

**PRACE NAUKOWE**

Uniwersytetu Ekonomicznego we Wrocławiu

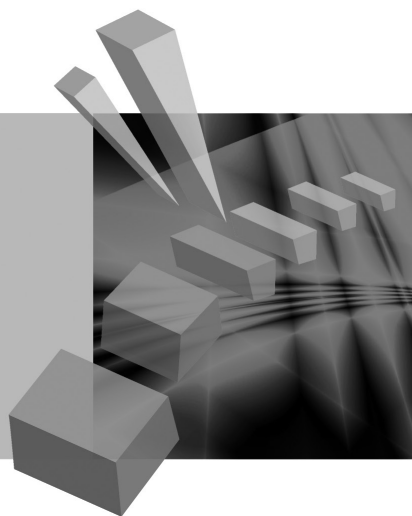
**RESEARCH PAPERS**

of Wrocław University of Economics

**279**

# Taksonomia 21

## Klasyfikacja i analiza danych – teoria i zastosowania



Redaktorzy naukowi

**Krzysztof Jajuga**

**Marek Walesiak**



Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu  
Wrocław 2013

Redaktor Wydawnictwa: Aleksandra Śliwka

Redaktor techniczny: Barbara Łopusiewicz

Korektor: Barbara Cibis

Łamanie: Małgorzata Czupryńska

Projekt okładki: Beata Dębska

Publikacja jest dostępna w Internecie na stronach:

[www.ibuk.pl](http://www.ibuk.pl), [www.ebscohost.com](http://www.ebscohost.com),

The Central and Eastern European Online Library [www.ceeol.com](http://www.ceeol.com),

a także w adnotowanej bibliografii zagadnień ekonomicznych BazEkon

[http://kangur.uek.krakow.pl/bazy\\_ae/bazekon/nowy/index.php](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](http://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 2013

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

**ISSN 1505-9332** (Taksonomia)

Wersja pierwotna: publikacja drukowana

Druk: Drukarnia TOTEM

## Spis treści

<b>Wstęp</b> .....	9
<b>Sabina Denkowska, Kamil Fijorek, Marcin Salamaga, Andrzej Sokolowski:</b> Sejm VI kadencji – maszynka do głosowania .....	11
<b>Barbara Pawelek, Adam Sagan:</b> Zmienne ukryte w modelach ekonomicznych – respecyfikacja modelu Kleina I .....	19
<b>Jan Paradysz:</b> Nowe możliwości badania koniunktury na rynku pracy .....	29
<b>Krzysztof Najman:</b> Samouczące się sieci GNG w grupowaniu dynamicznym zbiorów o wysokim wymiarze .....	41
<b>Kamila Migdał-Najman:</b> Zastosowanie jednowymiarowej sieci SOM do wyboru cech zmiennych w grupowaniu dynamicznym .....	48
<b>Aleksandra Matuszewska-Janica, Dorota Witkowska:</b> Zróżnicowanie płac ze względu na płeć: zastosowanie drzew klasyfikacyjnych .....	58
<b>Iwona Foryś, Ewa Putek-Szeląg:</b> Przestrzenna klasyfikacja gmin ze względu na sprzedaż użytków gruntowych zbywanych przez ANR w województwie zachodniopomorskim .....	67
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk:</b> Klasyfikacja internetowych rachunków bankowych z uwzględnieniem zmiennych symbolicznych.....	77
<b>Marta Jaročka:</b> Wpływ metody doboru cech diagnostycznych na wynik porządkowania liniowego na przykładzie rankingu polskich uczelni .....	85
<b>Anna Zamojska:</b> Badanie zgodności rankingów wyznaczonych według różnych wskaźników efektywności zarządzania portfelem na przykładzie funduszy inwestycyjnych.....	95
<b>Dorota Rozmus:</b> Porównanie dokładności taksonomicznej metody propagacji podobieństwa oraz zagregowanych algorytmów taksonomicznych opartych na idei metody <i>bagging</i> .....	106
<b>Ewa Wędrowska:</b> Wrażliwość miar dywergencji jako mierników niepodobieństwa struktur.....	115
<b>Katarzyna Wójcik, Janusz Tuchowski:</b> Wpływ automatycznego tłumaczenia na wyniki automatycznej identyfikacji charakteru opinii konsumenckich ...	124
<b>Małgorzata Misztal:</b> Ocena wpływu wybranych metod imputacji na wyniki klasyfikacji obiektów w modelach drzew klasyfikacyjnych.....	135
<b>Anna Czapkiewicz, Beata Basiura:</b> Badanie wpływu wyboru współczynnika zależności na grupowanie szeregów czasowych .....	146
<b>Tomasz Szubert:</b> Czynniki różnicujące poziom zadowolenia z życia oraz wartości życiowe osób sprawnych i niepełnosprawnych w świetle badań „Diagnozy społecznej” .....	154

