

PRACE NAUKOWE

Uniwersytetu Ekonomicznego we Wrocławiu

RESEARCH PAPERS

of Wrocław University of Economics

Nr 327

Taksonomia 22

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

Redaktorzy naukowci

Krzysztof Jajuga, Marek Walesiak



Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu
Wrocław 2014

Redaktor Wydawnictwa: Barbara Majewska

Redaktor techniczny: Barbara Łopusiewicz

Korektor: Barbara Cibis

Łamanie: Beata Mazur

Projekt okładki: Beata Dębska

Publikacja jest dostępna w Internecie na stronach:

www.ibuk.pl, www.ebscohost.com,

w Dolnośląskiej Bibliotece Cyfrowej www.dbc.wroc.pl,

The Central and Eastern European Online Library www.ceeol.com,

a także w adnotowanej bibliografii zagadnień ekonomicznych BazEkon

http://kangur.uek.krakow.pl/bazy_ae/bazekon/nowy/index.php

Informacje o naborze artykułów i zasadach recenzowania znajdują się

na stronie internetowej Wydawnictwa

www.wydawnictwo.ue.wroc.pl

Tytuł dofinansowany ze środków Narodowego Banku Polskiego

oraz ze środków Sekcji Klasyfikacji i Analizy Danych PTS

Kopiowanie i powielanie w jakiegokolwiek formie

wymaga pisemnej zgody Wydawcy

© Copyright by Uniwersytet Ekonomiczny we Wrocławiu

Wrocław 2014

ISSN 1899-3192 (Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu)

ISSN 1505-9332 (Taksonomia)

Wersja pierwotna: publikacja drukowana

Druk: Drukarnia TOTEM

Spis treści

Wstęp	9
Eugeniusz Gatnar , Balance of payments statistics and external competitiveness of Poland.....	15
Andrzej Sokolowski, Magdalena Czaja , Efektywność metody k -średnich w zależności od separowalności grup.....	23
Barbara Pawelek, Józef Pocięcha, Adam Sagan , Wielosektorowa analiza ukrytych przejść w modelowaniu zagrożenia upadłością polskich przedsiębiorstw	30
Elżbieta Gołata , Zróżnicowanie procesu starzenia i struktur demograficznych w Poznaniu i aglomeracji poznańskiej na tle wybranych dużych miast Polski w latach 2002-2011.....	39
Aleksandra Łuczak, Feliks Wysocki , Ustalanie systemu wag dla cech w zagadnieniach porządkowania liniowego obiektów	49
Marek Walesiak , Wzmacnianie skali pomiaru dla danych porządkowych w statystycznej analizie wielowymiarowej	60
Paweł Lula , Identyfikacja słów i fraz kluczowych w tekstach polskojęzycznych za pomocą algorytmu <i>RAKE</i>	69
Mariusz Kubus , Propozycja modyfikacji metody złagodzonego LASSO.....	77
Andrzej Bąk, Tomasz Bartłomowicz , Wielomianowe modele logitowe wyborów dyskretnych i ich implementacja w pakiecie <i>DiscreteChoice</i> programu R.....	85
Justyna Brzezińska , Wykorzystanie modeli logarytmiczno-liniowych do analizy bezrobocia w Polsce w latach 2004-2012.....	95
Andrzej Bąk, Marcin Pelka, Aneta Rybicka , Zastosowanie pakietu <i>dcMNM</i> programu R w badaniach preferencji konsumentów wódki	104
Barbara Batóg, Jacek Batóg , Analiza stabilności klasyfikacji polskich województw według sektorowej wydajności pracy w latach 2002-2010	113
Małgorzata Markowska, Danuta Strahl , Klasyfikacja europejskiej przestrzeni regionalnej ze względu na filary inteligentnego rozwoju z wykorzystaniem referencyjnego systemu granicznego.....	121
Kamila Migdał-Najman, Krzysztof Najman , Formalna ocena jakości odwzorowania struktury grupowej na mapie Kohonena	131
Kamila Migdał-Najman, Krzysztof Najman , Graficzna ocena jakości odwzorowania struktury grupowej na mapie Kohonena	139
Beata Basiura, Anna Czapkiewicz , Badanie jakości klasyfikacji szeregów czasowych	148
Michał Trzęsiok , Wybrane metody identyfikacji obserwacji oddalonych.....	157

