

PRACE NAUKOWE

Uniwersytetu Ekonomicznego we Wrocławiu

RESEARCH PAPERS

of Wrocław University of Economics

Nr 384

Taksonomia 24

**Klasyfikacja i analiza danych –
teoria i zastosowania**

Redaktorzy naukowi

Krzysztof Jajuga

Marek Walesiak



Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu
Wrocław 2015

Redaktor Wydawnictwa: Aleksandra Śliwka

Redaktor techniczny: Barbara Łopusiewicz

Korektor: Barbara Cibis

Łamanie: Beata Mazur

Projekt okładki: Beata Dębska

Tytuł dofinansowany ze środków Narodowego Banku Polskiego
oraz ze środków Sekcji Klasyfikacji i Analizy Danych PTS

Informacje o naborze artykułów i zasadach recenzowania
znajdują się na stronie internetowej Wydawnictwa
www.pracnaukowe.ue.wroc.pl
www.wydawnictwo.ue.wroc.pl

Publikacja udostępniona na licencji Creative Commons
Uznanie autorstwa-Użycie niekomercyjne-Bez utworów zależnych 3.0 Polska
(CC BY-NC-ND 3.0 PL)



© Copyright by Uniwersytet Ekonomiczny we Wrocławiu
Wrocław 2015

ISSN 1899-3192 (Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu)
e-ISSN 2392-0041 (Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu)
ISSN 1505-9332 (Taksonomia)

Wersja pierwotna: publikacja drukowana

Zamówienia na opublikowane prace należy składać na adres:
Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu
tel./fax 71 36 80 602; e-mail:econbook@ue.wroc.pl
www.ksiegarnia.ue.wroc.pl

Druk i oprawa: TOTEM

Spis treści

Wstęp.....	9
Krzysztof Jajuga, Józef Pociecha, Marek Walesiak: 25 lat SKAD.....	15
Beata Basiura, Anna Czapkiewicz: Symulacyjne badanie wykorzystania entropii do badania jakości klasyfikacji.....	25
Andrzej Bąk: Zagadnienie wyboru optymalnej procedury porządkowania liniowego w pakiecie <code>pllord</code>	33
Justyna Brzezińska: Analiza klas ukrytych w badaniach sondażowych.....	42
Grażyna Dehnel: Rejestr podatkowy oraz rejestr ZUS jako źródło informacji dodatkowej dla statystyki gospodarczej – możliwości i ograniczenia ..	51
Sabina Denkowska: Wybrane metody oceny jakości dopasowania w <i>Propensity Score Matching</i>	60
Marta Dziechciarz-Duda, Klaudia Przybysz: Zastosowanie teorii zbiorów rozmytych do identyfikacji pozafiskalnych czynników ubóstwa.....	75
Iwona Foryś: Potencjał rynku mieszkaniowego w Polsce w latach dekonjunktury gospodarczej.....	84
Eugeniusz Gatnar: Statystyczna analiza konwergencji krajów Europy Środkowej i Wschodniej po 10 latach członkostwa w Unii Europejskiej.....	93
Ewa Genge: Zaufanie do instytucji publicznych i finansowych w polskim społeczeństwie – analiza empiryczna z wykorzystaniem ukrytych modeli Markowa.....	100
Alicja Grześkowiak: Wielowymiarowa analiza uwarunkowań zaangażowania Polaków w kształcenie ustawiczne o charakterze pozaformalnym.....	108
Monika Hamerska: Wykorzystanie metod porządkowania liniowego do tworzenia rankingu jednostek naukowych.....	117
Bartłomiej Jefmański: Zastosowanie modeli IRT w konstrukcji rozmytego systemu wag dla zmiennych w zagadnieniu porządkowania liniowego – na przykładzie metody TOPSIS.....	126
Tomasz Józefowski, Marcin Szymkowiak: Wykorzystanie uogólnionej miary odległości do porządkowania liniowego powiatów województwa podkarpackiego w świetle funkcjonowania specjalnej strefy ekonomicznej Euro-Park Mielec.....	135
Krzysztof Kompa: Zastosowanie testów parametrycznych i nieparametrycznych do oceny sytuacji na światowym rynku kapitałowym przed kryzysem i po jego wystąpieniu.....	144
Mariusz Kubus: Rekurencyjna eliminacja cech w metodach dyskryminacji....	154