<b>Marcin Szymkowiak:</b> Konstrukcja estymatorów kalibracyjnych wartości globalnej dla różnych funkcji odległości .....	164
<b>Wojciech Roszka:</b> Szacowanie łącznych charakterystyk cech nieobserwowanych łącznie .....	174
<b>Justyna Brzezińska:</b> Metody wizualizacji danych jakościowych w programie <b>R</b> .....	182
<b>Agata Sielska:</b> Regionalne zróżnicowanie potencjału konkurencyjnego polskich gospodarstw rolnych w województwach po akcesji do Unii Europejskiej .....	191
<b>Mariusz Kubus:</b> Liniowy model prawdopodobieństwa z regularyzacją jako metoda doboru zmiennych .....	201
<b>Beata Basiura:</b> Metoda Warda w zastosowaniu klasyfikacji województw Polski z różnymi miarami odległości .....	209
<b>Katarzyna Wardzińska:</b> Wykorzystanie metody obwiedni danych w procesie klasyfikacji przedsiębiorstw .....	217
<b>Katarzyna Dębowska:</b> Modelowanie upadłości przedsiębiorstw oparte na próbach niezbilansowanych .....	226
<b>Danuta Tarka:</b> Wpływ metody doboru cech diagnostycznych na wyniki klasyfikacji obiektów na przykładzie danych dotyczących ochrony środowiska ..	235
<b>Artur Czech:</b> Zastosowanie wybranych metod doboru zmiennych diagnostycznych w badaniach konsumpcji w ujęciu pośrednim .....	246
<b>Beata Bal-Domańska:</b> Ocena relacji zachodzących między inteligentnym rozwojem a spójnością ekonomiczną w wymiarze regionalnym z wykorzystaniem modeli panelowych .....	255
<b>Mariola Chrzanowska:</b> <i>Ordinary kriging</i> i <i>inverse distance weighting</i> jako metody szacowania cen nieruchomości na przykładzie warszawskiego rynku .....	264
<b>Adam Depta:</b> Zastosowanie analizy wariancji w badaniu jakości życia na podstawie kwestionariusza SF-36v2 .....	272
<b>Maciej Beręsewicz, Tomasz Klimanek:</b> Wykorzystanie estymacji pośredniej uwzględniającej korelację przestrzenną w badaniach cen mieszkań .....	281
<b>Karolina Paradysz:</b> Benchmarkowa analiza estymacji dla małych obszarów na lokalnych rynkach pracy .....	291
<b>Anna Gryko-Nikitin:</b> Dobór parametrów w równoległych algorytmach genetycznych dla problemu plecakowego .....	301
<b>Tomasz Ząbkowski, Piotr Jałowiecki:</b> Zastosowanie reguł asocjacyjnych do analizy danych ankietowych w wybranych obszarach logistyki przedsiębiorstw przetwórstwa rolno-spożywczego .....	311
<b>Agnieszka Przedborska, Małgorzata Misztal:</b> Zastosowanie metod statystyki wielowymiarowej do oceny wydolności stawów kolanowych u pacjentów z chorobą zwyrodnieniową leczonych operacyjnie .....	321
<b>Dorota Perło:</b> Rozwój zrównoważony w wymiarze gospodarczym, społecznym i środowiskowym – analiza przestrzenna .....	331

<b>Ewa Putek-Szeląg, Urszula Gieraltowska, Analiza i diagnoza wielkości produkcji energii odnawialnej w Polsce na tle krajów Unii Europejskiej..</b>	342
--	-----

## Summaries

<b>Sabina Denkowska, Kamil Fijorek, Marcin Salamaga, Andrzej Sokolowski: VIth-term Sejm – a voting machine .....</b>	18
<b>Barbara Pawelek, Adam Sagan: Latent variables in econometric models – respecification of Klein I model .....</b>	28
<b>Jan Paradysz: New possibilities for studying the situation on the labour market .....</b>	40
<b>Krzysztof Najman: Self-learning neural network of GNG type in the dynamic clustering of high-dimensional data.....</b>	47
<b>Kamila Migdał-Najman: Applying the one-dimensional SOM network to select variables in dynamic clustering .....</b>	57
<b>Aleksandra Matuszewska-Janica, Dorota Witkowska: Gender wage gap: application of classification trees.....</b>	66
<b>Iwona Foryś, Ewa Putek-Szeląg: Spatial classification of communes by usable land traded by the APA in the Zachodniopomorskie voivodeship...</b>	76
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk: Classification of Internet banking accounts including symbolic variables .....</b>	84
<b>Marta Jarocka: The impact of the method of the selection of diagnostic variables on the result of linear ordering on the example of ranking of universities in Poland.....</b>	94
<b>Anna Zamojska: Empirical analysis of the consistency of mutual fund ranking for different portfolio performance measures.....</b>	105
<b>Dorota Rozmus: Comparison of accuracy of affinity propagation clustering and cluster ensembles based on bagging idea.....</b>	114
<b>Ewa Wędrowska: Sensitivity of divergence measures as structure dissimilarity measurements .....</b>	123
<b>Katarzyna Wójcik, Janusz Tuchowski: Machine translation impact on the results of the sentiment analysis .....</b>	134
<b>Małgorzata Misztal: Assessment of the influence of selected imputation methods on the results of object classification using classification trees ...</b>	145
<b>Anna Czapkiewicz, Beata Basiura: Simulation study of the selection of coefficient depending on the clustering time series.....</b>	153
<b>Tomasz Szubert: Factors differentiating the level of satisfaction with life and the life's values of people with and without disabilities in the light of the "Social Diagnosis" survey .....</b>	162
<b>Marcin Szymkowiak: Construction of calibration estimators of totals for different distance measures .....</b>	173