Grażyna Dehnel, Tomasz Klimanek , Taksonomiczne aspekty estymacji pośredniej uwzględniającej autokorelację przestrzenną w statystyce gospodarczej.....	167
Michał Bernard Pietrzak, Justyna Wilk , Odległość ekonomiczna w modelowaniu zjawisk przestrzennych z wykorzystaniem modelu grawitacji.....	177
Maciej Beręsewicz , Próba zastosowania różnych miar odległości w uogólnionym estymatorze Petersena	186
Marcin Szymkowiak, Tomasz Józefowski , Konstrukcja i praktyczne wykorzystanie estymatorów typu SPREE na przykładzie dwuwymiarowych tabel kontyngencji	195
Marcin Pelka , Klasyfikacja pojęciowa danych symbolicznych w podejściu wielomodelowym	202
Małgorzata Machowska-Szewczyk , Ocena klas w rozmytej klasyfikacji obiektów symbolicznych.....	210
Justyna Wilk , Problem wyboru liczby klas w taksonomicznej analizie danych symbolicznych.....	220
Andrzej Dudek , Metody analizy skupień w klasyfikacji markerów map Google	229
Ewa Roszkowska , Ocena ofert negocjacyjnych w słabo ustrukturyzowanych problemach negocjacyjnych z wykorzystaniem rozmytej procedury SAW	237
Marcin Szymkowiak, Marek Witkowski , Zastosowanie analizy korespondencji do badania kondycji finansowej banków spółdzielczych.....	248
Bartłomiej Jefmański , Budowa rozmytych indeksów satysfakcji klientów z zastosowaniem programu R.....	257
Karolina Bartos , Odkrywanie wzorców zachowań konsumentów za pomocą analizy koszykowej danych transakcyjnych	266
Joanna Trzęsiok , Taksonomiczna analiza krajów pod względem dzietności kobiet oraz innych czynników demograficznych	275
Beata Bal-Domańska , Próba identyfikacji większych skupisk regionalnych oraz ich konwergencja.....	285
Beata Bieszk-Stolorz, Iwona Markowicz , Wpływ zasiłku na proces poszukiwania pracy	294
Marta Dziechciarz-Duda, Klaudia Przybysz , Wykształcenie a potrzeby rynku pracy. Klasyfikacja absolwentów wyższych uczelni.....	303
Tomasz Klimanek , Problem pomiaru procesu dezagrarnizacji wsi polskiej w świetle wielowymiarowych metod statystycznych.....	313
Małgorzata Sej-Kolasa, Mirosława Sztemberg-Lewandowska , Wybrane metody analizy danych wzdluznych.....	321
Artur Zaborski , Zastosowanie miar odległości dla danych porządkowych do agregacji preferencji indywidualnych	330
Mariola Chrzanowska, Nina Drejerska, Iwona Pomianek , Zastosowanie analizy korespondencji do badania sytuacji mieszkańców strefy podmiejskiej Warszawy na rynku pracy.....	338

Katarzyna Wawrzyniak , Klasyfikacja województw według stopnia realizacji priorytetów Strategii Rozwoju Kraju 2007-2015 z wykorzystaniem wartości centrum wierszowego	346
---	-----