Marta Kuc: Wpływ sposobu definiowania macierzy wag przestrzennych na wynik porządkowania liniowego państw Unii Europejskiej pod względem poziomu życia ludności	163
Paweł Lula: Kontekstowy pomiar podobieństwa semantycznego	171
Iwona Markowicz: Model regresji Feldsteina-Horioki – wyniki badań dla Polski	182
Kamila Migdał-Najman: Ocena wpływu wartości stałej Minkowskiego na możliwość identyfikacji struktury grupowej danych o wysokim wymiarze	191
Małgorzata Misztal: O zastosowaniu kanonicznej analizy korespondencji w badaniach ekonomicznych.....	200
Krzysztof Najman: Zastosowanie przetwarzania równoległego w analizie skupień	209
Edward Nowak: Klasyfikacja danych a rachunkowość. Rozważania o relacjach	218
Marcin Pelka: Adaptacja metody <i>bagging</i> z zastosowaniem klasyfikacji pojęciowej danych symbolicznych.....	227
Józef Pocięcha, Mateusz Baryła, Barbara Pawelek: Porównanie skuteczności klasyfikacyjnej wybranych metod prognozowania bankructwa przedsiębiorstw przy losowym i nielosowym doborze prób	236
Agnieszka Przedborska, Małgorzata Misztal: Wybrane metody statystyki wielowymiarowej w ocenie jakości życia słuchaczy uniwersytetu trzeciego wieku	246
Wojciech Roszka: Konstrukcja syntetycznych zbiorów danych na potrzeby estymacji dla małych domen	254
Aneta Rybicka: Połączenie danych o preferencjach ujawnionych i wyrażonych	262
Elżbieta Sobczak: Poziom specjalizacji w sektorach intensywności technologicznej a efekty zmian liczby pracujących w województwach Polski	271
Andrzej Sokołowski, Grzegorz Harańczyk: Modyfikacja wykresu radarowego	280
Marcin Szymkowiak, Marek Witkowski: Wykorzystanie mediany do klasyfikacji banków spółdzielczych według stanu ich kondycji finansowej ..	287
Justyna Wilk, Michał B. Pietrzak, Roger S. Bivand, Tomasz Kossowski: Wpływ wyboru metody klasyfikacji na identyfikację zależności przestrzennych – zastosowanie testu <i>join-count</i>	296
Dorota Witkowska: Wykorzystanie drzew klasyfikacyjnych do analizy zróżnicowania płac w Niemczech	305
Artur Zaborski: Analiza niesymetrycznych danych preferencji z wykorzystaniem modelu punktu dominującego i modelu grawitacji.....	315

Summaries

Krzysztof Jajuga, Józef Pociecha, Marek Walesiak: XXV years of SKAD	24
Beata Basiura, Anna Czapkiewicz: Simulation study of the use of entropy to validation of clustering.....	32
Andrzej Bąk: Problem of choosing the optimal linear ordering procedure in the p_llord package.....	41
Justyna Brzezińska-Grabowska: Latent class analysis in survey research...	50
Grażyna Dehnel: Tax register and social security register as a source of additional information for business statistics – possibilities and limitations.....	59
Sabina Denkowska: Selected methods of assessing the quality of matching in Propensity Score Matching	74
Marta Dziechciarz-Duda, Klaudia Przybysz: Applying the fuzzy set theory to identify the non-monetary factors of poverty.....	83
Iwona Foryś: The potential of the housing market in Poland in the years of economic recessions.....	92
Eugeniusz Gatnar: Statistical analysis of the convergence of CEE countries after 10 years of their membership in the European Union.....	99
Ewa Genge: Trust to the public and financial institutions in the Polish society – an application of latent Markov models.....	107
Alicja Grześkowiak: Multivariate analysis of the determinants of Poles' involvement in non-formal lifelong learning	116
Monika Hamerska: The use of the methods of linear ordering for the creating of scientific units ranking.....	125
Bartłomiej Jefmański: The application of IRT models in the construction of a fuzzy system of weights for variables in the issue of linear ordering – on the basis of TOPSIS method	134
Tomasz Józefowski, Marcin Szymkowiak: GDM as a method of finding a linear ordering of districts of Podkarpackie Voivodeship in the light of the operation of the Euro-Park Mielec special economic zone	143
Krzysztof Kompa: Application of parametric and nonparametric tests to the evaluation of the situation on the world financial market in the pre- and post-crisis period.....	153
Mariusz Kubus: Recursive feature elimination in discrimination methods ...	162
Marta Kuc: The impact of the spatial weights matrix on the final shape of the European Union countries ranking due to the standard of living.....	170
Paweł Lula: The impact of context on semantic similarity.....	181
Iwona Markowicz: Feldstein-Horioka regression model – the results for Poland.....	190