<b>Wojciech Roszka:</b> Joint characteristics' estimation of variables not jointly observed.....	181
<b>Justyna Brzezińska:</b> Visualizing categorical data in $\mathbf{R}$ .....	190
<b>Agata Sielska:</b> Regional diversity of competitiveness potential of Polish farms after the accession to the European Union .....	200
<b>Mariusz Kubus:</b> Regularized linear probability model as a filter .....	208
<b>Beata Basiura:</b> The Ward method in the application for classification of Polish voivodeships with different distances.....	216
<b>Katarzyna Wardzińska:</b> Application of Data Envelopment Analysis in company classification process.....	225
<b>Katarzyna Dębowska:</b> Modeling corporate bankruptcy based on unbalanced samples .....	234
<b>Danuta Tarka:</b> Influence of the features selection method on the results of objects classification using environmental data.....	245
<b>Artur Czech:</b> Application of chosen methods for the selection of diagnostic variables in indirect consumption research.....	254
<b>Beata Bal-Domańska:</b> Assessment of relations occurring between smart growth and economic cohesion in regional dimension using panel models	263
<b>Mariola Chrzanowska:</b> Ordinary kriging and inverse distance weighting as methods of estimating prices based on Warsaw real estate market .....	271
<b>Adam Depta:</b> Application of analysis of variance in the study of the quality of life based on questionnaire SF-36v2 .....	280
<b>Maciej Beręsewicz, Tomasz Klimanek:</b> Using indirect estimation with spatial autocorrelation in dwelling price surveys.....	290
<b>Karolina Paradysz:</b> Benchmark analysis of small area estimation on local labor markets .....	300
<b>Anna Gryko-Nikitin:</b> Selection of various parameters of parallel evolutionary algorithm for knapsack problems .....	310
<b>Tomasz Ząbkowski, Piotr Jałowiecki:</b> Application of association rules for the survey of data analysis in the selected areas of logistics in food processing companies .....	320
<b>Agnieszka Przedborska, Małgorzata Misztal:</b> Using multivariate statistical methods to assess the capacity of the knee joint among the patients treated surgically for osteoarthritis .....	330
<b>Dorota Perło:</b> Sustainable development in the economic, social and environmental dimensions – spatial analysis.....	341
<b>Ewa Putek-Szeląg, Urszula Gieraltowska:</b> Analysis and diagnosis of the volume of renewable energy production in Poland compared to EU countries .....	352

**Marcin Szymkowiak**

Uniwersytet Ekonomiczny w Poznaniu

---

## KONSTRUKCJA ESTYMATORÓW KALIBRACYJNYCH WARTOŚCI GLOBALNEJ DLA RÓŻNYCH FUNKCJI ODLEGŁOŚCI

---

**Streszczenie:** W badaniach statystycznych jedną z metod umożliwiających redukcję obciążenia i zwiększenie precyzji szacunku na skutek występowania braków informacji jest kalibracja, której podstawy teoretyczne zostały zaproponowane przez Devilla i Särndala [1992]. W klasycznym ujęciu wyznaczanie wag kalibracyjnych oparte jest na odpowiednio dobranej funkcji odległości, która minimalizuje odległość między wyjściowymi wagami wynikającymi ze schematu losowania próby a tzw. wagami kalibracyjnymi. W artykule przedstawione zostały różne funkcje odległości, które można wykorzystać na etapie konstrukcji wag kalibracyjnych. W części empirycznej, z wykorzystaniem programu R i funkcji `calib` dostępnej w pakiecie `sampling`, pokazane zostało, w jaki sposób wyznaczać wagi kalibracyjne w badaniach z brakami odpowiedzi dla różnych funkcji odległości.

**Słowa kluczowe:** kalibracja, wagi kalibracyjne, estymatory kalibracyjne, braki odpowiedzi, funkcja odległości.

### 1. Wstęp

W badaniach statystycznych prowadzonych przez urzędy statystyczne braki odpowiedzi stanowią jeden z istotnych problemów, który wpływa na jakość zebranych danych, a w konsekwencji na cały proces estymacji. Jedną z metod umożliwiających redukcję obciążenia i zwiększenie precyzji szacunku na skutek występowania braków informacji jest kalibracja, której podstawy teoretyczne zostały zaproponowane przez Devilla i Särndala [1992].

Zgodnie z definicją zaproponowaną przez Lundströma i Särndala [Särndal, Lundström 2005; Särndal 2007] kalibracja to metoda polegająca na korygowaniu wag wyjściowych wynikających ze schematu losowania próby, tak aby spełnione były odpowiednie równania kalibracyjne w odniesieniu do zmiennych pomocniczych. W wyniku jej zastosowania najczęściej udaje się zredukować obciążenie i wariancję wykorzystywanych w uogólnianiu wyników estymatorów.

W klasycznym ujęciu wyznaczanie wag kalibracyjnych oparte jest na odpowiednio dobranej funkcji odległości, która minimalizuje odległość między wyjściowymi

wagami wynikającymi ze schematu losowania próby a tzw. wagami kalibracyjnymi. Wykorzystuje się przy tym funkcję odległości opartą na tzw. metryce chi-kwadrat. W artykule przedstawione zostaną inne funkcje odległości, które można wykorzystać na etapie konstrukcji wag kalibracyjnych. W części empirycznej, z wykorzystaniem pakietu R, przedstawiona zostanie metoda ich wyznaczania wraz z ich empiryczną oceną.