Summaries

Eugeniusz Gatnar , Statystyka bilansu płatniczego a konkurencyjność gospodarki Polski	22
Andrzej Sokółowski, Magdalena Czaja , Cluster separability and the effectiveness of k -means method	29
Barbara Pawelek, Józef Pocięcha, Adam Sagan , Multisectoral analysis of latent transitions in bankruptcy prediction models.....	38
Elżbieta Golata , Differences in the process of aging and demographic structures in Poznań and the agglomeration compared to selected Polish cities in the years 2002-2011	48
Aleksandra Łuczak, Feliks Wysocki , Determination of weights for features in problems of linear ordering of objects	59
Marek Walesiak , Reinforcing measurement scale for ordinal data in multivariate statistical analysis	68
Paweł Lula , Automatic identification of keywords and keyphrases in documents written in Polish.....	76
Mariusz Kubus , The proposition of modification of the relaxed LASSO method.....	84
Andrzej Bąk, Tomasz Bartłomowicz , Microeconomic multinomial logit models and their implementation in the <code>DiscreteChoice</code> R package .	94
Justyna Brzezińska , The analysis of unemployment data in Poland in 2004-2012 with application of log-linear models	103
Andrzej Bąk, Marcin Pelka, Aneta Rybicka , Application of the MMLM package of R software for vodka consumers preference analysis.....	112
Barbara Batóg, Jacek Batóg , Analysis of the stability of classification of Polish voivodeships in 2002-2010 according to the sectoral labour productivity	120
Małgorzata Markowska, Danuta Strahl , Classification of the European regional space in terms of smart growth pillars using the reference limit system.....	130
Kamila Migdał Najman, Krzysztof Najman , Formal quality assessment of group structure mapping on the Kohonen's map	138
Kamila Migdał Najman, Krzysztof Najman , Graphical quality assessment of group structure mapping on the Kohonen's map	147
Beata Basiura, Anna Czapkiewicz , Validation of time series clustering	156
Michał Trzęsiok , Selected methods for outlier detection.....	166

Grażyna Dehnel, Tomasz Klimanek , Taxonomic aspects of indirect estimation accounting for spatial correlation in enterprise statistics	176
Michał Bernard Pietrzak, Justyna Wilk , Economic distance in modeling spatial phenomena with the application of gravity model.....	185
Maciej Beręsewicz , An attempt to use different distance measures in the Generalized Petersen estimator	194
Marcin Szymkowiak, Tomasz Józefowski , Construction and practical using of SPREE estimators for two-dimensional contingency tables.....	201
Marcin Pelka , The ensemble conceptual clustering for symbolic data.....	209
Małgorzata Machowska-Szewczyk , Evaluation of clusters obtained by fuzzy classification methods for symbolic objects.....	219
Justyna Wilk , Problem of determining the number of clusters in taxonomic analysis of symbolic data	228
Andrzej Dudek , Clustering techniques for Google maps markers.....	236
Ewa Roszkowska , The evaluation of negotiation offers in ill structure negotiation problems with the application of fuzzy SAW procedure	247
Marcin Szymkowiak, Marek Witkowski , The use of correspondence analysis in analysing the financial situation of cooperative banks.....	256
Bartłomiej Jefmański , The construction of fuzzy customer satisfaction indexes using R program.....	265
Karolina Bartos , Discovering patterns of consumer behaviour by market basket analysis of the transactional data.....	274
Joanna Trzęsiok , Cluster analysis of countries with respect to fertility rate and other demographic factors	284
Beata Bal-Domańska , An attempt to identify major regional clusters and their convergence	293
Beata Bieszk-Stolorz, Iwona Markowicz , The influence of benefit on the job finding process	302
Marta Dziechciarz-Duda, Klaudia Przybysz , Education and labor market needs. Classification of university graduates	312
Tomasz Klimanek , The problem of measuring deagrarianisation process in rural areas in Poland using multivariate statistical methods.....	320
Małgorzata Sej-Kolasa, Mirosława Sztemberg-Lewandowska , Selected methods for an analysis of longitudinal data.....	329
Artur Zaborski , The application of distance measures for ordinal data for aggregation individual preferences	337
Mariola Chrzanowska, Nina Drejerska, Iwona Pomianek , Application of correspondence analysis to examine the situation of the inhabitants of Warsaw suburban area in the labour market	345
Katarzyna Wawrzyniak , Classification of voivodeships according to the level of the realization of priorities of <i>the National Development Strategy 2007-2015</i> with using the values of centroid of the rows	355

Beata Basiura, Anna Czapkiewicz

AGH Kraków

BADANIE JAKOŚCI KLASYFIKACJI SZEREGÓW CZASOWYCH

Streszczenie: Celem niniejszej pracy było zaprezentowanie wskaźnika jakości klasyfikacji z zastosowaniem entropii Renyiego na tle znanych wskaźników jakości grupowania wielowymiarowych szeregów czasowych. Punktem wyjścia były dane empiryczne. Podziału na grupy dokonano przy zastosowaniu algorytmu aglomeracji Warda, klasyfikacji k-średnich oraz klasyfikacji spektralnej. Otrzymane wyniki klasyfikacji zweryfikowano stosując wybrane indeksy określające poprawność klasyfikacji. Zaproponowany wskaźnik wydaje się obiecujący, ale wymaga zweryfikowania dla różnych rozkładów badanych szeregów.

