zależy od liczby formuł normalizacyjnych, liczby typów miar odległości i liczby

¹ Por. [13], s. 342–343; [21].

² Zob. [23].

metod klasyfikacji. Na podstawie zaproponowanego podejścia w artykule M. Walesiaka i A. Dudka³ scharakteryzowano program komputerowy cluster-Sim (opracowany w języku R oraz pomocniczo w języku C++) służący do realizacji wyodrębnionych ścieżek oraz wybrane wyniki obliczeń symulacyjnych.

1. Charakterystyka ścieżek w symulacyjnej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych

Punktem wyjścia analizy symulacyjnej jest macierz danych. W zależności od typu skali pomiaru zmiennych wyróżniono dziewięć ścieżek w procedurze symulacyjnej. Przy opracowywaniu poszczególnych ścieżek uwzględniono następujące elementy:

- typ skali pomiaru zmiennych w macierzy danych,
- typ formuły normalizacyjnej dla zmiennych mierzonych na skali przedziałowej i (lub) ilorazowej,
- miary odległości właściwe dla poszczególnych typów skal pomiaru zmiennych,
- typ metody klasyfikacji,
- miernik oceny jakości klasyfikacji.

Liczba rozpatrywanych wariantów procedury klasyfikacyjnej dla zmiennych mierzonych na skali:

1. Ilorazowej równa się 368 (11 formuł normalizacyjnych, 7 typów miar odległości, 8 metod klasyfikacji)⁴.

2. Przedziałowej (lub ilorazowej i przedziałowej) równa się 140 (5 formuł normalizacyjnych, 5 typów miar odległości, 8 metod klasyfikacji).

3. Porządkowej równa się 5 (miara odległości GDM2, 5 metod klasyfikacji).

³ Zob. [23].

⁴ Liczba wariantów w punktach 1, 2, 6, 7 nie wynika z prostego przemnożenia liczby formuł normalizacyjnych, typów miar odległości i metod klasyfikacji z uwagi na ograniczenia w ich stosowaniu.

4. Nominalnej wielostanowej równa się 5 (odległość Sokala i Michenera, 5 metod klasyfikacji).

5. Nominalnej binarnej równa się 50 (10 typów miar odległości, 5 metod klasyfikacji).

6. Ilorazowej bez normalizacji równa się 38 (7 typów miar odległości, 8 metod klasyfikacji).

7. Przedziałowej (lub ilorazowej i przedziałowej) bez normalizacji równa się 28 (5 typów miar odległości, 8 metod klasyfikacji).

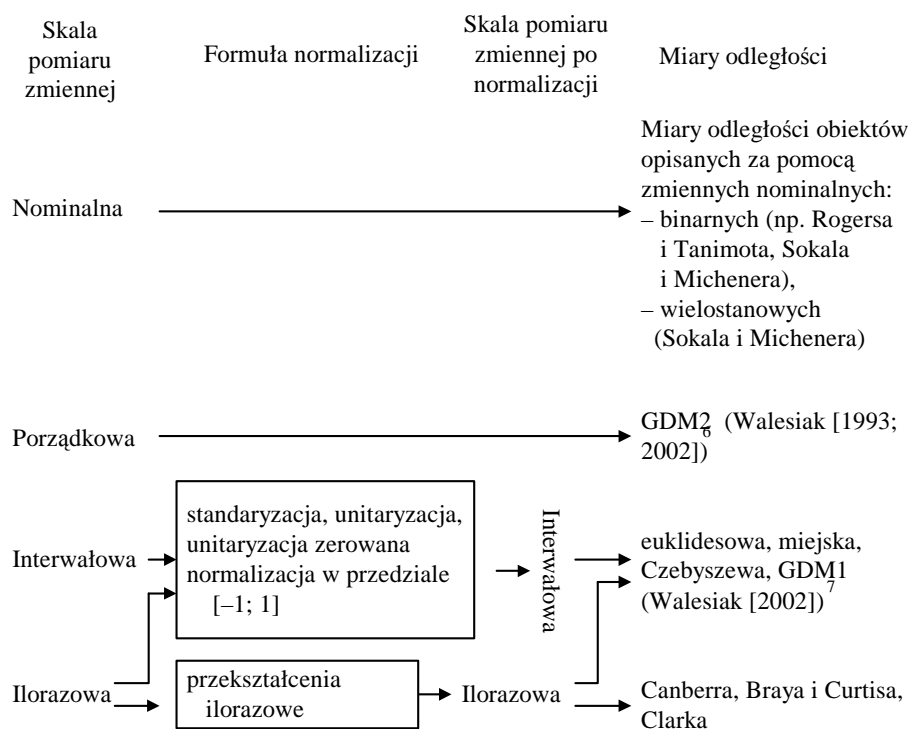
8. Ilorazowej z zastosowaniem metody k -średnich równa się 11 (11 formuł normalizacyjnych, 1 metoda klasyfikacji).