## 2. Teoretyczne podstawy kalibracji

Niech dana będzie  $N$ -elementowa populacja  $U = \{1, \dots, N\}$ . Z populacji tej losujemy zgodnie z określonym schematem losowania  $n$ -elementową próbę  $s \subseteq U$ . Niech  $\pi_i$  oznacza prawdopodobieństwo inkluzji  $i$ -tej jednostki do próby, tzn.

$\pi_i = P(i \in s)$  dla  $i = 1, \dots, N$ , a  $d_i = \frac{1}{\pi_i}$  będzie wagą odpowiadającą jednostce  $i$ .

Założmy, że celem badania jest oszacowanie wartości globalnej pewnej zmiennej  $y$ , określonej wzorem:

$$Y = \sum_{i=1}^N y_i, \quad (1)$$

gdzie  $y_i$  oznacza wartość zmiennej  $y$  dla  $i$ -tej jednostki badania,  $i = 1, \dots, N$ .

Klasycznym estymatorem wartości globalnej (1) jest znany z metody reprezentacyjnej estymator Horwitza-Thompsona, który wyraża się wzorem:

$$\hat{Y}_{HT} = \sum_s d_i y_i = \sum_{i=1}^n d_i y_i. \quad (2)$$

Jeżeli nie są znane wszystkie wartości zmiennej  $y$  dla jednostek wylosowanych do próby (na przykład na skutek braków odpowiedzi), estymator Horwitza-Thompsona charakteryzuje się znacznym obciążeniem i dużą wariancją. Wynika to na ogół z faktu, że braki odpowiedzi nie mają charakteru czysto losowego, a powstałe błędy wynikają z różnic pomiędzy respondentami i nierespondentami. Zmniejsza się ponadto efektywna liczebność próby, co w konsekwencji powoduje, że sumowanie we wzorze (2) nie odbywa się po zbiorze wszystkich jednostek, które miały wziąć udział w badaniu, a tylko po zbiorze respondentów  $r \subseteq s$ . Zakładać przy tym będziemy, że jest to zbiór  $m$ -elementowy, przy czym  $m \leq n$ . W efekcie ważona suma (2) jest najczęściej niedoszacowana w stosunku do prawdziwej wartości (1). W związku z tym wagi  $d_i$  powinny zostać odpowiednio skorygowane (skalibrowane), aby zniwelować obciążenie wynikające z braków odpowiedzi.



Oznaczmy przez  $w_i$  poszukiwaną wagę (tzw. wagę kalibracyjną) odnoszącą się do  $i$ -tego respondenta,  $i = 1, \dots, m$ . Naszym celem jest poszukanie wag  $w_i$  w taki sposób, aby były możliwie jak najbliższe co do wartości wyjściowym wagom  $d_i$  i aby niwelowały obciążenie będące konsekwencją występowania braków odpowiedzi. Konstrukcja wag kalibracyjnych uzależniona jest od wyboru odpowiedniej funkcji odległości. W literaturze przedmiotu na potrzeby wyznaczania wag kalibracyjnych przyjmuje się najczęściej tzw. funkcję odległości chi-kwadrat

$$D(\mathbf{w}, \mathbf{d}) = \frac{1}{2} \sum_{i=1}^m \frac{(w_i - d_i)^2}{d_i}, \quad (3)$$

gdyż umożliwia to uzyskanie wektora wag kalibracyjnych w jawnej postaci. Można pokazać [Szymkowiak 2007], że dla tej funkcji odległości wektor wag kalibracyjnych wyraża się wzorem:

$$w_i = d_i + d_i (\mathbf{X} - \hat{\mathbf{X}})^T \left( \sum_{i=1}^m d_i \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \mathbf{x}_i, \quad (4)$$

przy czym  $\mathbf{X}$  to wektor utworzony z wartości globalnej każdej zmiennej pomocniczej  $x_1, \dots, x_k$  tj.

$$\mathbf{X} = \left( \sum_{i=1}^N x_{i1}, \sum_{i=1}^N x_{i2}, \dots, \sum_{i=1}^N x_{ik} \right)^T, \quad (5)$$

$\hat{\mathbf{X}}$  jest wektorem złożonym z oszacowanych wartości globalnych zmiennych pomocniczych

$$\hat{\mathbf{X}} = \left( \sum_{i=1}^m d_i x_{i1}, \sum_{i=1}^m d_i x_{i2}, \dots, \sum_{i=1}^m d_i x_{ik} \right)^T, \quad (6)$$

a

$$\mathbf{x}_i = (x_{i1}, \dots, x_{ik})^T \quad (7)$$

jest wektorem złożonym z wartości wszystkich  $k$  zmiennych pomocniczych dla  $i$ -tego respondenta,  $i = 1, \dots, m$ . Estymator kalibracyjny wartości globalnej (1) wyraża się wówczas wzorem:

$$\hat{Y}_{\mathbf{X}} = \sum_{i=1}^m w_i y_i. \quad (8)$$

Jedną z zalet, wyznaczonych na podstawie funkcji odległości (3) wag kalibracyjnych (4), jest możliwość ich uzyskania wprost ze wzoru. Wagi te jednak w pewnych sytuacjach wykazują pewne niepożądane właściwości, tj. mogą być ujemne dla niektórych respondentów bądź odstające, tzn. znacznie różnić się od wag wyjściowych wynikających ze schematu losowania próby. Na wagi<sup>1</sup> nakłada się zatem czasami warunki ograniczające, tj. wymusza się, aby znajdowały się one w pewnym z góry określonym przedziale, co zapobiega występowaniu wag ujemnych i odstających. Takie podejście nie zapewnia zazwyczaj uzyskania wag w jawnej postaci i zachodzi potrzeba korzystania z metod iteracyjnych w poszukiwaniu wektora wag kalibracyjnych. Nie stanowi to oczywiście istotnej bariery, jednak w przypadku źle wyspecyfikowanych warunków ograniczających nie ma gwarancji, że algorytm poszukiwania wag kalibracyjnych, na które nałożono pewne restrykcje, osiągnie zbieżność. W związku z tym proponuje się wiele różnych funkcji odległości, uwzględniając przy tym (bądź nie) warunki ograniczające na iloraz wag  $w_i \setminus d_i$ . Prowadzi to, w zależności od przyjętej funkcji odległości, do wag kalibracyjnych w jawnej postaci bądź wymaga zastosowania metod numerycznych w poszukiwaniu przybliżonych rozwiązań. Poniżej przedstawiono omawiane najczęściej w literaturze przedmiotu funkcje odległości, które są wykorzystywane w badaniach statystycznych z brakami odpowiedzi w procesie poszukiwania wag kalibracyjnych [Deville, Särndal 1992; Pumputis 2005; Plikusas, Pumputis 2004].

W zależności od przyjętej funkcji odległości uzyskuje się różne postacie wag kalibracyjnych. W procesie ich poszukiwania wykorzystuje się przy tym metodę czynników nieoznaczonych Lagrange'a.

Kończącą postacią wag kalibracyjnych dla wybranych funkcji odległości (funkcja 1, 3, 6 i 7 w tab. 1) można znaleźć w pracy Pumputisa [2005] oraz Plikusasa i Pumputisa [2004]. W pracach tych można znaleźć również formuły na wariancję estymatorów kalibracyjnych w zależności od zastosowanej funkcji odległości. Szczegółowo omówione wyprowadzenie wag kalibracyjnych dla funkcji odległości chi-kwadrat można znaleźć również w pracy Szymkowiaka [2007]. Nie dla wszystkich jednak przedstawionych w tab. 1 funkcji odległości można wyznaczyć analityczną postać wag kalibracyjnych (na przykład dla funkcji 2). Należy ponadto podkreślić, że dla każdej z prezentowanych w tab. 1 funkcji odległości 1-9 można nałożyć ograniczenia na iloraz wag kalibracyjnych i wag wynikających ze schematu losowania próby. Przykładowo możemy przyjąć, że  $0,8 \leq w_i \setminus d_i \leq 1,2$ . Oznacza to, że wagi kalibracyjne  $w_i$  nie mogą się różnić od wag wynikających ze schematu losowania próby  $d_i$  o więcej niż 20%. Brak analitycznej postaci wag kalibracyjnych dla niektórych funkcji odległości oraz możliwość nałożenia na wagi pewnych ograniczeń powoduje, że należy stosować algorytmy iteracyjne w procesie

---

<sup>1</sup> Dokładniej na iloraz wag kalibracyjnych  $w_i$  i wag wynikających ze schematu losowania próby  $d_i$ , tj.  $w_i \setminus d_i$ .

poszukiwania optymalnych wag finalnych. Wykorzystuje się w tym celu najczęściej metodę Newtona-Raphsona bądź inne iteracyjne algorytmy rozwiązywania równań.

**Tabela 1.** Funkcje odległości w zagadnieniu wyznaczania wag kalibracyjnych

Nazwa	Postać funkcji odległości $D(\mathbf{w}, \mathbf{d})$
Funkcja odległości chi-kwadrat	$D_1(\mathbf{w}, \mathbf{d}) = \frac{1}{2} \sum_{i=1}^m \frac{(w_i - d_i)^2}{d_i}$
Zmodyfikowana funkcja odległości oparta na entropii	$D_2(\mathbf{w}, \mathbf{d}) = \sum_{i=1}^m \left( w_i \log \frac{w_i}{d_i} - w_i + d_i \right)$
Funkcja odległości Hellingera	$D_3(\mathbf{w}, \mathbf{d}) = 2 \sum_{i=1}^m \left( \sqrt{w_i} - \sqrt{d_i} \right)^2$
Funkcja odległości oparta na entropii	$D_4(\mathbf{w}, \mathbf{d}) = \sum_{i=1}^m \left( -d_i \log \frac{w_i}{d_i} + w_i - d_i \right)$
Zmodyfikowana funkcja odległości chi-kwadrat	$D_5(\mathbf{w}, \mathbf{d}) = \frac{1}{2} \sum_{i=1}^m \frac{(w_i - d_i)^2}{w_i}$
Funkcja odległości Plikusasa 1	$D_6(\mathbf{w}, \mathbf{d}) = \sum_{i=1}^m \left( \frac{w_i}{d_i} - 1 \right)^2$
Funkcja odległości Plikusasa 2	$D_7(\mathbf{w}, \mathbf{d}) = \sum_{i=1}^m \left( \frac{\sqrt{w_i}}{\sqrt{d_i}} - 1 \right)^2$
Raking	$D_8(\mathbf{w}, \mathbf{d}) = \sum_{i=1}^m \left( \frac{w_i}{d_i} \ln \frac{w_i}{d_i} - \frac{w_i}{d_i} + 1 \right)$
Logitowa funkcja odległości	$D_9(\mathbf{w}, \mathbf{d}) = A^{-1} \sum_{i=1}^m \left[ \left( \frac{w_i}{d_i} - L \right) \ln \frac{\frac{w_i}{d_i} - L}{1 - L} + \left( U - \frac{w_i}{d_i} \right) \ln \frac{U - \frac{w_i}{d_i}}{U - 1} \right]$
	przy czym $A = \frac{U - L}{(U - 1)(1 - L)}$ a L i U oznaczają dolne i górne ograniczenie na iloraz wag $w_i \setminus d_i$ .