Słowa kluczowe: wskaźnik klasyfikacji, entropia Renyiego, klasyfikacja szeregów czasowych.

1. Wstęp

Empiryczne badanie jakości struktury grupowej danych jest zadaniem niezwykle trudnym. W literaturze przedmiotu można znaleźć wiele różnych wskaźników miary jakości klasyfikacji. Wyróżnia się metody wzorcowe, w których uzyskany podział na grupy porównuje się z pewnym podziałem wzorcowym, oraz metody bezwzorcowe, w których ocena jakości klasyfikacji wykorzystuje jedynie informacje zawarte w danych. Większość wskaźników opiera się na miarach zwartości klas takich, jak na przykład wariancja wewnątrzgrupowa, oraz na miarach separowalności poszczególnych podgrup określających zmienność międzygrupową.

Istnieją trzy klasy wskaźników oceny jakości grupowania (por. Halkidi 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*). Większość z nich wykorzystuje miarę odległości pomiędzy obiektami i zazwyczaj lepiej ocenia grupy „eliptyczne”, natomiast w przypadku nietypowych podzbiorów danych wydaje się, że nie dają one wiarygodnych wyników.

Wszelkie decyzje podejmowane przez badaczy, związane z inwestycjami, rozwojem, wielkością sprzedaży i ceną, mogą być oparte na wartości szeregu czaso-

wego. Ocena wielkości zmian rynkowych i aspekty, do której grupy zaklasyfikować badane szeregi, zależy nie tylko od wyboru miary podobieństwa czy metody grupowania. Ostateczna decyzja dotycząca wyników grupowania powinna być podejmowana na podstawie oceny jakości uzyskanej klasyfikacji. W niniejszej pracy zaproponowano ocenę jakości klasyfikacji krótkich szeregów czasowych, wykorzystującą entropię Rényiego [Rényi 1961; Wędrowska 2010].

Prezentowana praca ma na celu porównanie wybranych wskaźników jakości grupowania wielowymiarowych szeregów czasowych. Punktem wyjścia były dane empiryczne. Podziału na grupy dokonano przy zastosowaniu algorytmu aglomeracji Warda, klasyfikacji k -średnich oraz klasyfikacji spektralnej. Otrzymane wyniki klasyfikacji zweryfikowano stosując wybrane wskaźniki określające poprawność klasyfikacji. Zaproponowany wskaźnik, wykorzystujący funkcję entropii, został przedstawiony na tle istniejących już w literaturze wskaźników klasyfikacji. Dla zbadania jakości prezentowanego wskaźnika wykonano krótkie badanie symulacyjne.

Celem pracy było pokazanie wskaźnika skonstruowanego na podstawie entropii Rényiego w stosunku do innych opisanych w literaturze i stosowanych wskaźników jakości klasyfikacji przy założeniu nieznanego rozkładu badanych szeregów [Milligan, Glenn 1981; Halkidi i in. 2010; Rendón i in. 2011; Walesiak, Gątnar (red.) 2009].

2. Wybrane wskaźniki klasyfikacji

Jak już wspomniano, wyróżnia się trzy klasy wskaźników jakości klasyfikacji. Wskaźniki zewnętrzne, w których uzyskana struktura porównywana jest z założoną, znaną z góry – ekspercką – strukturą danych, wskaźniki wewnętrzne, wykorzystujące jedynie informacje z analizowanego zbioru danych, oraz wskaźniki względne, gdy ocena struktury porównywana jest z grupowaniem uzyskanym za pomocą tego samego algorytmu, ale z założonymi innymi parametrami (np. inna liczba grup). Większość wskaźników opiera się na miarach spójności (zwartości) skupień: elementy każdego skupienia powinny odpowiednio blisko siebie (*cluster cohesion*, *cluster compactness*) oraz na miarach rozdzielenia skupień: klastry powinny być odpowiednio od siebie oddalone (*cluster separation*). Popularną miarą zwartości jest wariancja.

W niniejszej pracy wybrane zostały trzy najczęściej używane, wewnętrzne współczynniki: indeks Calińskiego i Harabasa [Calinski, Harabasz 1974], indeks Daviesa-Bouldina [Davies, Bouldin 1979] oraz Silhouette indeks [Rousseeuw 1987].