9. Przedziałowej (lub ilorazowej i przedziałowej) z zastosowaniem metody k -średnich równa się 5 (5 formuł normalizacyjnych, 1 metoda klasyfikacji).

Ścieżki 6 i 7 (dla danych metrycznych bez normalizacji) występują często w badaniach marketingowych opartych na danych otrzymanych na przykład ze skali Likerta lub semantycznej. Wprawdzie są to przykłady skal porządkowych, jednak z uwagi na to, że odstęp między kategoriami odpowiada w przybliżeniu jednakowym interwałom, traktuje się je w badaniach jako skale metryczne⁵.

Przy wyborze miar odległości obiektów opisanych zmiennymi mierzonymi na skali przedziałowej i (lub) ilorazowej należy wziąć pod uwagę zastosowaną formułę normalizacji wartości zmiennych. Klasyfikację formuł normalizacyjnych i miar podobieństwa obiektów z punktu widzenia skal pomiaru zmiennych przedstawiono na rysunku 1.

⁵ Zob. [7], s. 75.



Rys. 1. Klasyfikacja formuł normalizacyjnych i miar odległości obiektów z punktu widzenia skal pomiaru zmiennych

Źródło: [22].

Rozważania w artykule ograniczono do najczęściej wykorzystywanych metod klasyfikacji, czyli metody *k-medoids* i siedmiu hierarchicznych metod aglomeracyjnych opartych na macierzy odległości oraz metody *k-średnich* opartych na macierzy danych.

Indeksy oceny jakości klasyfikacji pozwalające na wyznaczenie optymalnej liczby klas dzieli się na globalne i lokalne⁸. Indeksy globalne $G(u)$ są oparte

⁶ Zob. [1]; [20].

⁷ *Ibidem*.

⁸ Zob. [6], s. 61.

na kompletnym zbiorze danych klasyfikacyjnych (zazwyczaj $(u = 2, \dots, n - 1)$, gdzie u oznacza liczbę klas, a n – liczbę obiektów) i z ich wykorzystaniem poszukuje się optymalnego podziału badanego zbioru obiektów na klasy.

Indeksy lokalne $L(u)$ są oparte tylko na podzbiore zbioru danych klasyfikacyjnych (zazwyczaj rozpatruje się dwa sąsiadujące podziały zbioru obiektów, czyli podziały na u oraz $u + 1$ klas, lub odwrotnie) i pozwalają na ocenę, czy dana klasa powinna być podzielona na dwie klasy (lub para klas powinna być połączona w jedną). Proces podziału (łączenia) jest kontynuowany do momentu osiągnięcia określonego progu lub odrzucenia określonej hipotezy zerowej.

Ze względu na to, że w zaproponowanej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych rozpatrywane są wszystkie warianty procedury klasyfikacyjnej danego zbioru obiektów (uzależnione od typu formuły normalizacyjnej, typu miary odległości i metody klasyfikacji), nie jest możliwe wykorzystanie kryteriów lokalnych do wyznaczenia optymalnej liczby klas,.