Źródło: opracowanie własne.

Istnieje wiele programów, w których zaimplementowany został algorytm wyznaczania wag kalibracyjnych w zależności od przyjętej postaci funkcji odległości i warunków ograniczających. Większość z nich została napisana w języku 4GL w systemie SAS (CLAN 97, CALMAR, GES). Wyjątek stanowi G-CALIB, który został oprogramowany w programie SPSS. Duże możliwości w tym zakresie oferuje

również program R, w którym kalibracja została szeroko opisana w dwóch pakietach: *survey* i *sampling*.

### 3. Kalibracja w programie R

W programie R podejście kalibracyjne zostało zaimplementowane w dwóch pakietach: *survey* autorstwa Thomasa Lumleya oraz *sampling*, którego autorami są Yves Tillé i Alina Matei.

**Tabela 2.** Funkcje pakietu *sampling* związane z podejściem kalibracyjnym

calib(Xs, d, total, method, bounds, description, maxiter) – funkcja programu R wyznaczająca wagi kalibracyjne oraz umożliwiającą ich ocenę	
Argumenty funkcji	
Xs	macierz ze zmiennymi kalibracyjnymi (pomocniczymi)
d	wektor wag wejściowych (podlegających kalibracji)
total	wektor wartości globalnych zmiennych pomocniczych
method	metoda kalibracji (dostępne są cztery funkcje odległości: linear, raking, truncated, logit)
bounds	ograniczenia na wagi kalibracyjne
description	jeśli description=TRUE, tworzony jest raport podsumowujący wagi wejściowe i kalibracyjne, domyślnie description=FALSE
maxiter	liczba iteracji w algorytmie poszukiwania wag kalibracyjnych (niewymagana dla pierwszej metody)
checkcalibration(Xs, d, total, g, EPS) – funkcja programu R umożliwiająca sprawdzenie, czy algorytm wyznaczania wag kalibracyjnych osiągnął zbieżność. Argumenty podobne jak w funkcji calib, przy czym dodatkowy argument EPS określa dopuszczalne różnice między wartościami globalnymi zmiennych pomocniczych a oszacowanymi wartościami globalnymi tych zmiennych na podstawie wag kalibracyjnych.	

Źródło: opracowanie własne.

Ze względu na podobne zastosowanie funkcji wyznaczających wagi kalibracyjne w obydwu pakietach oraz argumenty wejściowe przedstawiona zostanie jedynie funkcja *calib* z pakietu *sampling* i funkcje powiązane umożliwiające dodatkową analizę wag kalibracyjnych (por. tab. 2).

Jednym z bardzo ważnych argumentów funkcji *calib* jest *method*, który może przyjmować jedną z 4 wartości, tj. *linear*, *raking*, *truncated* oraz *logit* w zależności od przyjętej funkcji odległości. Ustalając *method*=„*linear*”, przyjmujemy funkcję odległości chi-kwadrat, *method*=„*raking*” oznacza przyjęcie funkcji odległości typu *raking*, *method*=„*truncated*” oznacza przyjęcie funkcji odległości chi-kwadrat, jednak dodatkowo zakłada się, że iloraz wag  $w_i \setminus d_i$  powinien znajdować się w góry określonym przedziale (określamy to za pomocą argumentu *bounds*), *method*=„*logit*” oznacza z kolei przyjęcie logitowej funkcji odległości.

#### 4. Przykład zastosowania pakietu *sampling*

Załóżmy, że celem pewnego badania jest oszacowanie łącznego miesięcznego dochodu osób w badanej populacji. Na potrzeby przykładu przyjęto, że dysponujemy informacjami z hipotetycznego badania reprezentacyjnego (por. tab. 3), w którym zebrano m.in. informacje na temat płci (zmienna plec: k – kobieta, m – mężczyzna), klasy miejscowości zamieszkania (zmienna klasa: m – miasto, w – wieś) oraz miesięcznego dochodu (zmienna dochod). Zakładamy przy tym, że łączny miesięczny dochód wszystkich osób wynosi 2 700 000 zł. Przyjmijmy ponadto, że próbę o liczebności  $n = 20$  wylosowano z populacji składającej się z  $N = 1000$  osób zgodnie ze schematem losowania prostego ze zwracaniem. Stąd wagi wejściowe  $d_i$  są równe  $N/n = 1000/20 = 50$ . Załóżmy ponadto, że w badanej populacji jest 500 mężczyzn i 500 kobiet oraz 600 osób z miasta i 400 ze wsi. Ponieważ dla części osób nie posiadamy informacji o ich miesięcznym dochodzie (NA – not available), ważona suma wynosząca 2 285 000 zł wyznaczona po zbiorze wszystkich respondentów, dla których znany jest dochód, zgodnie z formułą 2 jest niedoszacowana w stosunku do prawdziwej wartości. Zgodnie z ideą kalibracji należy skorygować wagi dla respondentów, którzy podali informacje o dochodzie, tak aby odtworzone zostały znane struktury dla zmiennych pomocniczych i zniwelowane zostało obciążenie będące konsekwencją braków odpowiedzi.

