

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 Pocięcha, 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 Pociecha, 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

Wojciech Roszka

Uniwersytet Ekonomiczny w Poznaniu
e-mail: wojciech.roszka@ue.poznan.pl

KONSTRUKCJA SYNTETYCZNYCH ZBIORÓW DANYCH NA POTRZEBY ESTYMACJI DLA MAŁYCH DOMEN

Streszczenie: Rozwijająca się gospodarka wymaga wsparcia informacyjnego w szczegółowych przekrojach. Liczebność próby w badaniach reprezentacyjnych uniemożliwia spełnienie tego postulatu, natomiast badania pełne, ze względu na koszty i czas realizacji, przeprowadzane są rzadko. Metody statystycznej integracji danych umożliwiają łączenie informacji z wielu źródeł. Dzięki ich zastosowaniu możliwe jest połączenie bogatych merytorycznie źródeł pochodzących z badań reprezentacyjnych ze zbiorami o pełnym pokryciu, zapewniającym możliwości estymacji dla małych domen. W artykule przedstawiony zostanie proces łączenia informacji z publikacji spisowych z informacjami z badań reprezentacyjnych w celu estymacji stopy bezrobocia w ujęciu powiatów.

Słowa kluczowe: statystyka małych obszarów, integracja danych, badania reprezentacyjne, spis powszechny.

DOI: 10.15611/pn.2015.384.27

1. Wstęp

Statystyka publiczna (jak również inne instytucje państwowe i prywatne) przeprowadzają badania reprezentacyjne pokrywające bardzo wiele aspektów społeczno-gospodarczych. Jednak ze względu na wysokie koszty przeprowadzenia takich badań, jak również czas ich trwania, liczebność próby jest zbyt mała¹, by możliwe było stosowanie estymatorów bezpośrednich dla małych domen. Badania pełne, takie jak spisy powszechny, umożliwiają tworzenie rzetelnych szacunków dla małych domen. Jednak ze względu na wysiłek organizacyjny, jak również koszty, przeprowadzany jest stosunkowo rzadko – z częstotliwością około dziesięcioletnią.

¹ Największe względem pokrycia badanie reprezentacyjne w Polsce – Badanie Aktywności Ekonomicznej Ludności – pokrywa około 1% populacji generalnej w ujęciu rocznym i jego budżet według PBSSP na 2015 rok wynosi 39 586 850 zł.

Wykorzystanie rejestrów administracyjnych jako źródła zasilania informacyjnego gospodarki jest coraz częściej podnoszonym postulatem [Dygaszewicz 2010]. Jednak ich integracja i harmonizacja z populacjami, definicjami i okresami referencyjnymi stosowanymi w statystyce publicznej jest zadaniem trudnym i wymagającym szeregu rozwiązań nie tylko technicznych, ale również prawnych.

Statystyczna integracja danych jest metodą umożliwiającą łączenie informacji z różnych źródeł w celu łącznej obserwacji wielu cech. Jednocześnie dzięki ich zastosowaniu możliwe jest połączenie bogatych merytorycznie źródeł pochodzących z badań reprezentacyjnych ze zbiorami o pełnym pokryciu, zapewniającym możliwości estymacji dla małych domen.

Celem niniejszego opracowania jest próba konstrukcji syntetycznego repozytorium danych jednostkowych o pełnym pokryciu na podstawie publikacji spisowych i zbiorów danych pochodzących z badań reprezentacyjnych w celu estymacji różnych charakterystyk dla małych domen. W artykule przedstawiona zostanie idea konstrukcji takiego zbioru, etapy harmonizacji i integracji różnych źródeł, a także badanie symulacyjne wraz z oceną jakości uzyskanych oszacowań.

2. Idea konstrukcji syntetycznego zbioru danych

Publikacje spisowe są powszechnie dostępne w formie predefiniowanych tablic zawierających informacje o różnych zjawiskach społeczno-ekonomicznych. Posiadając zagregowaną informację o liczbie jednostek współdzielących warianty określonych zmiennych, replikuje się taką kombinację charakterystyk N razy jako N rekordów w jednostkowym, syntetycznym zbiorze danych (por. tab. 1).