Pierwszy z nich oparty jest na zmienności międzygrupowej (*between groups* – BG) oraz zmienności wewnątrz grupowej (*within groups* – WG). Zmienność międzygrupowa jest ważoną sumą kwadratów odległości pomiędzy środkiem każdej klasy a środkiem całego zbioru. Wagami są wielkości klastrów.

$$CH = \frac{tr(BG) / (K - 1)}{tr(WG) / (N - K)}. \quad (1)$$

Zmienność wewnątrzgrupowa wyznaczana jest jako suma kwadratów odległości każdego elementu podzbioru od środka klasy.

Wybrany drugi indeks także wykorzystuje iloraz zwartości i separowalności klasy. Dla każdego skupienia wyznacza się w nim średnią odległość pomiędzy każdym punktem grupy a jej centrum (oznaczymy je jako δ_k i $\delta_{k'}$). Oznaczmy ko $\Delta_{kk'}$ odległość pomiędzy środkami skupienia k i skupienia k' . Następnie dla każdego podzbioru wyznacza się maksymalną wartość ilorazu: $\frac{\delta_k + \delta_{k'}}{\Delta_{kk'}}$. Wtedy indeks DB oblicza się według wzoru (2):

$$DB = \frac{1}{K} \sum_{k=1}^K \max_{k \neq k'} \frac{\delta_k + \delta_{k'}}{\Delta_{kk'}}. \quad (2)$$

Trzeci z indeksów, zaproponowany przez Rousseeuwa [1987], pozwala oceniać prawidłowość zaklasyfikowania poszczególnych obiektów do wyodrębnionych klas. Dany jest wzorem (3):

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i); b(i)\}}. \quad (3)$$

Określa się, że $a(i)$ to średnia odległość obiektu i od pozostałych obiektów należących do klasy C_k , $b(i)$ to minimalna ze średnich odległości obiektu i od obiektów należących do klasy $C_{k'}$, co zapisuje się następująco:

$$b(i) = \min_{k \neq k'} \{d_{iC_{k'}}\}, \text{ gdzie } d_{iC_k} = \frac{1}{n_k} \sum_{m \in C_k} d_{im}.$$

Na tej podstawie określona zostaje $S(P_r)$ jako prawidłowość wyodrębnienia poszczególnych klas (oznaczonych przez P_r) oraz wskaźnik $S(P)$, opisujący ogólną jakość klasyfikacji dane równaniami (5).

$$S(P_r) = \frac{1}{n_k} \sum_{i \in P_r} S(i), \quad S(P) = \frac{1}{k} \sum_r S(P_r). \quad (5)$$

We wszystkich wymienionych indeksach wykorzystano odległość zastosowaną do wyznaczenia klasyfikacji. W obliczeniach niektórych indeksów wykorzystano procedury z pakietu clusterSim programu R [Walesiak, Dudek 2012; R. Development Core Team 2005].

O lepszej jakości klasyfikacji mówią wyższe wartości indeksu Calińskiego i Harabasz, a także niższe indeksy Daviesa-Bouldina. Silhouette indeks z przedziału od 0,5 do 0,7 świadczy o poważnej strukturze klas, natomiast wartości wyższe niż 0,7 charakteryzują silną strukturę klas [Gatnar, Walesiak (red.) 2004].

3. Entropia

Niektóre sposoby określania miary podobieństwa wywodzą się z teorii informacji. Zakłada się, że takie grupowanie, które daje największy przyrost informacji, jest optymalne, gdyż odpowiada to małemu zróżnicowaniu kategorii w podzbiorach. Entropia jako miara zróżnicowania wydaje się funkcją bardzo uniwersalną, niezależną od charakteru zmiennych.