G.W. Milligan i M.C. Cooper⁹ przetestowali na podstawie zbiorów danych o znanej strukturze klas 30 indeksów pozwalających wyznaczyć optymalną liczbę klas. Do oceny wyników symulacji w zaproponowanej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych wprowadzono trzy najlepsze mierniki globalne z eksperymentu G.W. Milligana i M.C. Cooper: R.B. Calińskiego i J. Harabasza, F.B. Bakera i L.J. Huberta, L.J. Huberta i J.R. Levine, oraz dwa indeksy często wykorzystywane w literaturze w testach porównawczych¹⁰: Silhouette, W.J. Krzanowskiego i Y.T. Lai:

1. Indeks Calińskiego i Harabasza¹¹:

$$G1(u) = \frac{tr(\mathbf{B})/(u-1)}{tr(\mathbf{W})/(n-u)}, \quad G1(u) \in R_+, \quad (1)$$

gdzie:

⁹ Zob. [14].

¹⁰ Zob. np. [3]; [15]; [17]; [18].

¹¹ Zob. [2].

- B** – macierz kowariancji międzyklasowej,
W – macierz kowariancji wewnątrzklasowej,
tr – ślad macierzy,
u – liczba klas ($u = 2, \dots, n - 1$),
n – liczba obiektów.

Indeks Calińskiego i Harabasa jest nazywany pseudostatystyką F^{12} .

2. Indeks Gamma Bakera i Huberta¹³:

$$G2(u) = \frac{s(+)-s(-)}{s(+)+s(-)}, \quad G2(u) \in [-1, 1], \quad (2)$$

gdzie:

- s* (+) – liczba par odległości zgodnych,
s (-) – liczba par odległości niezgodnych,
u – liczba klas ($u = 2, \dots, n - 1$).

Przy obliczaniu indeksu Gamma¹⁴ porównuje się wszystkie odległości wewnątrzklasowe z wszystkimi odległościami międzyklasowymi. Liczba tych porównań wynosi więc $r \cdot c$, gdzie r (c) – liczba odległości wewnątrzklasowych (międzyklasowych). Jeśli odległość wewnątrzklasowa jest mniejsza (większa) niż odległość międzyklasowa, to parę taką uznajemy za zgodną (niezgodną). Odległości wewnątrzklasowe równe międzyklasowym nie są uwzględniane.

3. Indeks Huberta i Levine¹⁵:

$$G3(u) = \frac{D(u) - r \cdot D_{\min}}{r \cdot D_{\max} - r \cdot D_{\min}}, \quad D_{\min} \neq D_{\max}, \quad G3(u) \in (0, 1), \quad (3)$$

gdzie:

¹² Zob. [12], s. 291.

¹³ Zob. [1]; [8].

¹⁴ Por. [6], s. 62.

¹⁵ Zob. [9].

$D(u)$ – suma wszystkich odległości wewnątrzklasowych,
 r – liczba odległości wewnątrzklasowych,
 D_{min} – najmniejsza odległość wewnątrzklasowa,
 D_{max} – największa odległość wewnątrzklasowa,
 u – liczba klas ($u = 2, \dots, n - 2$).

4. Indeks Silhouette¹⁶:

$$S(u) = \sum_{i=1}^n S(i)/n, \quad S(u) \in [-1, 1], \quad (4)$$

gdzie:

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i); b(i)\}},$$

$i, k = 1, \dots, n$ – numer obiektu,

$a(i) = \sum_{k \in \{P_r \setminus i\}} d_{ik} / (n_r - 1)$ – średnia odległość obiektu i od pozostałych
 obiektów należących do klasy P_r ,

$$b(i) = \min_{s \neq r} \{d_{iP_s}\},$$

$d_{iP_s} = \sum_{k \in P_s} d_{ik} / n_s$ – średnia odległość obiektu i od obiektów należących do
 klasy P_s ,

$r, s = 1, \dots, u$ – numer klasy,

$u = 2, \dots, n - 1$ – liczba klas.