Kamila Migdal-Najman: The assessment of impact value of Minkowski's constant for the possibility of group structure identification in high dimensional data.....	199
Małgorzata Misztal: On the use of canonical correspondence analysis in economic research.....	208
Krzysztof Najman: The application of the parallel computing in cluster analysis.....	217
Edward Nowak: Data classification and accounting. A study of correlations	226
Marcin Pelka: The adaptation of bagging with the application of conceptual clustering of symbolic data.....	235
Józef Pocięcha, Mateusz Baryła, Barbara Pawelek: Comparison of classification accuracy of selected bankruptcy prediction methods in the case of random and non-random sampling technique.....	244
Agnieszka Przedborska, Małgorzata Misztal: Selected multivariate statistical analysis methods in the evaluation of the quality of life of the members of the University of the Third Age.....	253
Wojciech Roszka: Construction of synthetic data sets for small area estimation.....	261
Aneta Rybicka: Combining revealed and stated preference data.....	270
Elżbieta Sobczak: Specialization in sectors of technical advancement vs. effects of workforce number changes in Poland's voivodships.....	279
Andrzej Sokółowski, Grzegorz Harańczyk: Modification of radar plot.....	286
Marcin Szymkowiak, Marek Witkowski: Classification of cooperative banks according to their financial situation using the median.....	295
Justyna Wilk, Michał B. Pietrzak, Roger S. Bivand, Tomasz Kossowski: The influence of classification method selection on the identification of spatial dependence – an application of join-count test.....	304
Dorota Witkowska: Application of classification trees to analyze wages disparities in Germany.....	314
Artur Zaborski: Asymmetric preference data analysis by using the dominance point model and the gravity model.....	323

Beata Basiura, Anna Czapkiewicz

AGH Akademia Górniczo-Hutnicza

e-mail: bbasiura@zarz.agh.edu.pl

SYMULACYJNE BADANIE WYKORZYSTANIA ENTROPII DO BADANIA JAKOŚCI KLASYFIKACJI

Streszczenie: Celem prezentowanej pracy było zweryfikowanie, czy proponowany wskaźnik jakości klasyfikacji wykorzystujący entropię Renyi'ego może być stosowany do określenia liczby podgrup. Badanie przeprowadzono na danych sztucznie wygenerowanych. Zastosowano algorytm aglomeracji Warda, klasyfikacji k -średnich oraz klasyfikacji spektralnej. Na podstawie wybranych czterech indeksów zweryfikowano poprawność wykrywania struktury grupowej. Badano tylko liczbę grup, a nie przynależność poszczególnych elementów do grupy. Otrzymane wyniki potwierdzają skuteczność proponowanego wskaźnika w problemie dotyczącym weryfikacji liczby grup.

Słowa kluczowe: wskaźnik klasyfikacji, entropia Renyi'ego, klasyfikacja.

DOI: 10.15611/pn.2015.384.02

1. Wstęp

W literaturze dotyczącej metod klasyfikacji można znaleźć wiele różnych wskaźników miary jakości klasyfikacji. Wyróżnia się podział na trzy klasy wskaźników oceny jakości grupowania (por. [Halikidi i in. 2001; Baarsch, Celebi 2012; Rendón i in. 2011; Walesiak, Dudek 2012]): wskaźniki oparte na kryteriach zewnętrznych (*external criteria, external validation*), wskaźniki oparte na kryteriach wewnętrznych (*internal criteria, internal validation*) oraz wskaźniki oparte na kryteriach względnych (*relative criteria, relative validation*).

W konstrukcji wskaźników wykorzystuje się miarę odległości (podobieństwa) pomiędzy obiektami. Przykładem takich indeksów są między innymi indeks Calińskiego i Harabasa, indeks Davies-Bouldina czy Silhouette indeks. Wykorzystanie miary zwartości i separowalności zazwyczaj wpływa na lepszą ocenę grup danych o rozkładach eliptycznych.

W prezentowanej pracy rozważono wskaźnik miary klasyfikacji danych, do konstrukcji którego wykorzystano własności entropii Renyi'ego. Pewne własności

tego wskaźnika były prezentowane w pracy Basiury i Czapkiewicz [2014]. W pracy tej między innymi zaprezentowano wstępne badanie symulacyjne, które pokazało, że dla wielowymiarowych rozkładów normalnych wskaźnik ten poprawnie wykrywa strukturę grupową danych.