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

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

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

$$H(x) = -\log(\int f^2(x) dx). \quad (6)$$

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 w oparciu o estymację jądrową [Liang i in. 2011; Jensen i in. 2003]. Niech:

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

W naszych badaniach została wykorzystana funkcja jądrowa Gaussa, określona wzorem (7).

$$W_{\sigma^2}(x, x_i) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{(x-x_i)^T(x-x_i)}{2\sigma^2}\right). \quad (7)$$

Można pokazać, że entropię 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 klastrów: $C_k, k = 1, \dots, K$, w których pojawia się N_k obiektów. Entropię w k -tym klastrze można zdefiniować jako:

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

Wskaźnik postaci:

$$V(C_1, C_2, \dots, C_K) = \sum_{i=1}^K \frac{N_k}{N} H(C_k) \quad (8)$$

mógłby 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

Rényiego, otrzymujemy równanie (8). Ponadto zdefiniujemy entropię pomiędzy grupami jako:

$$H(C_1, C_2, \dots, C_K) = -\log \frac{1}{2 \prod_{k=1}^K N_k} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N W_{2\sigma^2}(x_j, x_i). \quad (9)$$

Jeśli klastry są dobrze wybrane, wskaźnik ten powinien być odpowiednio duży [Jenssen i in. 2003]. Jeśli zatem zastosujemy iloraz:

$$V = \frac{H(C_1, C_2, \dots, C_K)}{V(C_1, C_2, \dots, C_K)}, \quad (10)$$

to otrzymamy wskaźnik jakości grupowania, zgodny z ideą tworzenia wskaźnika CH, ale może lepiej działający dla grup z różnych rozkładów, w szczególności dla grup nietypowych.

4. Badanie empiryczne

Własności wskaźnika przedstawionego wzorem (10) zostały zaprezentowane na tle innych wybranych indeksów dla danych empirycznych. W badaniu empirycznym porównane zostały wyniki klasyfikacji dla czterech zestawów danych. Pod uwagę zostały wzięte dane empiryczne w postaci krótkich szeregów czasowych. W przykładzie pierwszym i drugim zastosowano klasyfikację 24 spółek branży informatycznej notowanych na GPW w Warszawie. W przykładzie pierwszym brano pod uwagę zysk netto tych spółek w latach (na zakończenie roku kalendarzowego) 2004-2011. W drugim spółki te klasyfikowano pod kątem wartości ich przychodów w tym samym okresie. Długość szeregu określona została na $d = 8$. W przykładzie trzecim rozważono szeregi czasowe zawierające wartości procentowe miesięcznej inflacji (zmiany miesiąc/miesiąc poprzedni) z 22 wybranych krajów w okresie od stycznia 2010 do czerwca 2013. W tym przypadku wymiar każdego szeregu wyniósł 42 pomiary. Ostatni zbiór danych zawiera liczbę rejestracji samochodów dokonywaną na terenie 16 wybranych państw w poszczególnych miesiącach roku 2011. Rozważano 12 pomiarów dla każdego szeregu. Uzyskane wartości wskaźników jakości klasyfikacji zaprezentowano w tabelach 1, 2, 3 i 4.

Tabela 1. Wskaźniki grupowania zysku netto w latach 2004-2011 spółek branży informatycznej

Liczba grup	Ward				<i>k</i> -means				spectral			
	CH	DB	S	E	CH	DB	S	E	CH	DB	S	E
2	17,20	1,05	0,41	1,97	17,97	1,11	0,41	1,96	16,99	1,13	0,40	1,73
3	17,76	1,06	0,41	1,83	17,96	1,05	0,42	1,83	13,35	1,28	0,37	1,90
4	16,85	0,99	0,40	1,97	13,05	1,18	0,27	1,62	15,94	0,99	0,39	1,93

Źródło: opracowanie własne.

W tabeli 1 porównane zostały wyniki uzyskane przy zastosowaniu klasyfikacji hierarchicznej metodą Warda, klasyfikacji k -średnich oraz klasyfikacji spektralnej z wybraną miarą odległości euklidesowej. Najlepsze wskaźniki zostały wyłuszczone. Należy zwrócić uwagę na zgodność współczynnika konstruowanego w oparciu o entropię ze wskaźnikiem CH, ale tylko w przypadku metody k -średnich. W przypadku metody Warda i klasyfikacji spektralnej wnioski są rozbieżne. Dla tej ostatniej wskaźnik entropii jest zgodny ze wskaźnikiem DB, a dla metody aglomeracyjnej Warda daje wynik taki sam jak indeks S. Wybierając najlepszą klasyfikację na podstawie wskaźnika CH, otrzymujemy podział na dwie grupy metodą k -średnich. Indeks DB wskazuje na klasyfikację metodą Warda i klasyfikację spektralną z podziałem na cztery grupy. Wskaźnik S jako najlepszą określa metodę k -średnich z liczbą czterech grup. Natomiast według proponowanego indeksu E najlepsze byłaby grupowanie metodą Warda z podziałem na dwie grupy.