5. Indeks Krzanowskiego i Lai¹⁷:

$$KL(u) = \left| \frac{DIFF_u}{DIFF_{u+1}} \right|, \quad KL(u) \in R_+, \quad (5)$$

$$DIFF_u = (u - 1)^{2/m} tr \mathbf{W}_{u-1} - u^{2/m} tr \mathbf{W}_u,$$

¹⁶ Zob. [10]; [16].

¹⁷ Zob. [11].

gdzie:

\mathbf{W} – macierz kowariancji wewnątrzklasowej,

u – liczba klas ($u = 2, \dots, n - 3$),

n – liczba obiektów, m – liczba zmiennych.

Indeksy $G1(u)$ i $KL(u)$ są oparte na macierzy danych, natomiast indeksy $G2(u)$, $G3(u)$ i $S(u)$ na macierzy odległości. Maksymalna wartość $G1(u)$, $G3(u)$, $S(u)$ i $KL(u)$ oraz minimalna $G2(u)$ wskazuje najlepszy podział zbioru obiektów, a zarazem wyznacza liczbę klas.

Szczegółowe charakterystyki ścieżek zaprezentowano w tabeli 1.

3. Podsumowanie

W typowym studium klasyfikacyjnym etapy wyboru formuły normalizacji wartości zmiennych, miary odległości, metody klasyfikacji oraz ustalenia liczby klas mają zwykle arbitralny charakter. Zaletą tego podejścia jest obiektywizacja problemu ich wyboru. Uzyskuje się to w wyniku przeprowadzenia symulacyjnej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych z wykorzystaniem programu clusterSim¹⁸. Miernikami oceny wszystkich procedur klasyfikacyjnych badanego zbioru obiektów są globalne indeksy oceny jakości klasyfikacji pozwalające na wyznaczenie optymalnej liczby klas.

Zaprezentowane podejście ma pewne ograniczenia:

- a) w literaturze jest ponad 40 mierników oceny jakości klasyfikacji; w zaprezentowanym podejściu możliwe było uwzględnienie tylko indeksów globalnych;
- b) spośród indeksów globalnych uwzględniono pięć najważniejszych, jednak ostateczny wybór jednego z nich nadal jest arbitralny.

¹⁸ Zob. [23].

Ścieżki w symulacyjnej optymalizacji wyboru procedury klasyfikacyjnej dla danego typu danych

Nr	Etapy typowej procedury klasyfikacyjnej	Numer ścieżki w procedurze symulacyjnej										
		1	2	3	4	5	6	7	8	9		
I	Wybór obiektów i zmiennych	macierz danych $[X_{ij}]$										
II	Skala pomiaru zmiennych	ilorazowa	ilorazowa	przedziałowa ¹	porządkowa	nominalna wielostanowa	binarna	ilorazowa	przedziałowa ¹	ilorazowa	przedziałowa ¹	
	Wybór formuły normalizacji ²	n6–n11	n1–n5	n1–n5	NA	NA		bez normalizacji	n6–n11/ n1–n5	n1–n5		
	Skala pomiaru zmiennych po normalizacji	ilorazowa	przedziałowa	przedziałowa	porządkowa	nominalna wielostanowa	binarna	ilorazowa	przedziałowa ¹	ilorazowa/ przedziałowa	przedziałowa	
III	Wybór miary odległości ³	d1–d7	d1–d5	d1–d5	d8	d9	b1–b10	d1–d7	d1–d5	NA		
IV	Wybór metody klasyfikacji	1. Pojedynczego połączenia 2. Kompletnego połączenia		3. Średniej klasowej 4. Ważonej średniej klasowej		5. <i>k-medoids</i> (pam) 6. Warda ⁴		7. Centroidalna ⁴ 8. Medianowa ⁴		<i>k</i> -średnich		
V	Liczba możliwości	$\frac{[(6 \times 7 \times 5) + (6 \times 1 \times 3)] + [(5 \times 5 \times 5) + (5 \times 1 \times 3)]}{368}$		$\frac{(5 \times 5 \times 5) + (5 \times 1 \times 3)}{140}$		1 x 5 = 5	1 x 5 = 5	10 x 5 = 50	$\frac{(7 \times 5) + (1 \times 3)}{38}$	$\frac{(5 \times 5) + (1 \times 3)}{28}$	11	5
	Miernik jakości klasyfikacji	1. Caliński & Harabasz (G1) 2. Baker & Hubert (G2) 3. Hubert & Levine (G3) 4. Silhouette (S) 5. Krzanowski & Lai (KL)			1. NA 2. G2 3. G3 4. S 5. NA			1. G1 2. G2 3. G3 4. S 5. KL		1. G1 2. NA 3. NA 4. NA 5. KL		