W niniejszej pracy dokonano dalszej analizy własności wskaźnika opartego na entropii Reny'ego. Celem pracy było porównanie tego wskaźnika z wybranymi wskaźnikami opisanymi w literaturze przy założonej strukturze danych [Milligan Glenn 1981; Halkidi i in. 2010; Rendón i in. 2011; Walesiak, Gatnar 2009; Walesiak 2013]. Podstawą przeprowadzonych badań były symulacyjnie wygenerowane przykłady. Podziału na grupy dokonano przy zastosowaniu algorytmu aglomeracji Warda, podstawowej klasyfikacji k -średnich oraz klasyfikacji spektralnej [Walesiak, Dudek 2012]. Otrzymane wyniki klasyfikacji zweryfikowano, stosując wybrane cztery wskaźniki określające poprawność klasyfikacji.

2. Wybrane wskaźniki klasyfikacji

Do analizy porównawczej wybrano wskaźniki wewnętrzne. Do nich należy między innymi: indeks Calińskiego i Harabasa [1974], indeks Davies-Bouldina [Davies, Bouldin 1979] oraz Silhouette indeks [Rousseeuw 1987].

Indeks Calińskiego i Harabasa (CH) wykorzystuje iloraz zmienności międzygrupowej (Between Groups – BG) oraz zmienności wewnątrzgrupowej (Within Groups – WG). Zmienność międzygrupowa jest ważoną sumą kwadratów odległości pomiędzy środkiem każdej klasy a środkiem całego zbioru. Wagami są wielkości analizowanych klas. Natomiast zmienność wewnątrzgrupowa wyznaczana jest jako suma kwadratów odległości każdego elementu podzbioru od środka klasy.

W indeksie Davies-Bouldina (DB) dla każdego skupienia wyznacza się średnią odległość pomiędzy każdym punktem grupy a jej centrum (oznaczymy je jako δ_k i $\delta_{k'}$) oraz odległość pomiędzy środkami skupienia k i skupienia k' (oznaczymy jako $\Delta_{kk'}$). Następnie dla każdego podzbioru wyznacza się maksymalną wartość zu: $\frac{\delta_k + \delta_{k'}}{\Delta_{kk'}}$. Indeks DB to średnia wartość z ilorazów po wszystkich podzbiórach.

Silhouette indeks zaproponowany został przez Rousseeuwa (1987). Indeks ten pozwala oceniać prawidłowość zaklasyfikowania poszczególnych obiektów do wyodrębnionych klas na podstawie następującej reguły. Niech:

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (1)$$

Rozważmy dwie klasy: C_k oraz $C_{k'}$. Wielkość $a(i)$ określa średnią odległość obiektu i od pozostałych obiektów należących do klasy C_k , natomiast wielkość $b(i)$ określa minimalną odległość i od obiektów należących do klasy $C_{k'}$. Indeks $S(P_r)$ określający prawidłowość wyodrębnienia klasy jest średnią z wartości $S(i)$ dla po-

szczególnych elementów klasy, natomiast średnia z $S(P_r)$ po wszystkich klasach jest indeksem sylwetkowym.

O lepszej jakości klasyfikacji mówią wyższe wartości indeksu Calińskiego i Harabasa, wyższe wartości indeksu Silhouette oraz niższe indeksy Davies-Bouldina.

3. Wskaźnik klasyfikacji na podstawie miary entropii

Pojęcie entropii wprowadził Shannon w 1948 r., następnie w drugiej połowie ubiegłego wieku pojawiło się wiele uogólnień probabilistycznej miary tej entropii. Węgierski matematyk Alfred Rényi [Rényi 1961] zaproponował następujące uogólnienie pojęcia entropii:

$$H(x) = \frac{1}{1-\alpha} \log(\int f^\alpha(x) dx), \quad \alpha > 0, \alpha \neq 1. \quad (2)$$

W szczególności dla $\alpha = 2$ otrzymuje się:

$$H(x) = -\log(\int f^2(x) dx). \quad (3)$$

Niech $\{x_1, \dots, x_N\}$, gdzie x_i jest d -wymiarowym obiektem, będzie zbiorem danych niezależnych o tym samym rozkładzie $f(x)$. Jeśli nie znamy rozkładu danej funkcji, to do jej estymacji można zastosować metodę nieparametryczną na podstawie estymacji jądrowej [Liang i in. 2011; Jensen i in. 2003]. Niech:

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^N W_{\sigma^2}(x, x_i).$$

Przy czym $W_{\sigma^2}(x, x_i)$ jest oknem Parzena, natomiast σ^2 określa szerokość okna w zależności od rozmiaru danych. W naszych badaniach została wykorzystana funkcja jądrowa Gaussa, określona wzorem (4), w którym parametr σ oznacza optymalnej wielkości okno:

$$W_{\sigma^2}(x, x_i) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{(x-x_i)^T(x-x_i)}{2\sigma^2}\right). \quad (4)$$

Można pokazać, że entropię całego układu można wyznaczyć jako:

$$H = -\log \frac{1}{N^2} \sum_{j=1}^N \sum_{i=1}^N W_{2\sigma^2}(x_j, x_i).$$

Niech dane będą podzielone na K klastrow: C_k dla $k = 1, 2, \dots, K$, w których pojawia się N_k obiektów. Entropię w k -tym klastrze można zdefiniować jako:

$$H(C_k) = -\log \frac{1}{N_k^2} \sum_{j=1}^{N_k} \sum_{i=1}^{N_k} W_{2\sigma^2}(x_j, x_i). \quad (5)$$

Wskaźnik postaci:

$$V(C_1, C_2, \dots, C_K) = \sum_{i=1}^K \frac{N_k}{N} H(C_K) \quad (6)$$

mógły być interpretowany jako wskaźnik entropii wewnątrzgrupowej. W literaturze pojawiła się taka ważona suma entropii w każdym klastrze, ale wyznaczana dla danych dyskretnych [Rendón i in. 2011]. Stosując to rozumowanie dla entropii Reny'ego, otrzymujemy (6). Entropię pomiędzy grupami zdefiniujemy jako:

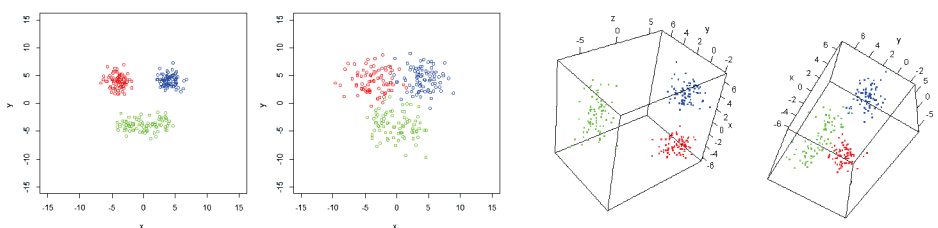
$$H(C_1, C_2, \dots, C_K) = -\log \frac{1}{2 \prod_{k=1}^K N_k} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N W_{2\sigma^2}(x_j, x_i). \quad (7)$$

Jeśli klasy są dobrze utworzone, wskaźnik ten powinien być duży [Jensen i in. 2003]. Jako wskaźnik klasyfikacji można rozważyć zatem następujący iloraz:

$$E = \frac{H(C_1, C_2, \dots, C_K)}{V(C_1, C_2, \dots, C_K)}. \quad (8)$$

4. Badanie empiryczne

W pierwszym eksperymencie symulacyjnym rozważane były dane pochodzące z wielowymiarowego rozkładu normalnego. W przypadku rozkładów dwu- i trójwymiarowych zadeklarowano klasyfikacje na trzy grupy. Liczba elementów w poszczególnych grupach w pierwszym eksperymencie była taka sama, ale zbadano też przypadki różnej liczebności klas. Grupy różniły się przede wszystkim środkami ciężkości. Przy generowaniu danych wielowymiarowych rozważano 20-, 30- i 36- wymiarowe wektory z podziałem na dwa, trzy i cztery skupienia. Badanie przeprowadzono dla różnych liczebności prób (20, 40 i 60 elementów). W każdym eksperymencie tak dobrano środki ciężkości klas i rozrzut elementów, aby rozważane zbiory były mniej lub bardziej separowalne. Opis modeli symulacyjnych prezentuje tab. 1, a przykładowe dane zaprezentowano na rys. 1.



Rys. 1. Przykłady danych symulacyjnych dwuwymiarowych i trójwymiarowych

Źródło: opracowanie własne.

W tabeli 2 porównane zostały wyniki uzyskane przy zastosowaniu klasyfikacji hierarchicznej metodą Warda, klasyfikacji k -średnich oraz klasyfikacji spektralnej z wybraną miarą odległości euklidesowej dla eksperymentów o numerach 1, 2, 3 i 4 z tab. 1. Przy stosunkowo słabo separowalnych grupach otrzymuje się niski procent poprawnie określonych podgrup. Na tle badanych wskaźników indeks E wypada korzystnie w grupie 1.