Tabela 2 zawiera wyniki uzyskane dla zbioru danych zawierającego wartości przychodów wybranych spółek branży informatycznej. We wszystkich grupowaniach wartości wskaźników wskazały na wybór podziału zbioru spółek na dwa zbiory. Należy zwrócić uwagę, że indeks entropii wskazał na grupowanie analogicznie jak pozostałe wskaźniki jakości klasyfikacji.

Tabela 2. Wskaźniki grupowania przychodów w latach 2004-2011 spółek branży informatycznej

Liczba grup	Ward				k -means				spectral			
	CH	DB	S	E	CH	DB	S	E	CH	DB	S	E
2	21,00	0,66	0,57	2,24	21,00	0,66	0,57	2,24	21,00	0,66	0,57	2,24
3	18,46	1,13	0,31	1,42	18,46	1,13	0,31	1,42	3,55	2,96	0,10	1,36
4	14,95	1,28	0,24	1,47	15,34	1,14	0,26	1,46	13,82	1,03	0,22	1,42

Źródło: opracowanie własne.

Tabela 3. Wskaźniki grupowania inflacji z 22 krajów w okresie od stycznia 2010 do czerwca 2013

Liczba grup	Ward				k -means				spectral			
	CH	DB	S	E	CH	DB	S	E	CH	DB	S	E
2	7,57	1,54	0,24	1,97	7,57	1,54	0,24	1,97	5,91	1,80	0,16	1,66
3	5,63	1,70	0,14	1,72	5,63	1,29	0,24	1,99	3,59	2,25	0,11	1,70
4	5,07	1,52	0,16	1,74	4,69	1,55	0,14	1,60	2,29	2,37	0,05	1,58

Źródło: opracowanie własne.

Przy klasyfikacji wskaźnika inflacji proponowany indeks wykorzystujący miarę entropii dla algorytmu Warda ustala liczbę klas podobnie jak indeks CH. Przy zastosowaniu metody k -średnich wskaźnik E pozwala na określenie, że najlepszą klasyfikację otrzyma się przy podziale na trzy grupy, zgodnie z indeksem DB. Niestety, pozostałe dwa wskaźniki proponują podział zbioru szeregów na dwie grupy. W przypadku grupowania metodą spektralną wskaźnik entropii jako najlepszą kla-

syfikację proponuje podział na trzy grupy, w przeciwieństwie do pozostałych wskaźników, które określają podział 22 państw na dwie grupy jako lepszy.

Tabela 4. Wskaźniki grupowania dla miesięcznej liczby rejestrowanych samochodów w 16 krajach w roku 2011

Liczba grup	Ward				<i>k</i> -means				spectral			
	CH	DB	S	E	CH	DB	S	E	CH	DB	S	E
2	5,76	1,66	0,23	1,51	6,01	1,44	0,23	1,75	3,89	1,88	0,15	1,73
3	5,16	1,37	0,20	1,61	6,03	1,24	0,22	1,79	4,02	1,69	0,14	1,74
4	6,10	1,20	0,23	1,72	4,87	1,34	0,17	1,77	3,29	1,44	0,11	1,66

Źródło: opracowanie własne.

W klasyfikacji Warda wszystkie indeksy wybrały jako najlepszą liczbę czterech podgrup. Klasyfikacja metodą *k*-średnich jest mniej jednoznaczna, ale indeks E daje wynik zgodny z wskaźnikami CH i DB. W klasyfikacji spektralnej indeks oparty na entropii jest zgodny jedynie z indeksem CH. Indeksy CH i DB określają jako najlepszą klasyfikację metodą aglomeracji Warda z podziałem na cztery grupy, natomiast według wskaźnika S najlepszy wynik daje metoda *k*-średnich z podziałem na dwie grupy. Indeks E proponuje natomiast wybrać metodę *k*-średnich, ale z podziałem na trzy podzbiory.