¹ Lub ilorazowa i przedziałowa.

² n1 (n2) – standaryzacja klasyczna (Webera), n3 – unitaryzacja, n4 – unitaryzacja zerowana, n5 – normalizacja w przedziale [-1; 1], n6–n11 – przekształcenia ilorazowe.

³ d1 – miejska, d2 – euklidesowa, d3 – Czebyszewa, d4 – kwadrat euklidesowej, d5 – GDM1, d6 – Canberra, d7 – Braya-Curtisa, d8 – GDM2, d9 – Sokala i Michenera dla zmiennych nominalnych; odległości dla zmiennych binarnych (dostępne w procedurze dist.binary): b1 = Jaccard; b2 = Sokal & Michener; b3 = Sokal & Sneath (1); b4 = Rogers & Tanimoto; b5 = Czekanowski; b6 = Gower & Legendre (1); b7 = Ochiai; b8 = Sokal & Sneath (2); b9 = Phi of Pearson; b10 = Gower & Legendre (2).

⁴ Metody klasyfikacji przyjmujące założenie, że odległości między obiektami zostały wyznaczone za pomocą kwadratu odległości euklidesowej, tylko bowiem w tym przypadku metody te mają interpretację geometryczną, zgodną z ich nazwami.

NA – nie stosuje się.

Źródło: opracowanie własne (opisy metod znajdują się m.in. w następujących pracach: [4]; [5]; [6]).

Literatura

1. Baker F.B., Hubert L.J.: *Measuring the power of hierarchical cluster analysis*. „Journal of the American Statistical Association” 1975, vol. 70, No 349.
2. Caliński R.B., Harabasz J.: *A dendrite method for cluster analysis*. „Communications in Statistics” 1974, vol. 3.
3. Dudoit S., Fridlyand J.: *A prediction-based resampling method for estimating the number of clusters in a dataset*. „Genome Biology” 2002, vol. 3, No 7.
4. Everitt B.S., Landau S., Leese M.: *Cluster analysis*. Edward Arnold, London 2001.
5. Gatnar E., Walesiak M.: *Metody statystycznej analizy wielowymiarowej w badaniach marketingowych*. Wyd. Naukowe Akademii Ekonomicznej we Wrocławiu, Wrocław 2004.
6. Gordon A.D.: *Classification*. Chapman and Hall/CRC, London 1999.
7. Górniak J.: *My i nasze pieniądze*. Wyd. Aureus, Kraków 2000.
8. Hubert L.J. *Approximate evaluation technique for the single-link and complete-link hierarchical clustering procedures*. „Journal of the American Statistical Association” 1974, vol. 69, No 347.
9. Hubert L.J., Levine J.R.: *Evaluating object set partitions: free sort analysis and some generalizations*, „Journal of Verbal Learning and Verbal Behaviour” 1976, vol. 15.
10. Kaufman L., Rousseeuw P.J.: *Finding groups in data: an introduction to cluster analysis*. Wiley, New York 1990.
11. Krzanowski W.J., Lai Y.T.: *A criterion for determining the number of groups in a data set using sum of squares clustering*. „Biometrics” 1985, No 44.
12. Lattin J.M., Carroll J.D., Green P.E.: *Analyzing multivariate data*. Brooks/Cole, Pacific Grove 2003.
13. Milligan G.W.: *Clustering validation: results and implications for applied analyses*. W: *Clustering and classification*. Red. P. Arabie, L.J. Hubert, G. de Soete. World Scientific, Singapore 1996.