Tabela 1. Charakterystyka modeli symulacyjnych danych o wielowymiarowym rozkładzie normalnym

Nr	Liczba klas	Liczba zmiennych	Liczba pomiarów	Środki ciężkości klas	Macierz kowariancji
1	3	2	100	(-6,4) (6, -4) (6, 6)	$\begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}, \begin{bmatrix} 8 & 0.9 \\ 0.9 & 8 \end{bmatrix}$
2	3	2	50	(0, 0), (1.5, 7), (3, 14)	$\begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}, \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$
3	3	3	30,50,50	(-4,4,0), (4,4,0), (0,-4,5)	macierz jednostkowa
4	3	3		(1.5,6,-3), (3,12,-6), (4.5,18,-9)	$\begin{bmatrix} 1 & -0.9 & -0.9 \\ -0.9 & 1 & 0.9 \\ -0.9 & 0.9 & 1 \end{bmatrix}$
5	2	20	20, 40	(4,4,...4,0,...,0), (0,...0,4,...,4)	macierz jednostkowa
6	3	30	60	(4,...,4,0,...,0,0,...,0) (0,...,0,4,...,4,0,...,0) (0,...,0,0,...,0,4,...,4)	macierz jednostkowa
7	4	36	36	(4,...,4,0,...,0,0,...,0,0,...,0), (0,...,0,4,...,4,0,...,0,0,...,0) (0,...,0,0,...,0,4,...,4,0,...,0) (0,...,0,0,0,...,0,0,...,0,4,...,4)	macierz jednostkowa

W punkcie 2) wykorzystano model 13 z pakietu clusterSim [Walesiak 2013], natomiast w punkcie 4) model 5 z tego pakietu; dane z punktów 5, 6 i 7 były generowane przy założeniu niezależności.

Źródło: opracowanie własne.

Tabela 2. Liczba poprawnie określonych podgrup (w procentach)

Indeks	Metody klasyfikacji											
	metoda Warda				metoda <i>k</i> -średnich				metoda spektralna			
	Nr 1	Nr 2	Nr 3	Nr 4	Nr 1	Nr 2	Nr 3	Nr 4	Nr 1	Nr 2	Nr 3	Nr 4
CH	45	13	98	68	16	15	99	76	13	17	98	65
DB	15	15	95	76	14	25	97	65	52	14	99	62
S	17	10	98	85	5	14	96	82	3	15	95	83
E	45	14	99	69	55	16	86	67	56	16	98	67

CH – indeks Calińskiego i Harabasa; DB – indeks Davies-Bouldina; S – Silhouette indeks; E – indeks entropii; nr to numer modelu z tab. 1.

Źródło: opracowanie własne.

Wyniki eksperymentów o numerach 5, 6 i 7 nie zostały zestawione w tabeli, ponieważ dotyczyły bardzo dobrze separowalnych danych, których grupy są bardzo dobrze wykrywane przez wszystkie wskaźniki. Liczba poprawnie określonych podgrup wahała się od 96% do 100%. Także wskaźnik *E* na równi z pozostałymi trzema we wszystkich przypadkach dawał poprawne wyniki w 97-99%. Zatem można zauważyć, że indeks klasyfikacji konstruowany na podstawie entropii Reny'ego (*E*) ma bardzo dobre własności. Jego skuteczność w wykrywaniu liczby klas jest porównywalna ze skutecznością klasycznych indeksów.

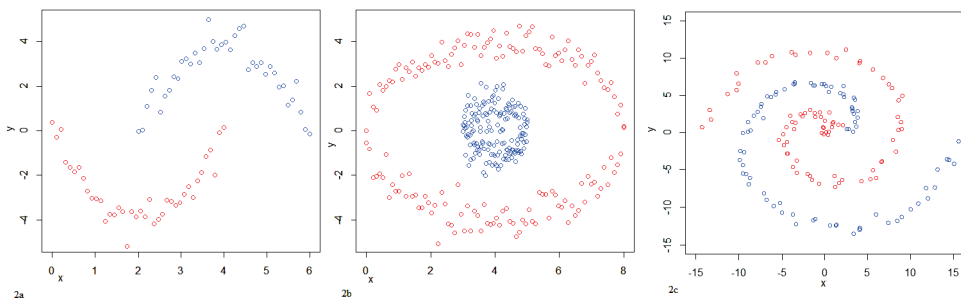
W drugim zestawie modeli symulacyjnych (tab. 3) rozważono niekorelowane dane o rozkładzie skośnym t -Studenta. Przyjęto dwa parametry skośności $\eta = 1.7$ i $\eta = 0.6$. W ten sposób osłabione zostało założenie o eliptyczności rozkładów.