Tabela 3. Przykładowy zbiór danych<sup>2</sup>

Lp.	plec	klasa	dochod	$d_i$	x1	x2	x3	w1	w2	w3	w4
1	m	m	2000	50	0	1	1	63,39	63,45	62,50	62,96
2	k	w	2500	50	1	0	0	54,24	54,31	53,12	53,70
4	m	w	4000	50	0	1	0	61,02	60,92	62,50	61,73
5	m	m	1500	50	0	1	1	63,39	63,45	62,50	62,96
6	m	m	3500	50	0	1	1	63,39	63,45	62,50	62,96
7	k	m	3700	50	1	0	1	56,61	56,55	57,50	57,04
8	k	w	5500	50	1	0	0	54,24	54,31	53,12	53,70
9	k	m	2400	50	1	0	1	56,61	56,55	57,50	57,04
10	m	m	2200	50	0	1	1	63,39	63,45	62,50	62,96
11	k	w	2800	50	1	0	0	54,24	54,31	53,12	53,70
12	m	w	3200	50	0	1	0	61,02	60,92	62,50	61,73
13	k	m	1600	50	1	0	1	56,61	56,55	57,50	57,04
15	m	m	1900	50	0	1	1	63,39	63,45	62,50	62,96
16	k	w	2100	50	1	0	0	54,24	54,31	53,12	53,70
17	k	m	1400	50	1	0	1	56,61	56,55	57,50	57,04
18	m	w	2500	50	0	1	0	61,02	60,92	62,50	61,73
20	k	m	2900	50	1	0	1	56,61	56,55	57,50	57,04

Źródło: opracowanie własne.

<sup>2</sup> W zbiorze danych ograniczono się jedynie do podania informacji o respondentach, tj. osobach, dla których znany był ich dochód. Brakuje więc osób o liczbie porządkowej 3, 14 i 19, dla których dochod=NA.

W rozważanym przykładzie w charakterze zmiennych pomocniczych wykorzystano płeć oraz klasę miejscowości. Utworzono przy tym 3 zmienne dychotomiczne ( $x_1 - 1$ , jeżeli osoba jest kobietą, 0 – w przeciwnym wypadku;  $x_2 - 1$ , jeżeli osoba jest mężczyzną, 0 – w przeciwnym przypadku;  $x_3 - 1$ , jeżeli osoba jest z miasta, 0 – w przeciwnym wypadku). Tak utworzone zmienne zagwarantują sumowalność odpowiednich wag kalibracyjnych  $w_i$  do faktycznej liczby kobiet i mężczyzn w populacji, a także do liczby osób zamieszkujących miasto oraz wieś<sup>3</sup>.

Na potrzeby wyznaczania wag kalibracyjnych wykorzystano funkcję *calib* zaimplementowaną w programie R w pakiecie *sampling*. W pierwszej kolejności utworzono odpowiednie zbiory wejściowe (por. argumenty tej funkcji zawarte w tab. 2), a następnie wyznaczono wagi kalibracyjne dla wszystkich 4 funkcji odległości, które obsługiwane są przez funkcję *calib*. Składnia poleceń umożliwiająca uzyskanie wag kalibracyjnych jest następująca:

```
library(sampling)
# wczytanie danych
dane <- read.csv("d:/dane.csv", header = TRUE, sep = ";", dec =
",")
# Utworzenie zmiennych pomocniczych
dane$x1 <- ifelse(dane$plec == "k", 1, 0)
dane$x2 <- ifelse(dane$plec == "m", 1, 0)
dane$x3 <- ifelse(dane$klasa == "m", 1, 0)
# ograniczenie zbioru danych do respondentów tj. osób, dla których
znany jest dochód
dane_wej <- subset(dane, !is.na(dane$dochod))
# Utworzenie macierzy Xs
xs <- cbind(dane_wej$x1, dane_wej$x2, dane_wej$x3)
# Utworzenie wektora wartości globalnych
total <- t(cbind(500, 500, 600))
# Utworzenie wektora wag kalibracyjnych dla 4 funkcji odległości
```