14. Milligan G.W., Cooper M.C.: *An examination of procedures for determining the number of clusters in a data set.* „Psychometrika” 1985, No 2.
15. Mufti G.B., Bertrand P., El Moubarki L.: *Determining the number of groups from measures of cluster stability.* W: *Applied Stochastic Models and Data Analysis.* Red. J. Janssen, P. Lenca. ENST Bretagne, Brest 2005.
16. Rousseeuw P.J.: *Silhouettes: a graphical aid to the interpretation and validation of cluster analysis.* „Journal of Computational and Applied Mathematics” 1987, No 20.
17. Sugar C.A., James G.H.: *Finding the number of clusters in a dataset: an information-theoretic approach.* „Journal of the American Statistical Association” 2003, vol. 98, No 463.
18. Tibshirani R., Walther G., Hastie T.: *Estimating the number of clusters in a data set via the gap statistic.* „Journal of the Royal Statistical Society” 2001, ser. B, vol. 63, part 2.
19. Walesiak M.: *Statystyczna analiza wielowymiarowa w badaniach marketingowych.* Prace Naukowe Akademii Ekonomicznej we Wrocławiu nr 654, seria Monografie i Opracowania nr 101. Wrocław 1993.
20. Walesiak M.: *Uogólniona miara odległości w statystycznej analizie wielowymiarowej.* Wyd. Akademii Ekonomicznej, Wrocław 2002.
21. Walesiak M.: *Rekomendacje w zakresie strategii postępowania w procesie klasyfikacji zbioru obiektów.* W: *Przestrzenno-czasowe modelowanie i prognozowanie zjawisk gospodarczych.* Red. A. Zeliaś. Wyd. Akademii Ekonomicznej w Krakowie, Kraków 2005.
22. Walesiak M.: *Uogólniona miara odległości w statystycznej analizie wielowymiarowej.* Wydanie II rozszerzone. Wyd. Akademii Ekonomicznej we Wrocławiu, Wrocław 2006.
23. Walesiak M., Dudek A.: *Symulacyjna optymalizacja wyboru procedury klasyfikacyjnej dla danego typu danych – oprogramowanie komputerowe i wyniki badań.* W: *Klasyfikacja i analiza danych – teoria i zastosowania.* Red. K. Jajuga, M. Walesiak. Prace Naukowe Akademii Ekonomicznej we Wrocławiu (w redakcji).

**DETERMINATION OF OPTIMAL CLUSTERING PROCEDURE
FOR A DATA SET
– THE CHARACTERISATION OF THE PROBLEM**

Summary

In typical cluster analysis study eight major steps are distinguished (see Milligan [1996], 342–343; Walesiak [2005]). Four of them represent the critical steps: decisions concerning variable normalisation formula, selection of a distance measure, selection of clustering method, determining the number of clusters.

The article presents determination of optimal clustering procedure for a data set by varying all combinations of normalization formulas, distance measures, and clustering methods. Nine paths of simulation was separated depends on variable scale of measurement in a data set. Based on this approach in article of Walesiak and Dudek [2005] the clusterSim computer program written in R and C++ languages was proposed.

Translated by Marek Walesiak, Andrzej Dudek