Tabela 3. Charakterystyka modeli symulacyjnych danych o rozkładach brzegowych skośnych

Nr	Liczba klas	Liczba zmiennych	Liczba pomiarów	Środki ciężkości klas	Macierz kowariancji i przyjęte stopnie swobody i parametry skośności
8	2	20	20, 40	$(1,1,\dots,1,0,\dots,0),$ $(0,\dots,0,1,\dots,1)$	macierz jednostkowa, $df = 5, \eta = 1.7$
9	2	20	20, 40	$(4,4,\dots,4,0,\dots,0),$ $(0,\dots,0,4,\dots,4)$	macierz jednostkowa, $df = 5, \eta = 0.6$
10	3	30	60	$(1,\dots,1,0,\dots,0,0,\dots,0)$ $(0,\dots,0,1,\dots,1,0,\dots,0)$ $(0,\dots,0,0,\dots,0,1,\dots,1)$	macierz jednostkowa, $df = 5, \eta = 1.7$
11	3	30	60	$(4,\dots,4,0,\dots,0,0,\dots,0)$ $(0,\dots,0,4,\dots,4,0,\dots,0)$ $(0,\dots,0,0,\dots,0,4,\dots,4)$	macierz jednostkowa, $df = 5, \eta = 0.6$
12	4	36	36	$(1,\dots,1,0,\dots,0,0,\dots,0,0,\dots,0),$ $(0,\dots,0,1,\dots,1,0,\dots,0,0,\dots,0)$ $(0,\dots,0,0,\dots,0,1,\dots,1,0,\dots,0)$ $(0,\dots,0,0,\dots,0,0,\dots,0,1,\dots,1)$	macierz jednostkowa, $df = 5, \eta = 1.7$
13	4	36	36	$(4,\dots,4,0,\dots,0,0,\dots,0,0,\dots,0),$ $(0,\dots,0,4,\dots,4,0,\dots,0,0,\dots,0)$ $(0,\dots,0,0,\dots,0,4,\dots,4,0,\dots,0)$ $(0,\dots,0,0,\dots,0,0,\dots,0,4,\dots,4)$	macierz jednostkowa, $df = 5, \eta = 0.6$

Dane generowane z założeniem niezależności i uwzględnieniem skośności ($\eta = 0.6, \eta = 1.7$) w rozkładzie brzegowym t -Studenta.

Źródło: opracowanie własne.



Rys. 2. Przykładowe wygenerowane dane nieklasyczne: 2a – dane typu Worms, 2b – dane typu koła, 2c – dane typu spirala

Źródło: opracowanie własne.

W wyniku eksperymentów od 8 do 13 liczba poprawnie określonych podgrup wahała się od 90% do 100% dla wszystkich wskaźników. W eksperymencie nr 8 i 10 najlepszym wynikiem było 96% (dla E i S w metodzie Warda), przy czym pozostałe wartości były niewiele niższe. W pozostałych grupach eksperymentów symulacyjnych (9,11,12,13) wszystkie wskaźniki, także wskaźnik E , dawały poprawne wyniki w 97-99%.

W trzecim eksperymencie weryfikowano wartości wskaźników na podstawie nieklasycznych zbiorów danych. Dane te były losowo zaburzane w taki sposób, aby nie utraciły swojej struktury grupowej. Przykładowe wygenerowane zbiory danych przedstawia rys. 2.

Niestety w tych eksperymentach symulacyjnych wszystkie cztery badane wskaźniki wypadają słabo. Jedynie dane typu „Worms” są dobrze oceniane, bowiem w tym przypadku liczba poprawnie określonych podgrup wynosiła od 78% do 98%. Najlepiej wypadła metoda aglomeracyjna Warda i metoda spektralna. Wyniki poprawnie określonej liczby podgrup dla danych typu dwa koła wynosiły od 18% do 65%. Najlepiej w tym przypadku wypadła metoda spektralna: wskaźnik CH – 65%, a wskaźnik S – 56%. Badany wskaźnik E tylko w ok. 18% poprawnie wykrywała podział na dwie grupy przy wszystkich rozważanych metodach klasyfikacji.

Wybór dwóch klas przy danych typu „spirala” potwierdzał się tylko w 10% dla wszystkich badanych indeksów.