---

<sup>3</sup> Z formalnego punktu widzenia do opisu płci wystarczyłoby wziąć jedną zmienną dychotomiczną przyjmującą na przykład wartość 1 dla kobiet i 0 dla mężczyzn. W podejściu kalibracyjnym zagwarantowałyby to jedynie sumowalność wag kalibracyjnych do liczby kobiet w całej populacji, stąd wagi kalibracyjne dla mężczyzn nie musiałyby się sumować do liczby mężczyzn w populacji. Przyjęcie dwóch zmiennych dychotomicznych dla płci zagwarantuje sumowalność wag dla kobiet i mężczyzn do znanych wartości w populacji, a także zapewni sumowalność wag do łącznej liczby osób w populacji, tj. do 1000. Z tego względu dla zmiennej opisującej klasę miejscowości wystarczy utworzyć już tylko jedną zmienną dychotomiczną. Wówczas sumowalność wag kalibracyjnych do wszystkich osób w populacji i sumowalność wag dla osób z miasta do wszystkich osób zamieszkujących miasto wymusi sumowalność wag dla osób ze wsi do łącznej liczby osób mieszkających na wsi. Należy także podkreślić, że wagi korygowane są tylko dla tych respondentów, dla których znany jest miesięczny dochód (we wzorach na funkcje odległości D sumowanie odbywa się po zbiorze respondentów).

```
dane_wej$w1 <- dane_wej$di * calib(Xs, d = dane_wej$di, total, method = "linear")
dane_wej$w2 <- dane_wej$di * calib(Xs, d = dane_wej$di, total, method = "raking")
dane_wej$w3 <- dane_wej$di * calib(Xs, d = dane_wej$di, total, method = "truncated", bounds = c(0.75, 1.25))
dane_wej$w4 <- dane_wej$di * calib(Xs, d = dane_wej$di, total, method = "logit", bounds = c(0.7, 1.3))
print(dane_wej)
```

W wyniku zastosowanej składni poleceń otrzymano raport końcowy w programie R z wagami kalibracyjnymi dla wszystkich 4 funkcji odległości, który przedstawia tab. 3.

Wszystkie wagi kalibracyjne wyznaczone na podstawie 4 różnych formuł pozwalają odtwarzać znane struktury demograficzne na poziomie całej populacji w odniesieniu do zmiennych płeć i klasa miejscowości zamieszkania. Oszacowane na ich podstawie, z wykorzystaniem czterech różnych funkcji odległości, łączne dochody kształtowały się odpowiednio:  $\hat{Y}_X^{linear} = 2\ 674\ 475$  zł;  $\hat{Y}_X^{raking} = 2\ 674\ 422$  zł;  $\hat{Y}_X^{truncated} = 2\ 674\ 879$  zł;  $\hat{Y}_X^{logit} = 2\ 674\ 475$  zł. Wagi te więc, bez względu na przyjętą funkcję odległości, pozwalają dodatkowo redukować obciążenie będące konsekwencją braków odpowiedzi.

## 5. Podsumowanie

W artykule zaprezentowano metody wyznaczania wag kalibracyjnych dla różnych funkcji odległości, w przypadku gdy w badaniu występują braki odpowiedzi. Przedstawiona technika umożliwiła redukcję obciążenia w sytuacji niepełnych danych i w związku z tym może być użyteczna w działalności wszystkich tych instytucji, które na co dzień zajmują się opracowywaniem wyników na podstawie badań ankietowych, do których jednostki dobierane są zgodnie z określonym schematem losowania. Dotyczyć to będzie przede wszystkim ośrodków badania opinii publicznej i Głównego Urzędu Statystycznego, które wykorzystują w wielu badaniach metodę reprezentacyjną na etapie projektowania próby oraz uogólniania wyników na podstawie danych, w których występują braki odpowiedzi.

## Literatura

- Deville J-C., Särndal C-E. (1992), *Calibration estimators in survey sampling*, „Journal of the American Statistical Association”, vol. 87, pp. 376-382.
- Pumpūtis D. (2005), *Calibrated estimators under different distance measures*, Proceedings of the Workshop on Survey Sampling Theory and Methodology 2005, pp. 137-141.

- Plikusas A., Pumputis D. (2004), *Calibrated estimators of totals under different distance measures*, "Lietuvos Matematikos Rinkiny" 2004, vol. 44, special issue, pp. 572-576.
- Särndal C-E., Lundström S. (2005), *Estimation in Surveys with Nonresponse*, John Wiley & Sons, Ltd.
- Särndal C-E. (2007), *The calibration approach in survey theory and practice*, "Survey Methodology", vol. 33, no. 2, pp. 99-119.
- Szymkowiak M. (2007), *Przyczynek do kalibracji w badaniach statystycznych z brakami odpowiedzi*, [w:] *Kapitał ludzki i wiedza w gospodarce. Wyzwania XXI wieku*, E. Panek (red.), Zeszyty Naukowe nr 96, Wydawnictwo Akademii Ekonomicznej w Poznaniu, Poznań 2007, s. 194-204.

## CONSTRUCTION OF CALIBRATION ESTIMATORS OF TOTALS FOR DIFFERENT DISTANCE MEASURES

**Summary:** Missing data are one of the major types of non-random errors in statistical surveys. One of the methods proposed by Deville and Särndal [1992] which is designed to offset the negative effect of missing data is calibration, which is successfully used in practice by statistical offices of many countries. In its classical form calibration is a method in which calibrated weights are computed by minimizing a distance measure between the initial sampling weights and new weights, which need to satisfy certain calibration constraints. The main goal of this paper is to present the construction of calibration estimators of totals for different types of distance measures. Its empirical part, based on the calib function, which is available in R program in the sampling package, is devoted to the method of finding calibration weights in surveys with nonresponse for different distance measures.

**Keywords:** calibration, calibration weights, calibration estimators, nonresponse, distance measures.