Tabela 5. Wskaźniki grupowania dla wygenerowanych modelowych danych

Liczba grup	Ward				<i>k</i> -means				spectral			
	CH	DB	S	E	CH	DB	S	E	CH	DB	S	E
2	13,29	1,09	0,29	2,89	12,81	1,16	0,28	2,75	11,56	1,71	0,26	2,78
3	19,91	1,02	0,41	2,95	19,97	1,02	0,42	3,10	17,42	1,11	0,39	2,95
4	40,96	0,67	0,55	3,12	40,96	0,67	0,55	3,12	40,96	0,67	0,55	3,12
5	33,2	0,87	0,47	2,05	31,47	1,07	0,42	1,66	33,21	0,87	0,47	2,05

Źródło: opracowanie własne.

Aby zbadać przydatność proponowanego indeksu do określenia jakości klastrowania, wykonano krótkie badanie symulacyjne, które jest wstępem do dalszych badań. Wygenerowano zestawy danych pochodzących z populacji wielowymiarowych rozkładów normalnych o określonej liczbie klas, a następnie wspomnianymi metodami dokonano ich klasyfikacji i wyznaczono omawiane wskaźniki. Zazwyczaj wyniki indeksów dobrze wykrywały strukturę grupową wygenerowanych danych. W tabeli 5 zawarto wyniki dla wygenerowanych 36 wektorów o długości 20 pomiarów. Dane pochodziły z populacji podzielonej na cztery klasy.

5. Podsumowanie

W pracy przedstawiono próbę określenia jakości klasyfikacji za pomocą wskaźnika opartego na entropii Rényiego. Uzyskane wyniki dla wybranych danych empirycznych wskazują na podobieństwo do istniejących wskaźników klasyfikacji. Proponowany indeks E najczęściej dawał wyniki zgodne ze wskaźnikami CH i DB, ale były przypadki, w których zachowywał się inaczej. Krótkie badanie symulacyjne pokazało, że dla wielowymiarowych rozkładów normalnych wskaźnik E, na równi z innymi wskaźnikami, poprawnie wykrywa strukturę grupową danych. Zaproponowany wskaźnik mógłby być wykorzystany jako miara jakości klasyfikacji w przypadku nieznanego rozkładów szeregów czasowych i jest zachętą do dalszych badań.

Literatura

- Baarsch J., Celebi M.C. (2012), *Investigation of Internal Validity Measures for K-means Clustering*, IMECS, Hong Kong.
- Davies D., Bouldin D. (1979), *A Cluster Separation Measure*, IEEE Transactions on Pattern Analysis and Machine Intelligence 1(2), s. 224-227.
- Calinski R.B., Harabasz J. (1974), *A Dendrite Method for Cluster Analysis*, Communications in Statistics – Theory and Methods 3(1), s. 1-27.
- Gatnar E., Walesiak M. (red.) (2004), *Metody statystycznej analizy wielowymiarowej w badaniach marketingowych*, Wydawnictwo AE, Wrocław.
- Halkidi M., Yannis B., Vazirgiannis M. (2001), *On Clustering Validation Techniques*, „Journal of Intelligent Information Systems”, 17, 2/3, s. 107-145.
- Jenssen R., Hild K.E., Erdogmus D., Principe J.C., Eltoft T. (2003), *Clusterin Using Renyi's Entropy*, Proceedings of the International Joint Conference on Neural Networks, Vol. 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, Vol. 45, s. 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*, „International 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, s. 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), s. 53-65.
- Walesiak M., Gatnar E. (2009), *Statystyczna analiza danych z wykorzystaniem programu R*, Wyd. Naukowe PWN, Warszawa.
- Walesiak M., Dudek M. (2012), *Package 'clusterSim' in R project*, <http://keii.ue.wroc.pl/clusterSim/index.html> (30.08.2013).
- Wędrowska E. (2010), *Wykorzystanie entropii Shanona i jej uogólnień do badania rozkładu prawdopodobieństwa zmiennej losowej dyskretnej*, „Przegląd Statystyczny”, LVII, 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, <http://www.R-project.org> (30.08.2013).

VALIDATION OF TIME SERIES CLUSTERING

Summary: The aim of this paper is to present a quality index classification using Renyi's entropy against known quality indicators grouping of multidimensional time series. The starting point was the 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. The proposed index seems to be promising but it needs to be verified for various distributions of time series.

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