5. Zakończenie

W pracy przedstawiono badanie przydatności wskaźnika jakości klasyfikacji opartego na własnościach entropii Reny’ego. Uzyskane wyniki dla wybranych danych symulacyjnych wskazują na dużą przydatność tego wskaźnika do określenia liczby klas. Indeks E generalnie jest podobny do innych klasycznych wskaźników. Badanie symulacyjne pokazało, że dla dobrze separowalnych wskaźnik E , na równi z innymi indeksami, poprawnie wykrywa strukturę grupową danych. Jednakże w przypadku danych o bardzo rozbudowanej strukturze indeks ten, na równi z pozostałymi indeksami, nie wybiera właściwie liczby klas.

Wydaje się zatem, że wskaźnik E mógłby być wykorzystany jako miara jakości klasyfikacji na równi z pozostałymi klasycznym wskaźnikami. Może służyć jako informacja wspomagająca wybór właściwej decyzji.

Literatura

- Baarsch J., Celebi M.C., 2012, *Investigation of Internal Validity Measures for K-Means Clustering*, IMECS 2012, Hong Kong.
- Basiura B., Czapkiewicz A., 2014, *Badanie jakości klasyfikacji szeregów czasowych*, PN 327 Taksonomia 22, *Klasyfikacja i analiza danych – teoria i zastosowania*, Jajuga K., Walesiak M. (red.), Uniwersytet Ekonomiczny we Wrocławiu.

- Davies D., Bouldin D., 1979, *A cluster separation measure*, IEEE Transactions on Pattern Analysis and Machine Intelligence 1(2): 224-227.
- Calinski R.B., Harabasz J., 1974, *A dendrite method for cluster analysis*, Communications in Statistics – Theory and Methods 3(1): 1-27.
- Halkidi M., Yanniss B., Vazirgiannis M., 2001, *On clustering validation techniques*, Journal of Intelligent Information Systems, 17:2/3, 107-145.
- Jenssen R., Hild K.E., Erdogmus D., Principe J.C., Eltoft T., 2003, *Clusterin using renyi's entropy*, Neural Networks, Proceedings of the International Joint Conference on (Volume:1).
- Liang J., Zhao X., Li D., Cao F., Dang C., 2011, *Determining the number of clusters using information entropy for Mixed Data*, Patter Recognition, v. 45, 2251-2265.
- Milligan G., Glenn W., 1981, *A Monte Carlo study of thirty internal criterion measures for cluster analysis*, Psychometrika 46(2): 187-199.
- Rendón E., Abundez I., Arizmendi A., Quiroz E.M., 2011, *Internal versus external cluster validation indexes*, Intenational Journal of Computers and Communications, no. 1, vol. 5.
- Rényi A., 1961, *On measures of information and entropy*, Proceedings of the fourth Berkeley Symposium on Mathematics, Statistics and Probability 1960, pp. 547-561.
- Rousseeuw P.J., 1987 *Silhouettes: A graphic aid to the interpretation and validation of cluster analysis*, Journal of Computational and Applied Mathematics 20(1): 53-65.
- Walesiak M., Gatnar E., 2009, *Statystyczna analiza danych z wykorzystaniem programu R*, PWN, Warszawa.
- Walesiak M., 2013, *Zagadnienie doboru liczby klas w klasyfikacji spektralnej*, Taksonomia 20 – Klasyfikacja i analiza danych – teoria i zastosowania, UE, Wrocław.
- Walesiak M., Dudek M., 2012, *Package 'clusterSim' in R project*, <http://keii.ue.wroc.pl/clusterSim/index.html>.
- Wędrowska E., 2011, *Wykorzystanie entropii Shanona i jej uogólnień do badania rozkładu prawdopodobieństwa zmiennej losowej dyskretnej*, Przegląd Statystyczny, RLVII, zeszyt 4.
- R Development Core Team (2005). *R: A language and environment for statistical computing, reference index version 2.12.2 (2011-02-25)* R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>.
- Kovács F., Legány C., & Babos A., 2005, *Cluster validity measurement techniques. Proceedings of the 6th International Symposium of Hungarian Researchers on Computational Intelligence*, Budapest, Nov. 2005, 18-19.

SIMULATION STUDY OF THE USE OF ENTROPY TO VALIDATION OF CLUSTERING

Summary: The aim of this paper is to present a quality index classification using Renyi entropy against known quality indicators grouping of multidimensional time series. The study was conducted on artificially generated data and empirical data. The division into groups was made by using Ward's agglomeration algorithm, *k*-means method's and spectral clustering. The results were verified using the selected indices of clustering validation.

Keywords: clustering validation, Renyi's entropy, clustering.