Problemem w takim podejściu jest fakt, że tablice spisowe bardzo często nie zawierają łącznej obserwacji wielu zmiennych. W Banku Danych Lokalnych, jak również w innych formach publikacji wyników publikowana jest łączna obserwacja co najwyżej trzech-czterech zmiennych (np. ludność według grup wieku i poziomu wykształcenia w ujęciu powiatu zamieszkania). Charakterystyki łącznego rozkładu większej liczby zmiennych nie są publikowane.

W celu utworzenia możliwie kompleksowego jednostkowego zbioru danych, zawierającego możliwie wiele charakterystyk, wykorzystuje się zbiory danych pochodzące z badań reprezentacyjnych, zawierające łączną obserwację długiego wektora cech. W celu dostosowania rozkładu analizowanych cech do ograniczeń spisowych stosuje się różne miary przeważania [Anderson 2013]. Następnie dokonuje się replikacji przeważonych i zagregowanych zbiorów danych podobnie jak pokazano w tab. 1.

Utworzenie syntetycznego zbioru danych jednostkowych umożliwia nie tylko możliwość bardziej elastycznego tworzenia zestawień, ale również może służyć jako podstawa dołączania informacji z badań reprezentacyjnych [Rahman 2008].

Tabela 1. Dezagregacja tablic spisowych w jednostkowy zbiór danych

Zagregowana tabela spisowa

Powiat	Płeć	Wiek	Wyksz.	N
3064	K	20-24	wyższe	100
2861	M	30-35	zawodowe	50
...



Syntetyczny jednostkowy zbiór danych

Lp.	Powiat	Płeć	Wiek	Wyksz.
1	3064	K	20-24	wyższe
2	3064	K	20-24	wyższe
...
100	3064	K	20-24	wyższe
101	2861	M	30-35	zawodowe
...
150	2861	M	30-35	zawodowe
...

Źródło: opracowanie własne.

Tabela 2. Imputacja metodą syntetycznej rekonstrukcji

Kroki	1.	2.	...	Ostatni																								
Wiek, płeć, stan cywilny	Wiek: 18 Płeć: M St. cyw.: żonaty	Wiek: 34 Płeć: K St. cyw.: rozwiedziona	...	Wiek: 87 Płeć: M St. cyw.: wdowiec																								
Prawdopodobieństwo rodzaju aktywności ekonomicznej dla danych wieku, płci i stanu cywilnego	<table border="1"> <tr><td>Prawd.</td><td>Przedziały skum. prawd.</td></tr> <tr><td>Prac.: 0,4</td><td>(0,0 - 0,4)</td></tr> <tr><td>Bezr.: 0,3</td><td>(0,4 - 0,7)</td></tr> <tr><td>Bier.: 0,3</td><td>(0,7 - 1)</td></tr> </table>	Prawd.	Przedziały skum. prawd.	Prac.: 0,4	(0,0 - 0,4)	Bezr.: 0,3	(0,4 - 0,7)	Bier.: 0,3	(0,7 - 1)	<table border="1"> <tr><td>Prawd.</td><td>Przedziały skum. prawd.</td></tr> <tr><td>Prac.: 0,6</td><td>(0,0 - 0,6)</td></tr> <tr><td>Bezr.: 0,3</td><td>(0,6 - 0,9)</td></tr> <tr><td>Bier.: 0,1</td><td>(0,9 - 1)</td></tr> </table>	Prawd.	Przedziały skum. prawd.	Prac.: 0,6	(0,0 - 0,6)	Bezr.: 0,3	(0,6 - 0,9)	Bier.: 0,1	(0,9 - 1)	...	<table border="1"> <tr><td>Prawd.</td><td>Przedziały skum. prawd.</td></tr> <tr><td>Prac.: 0,0</td><td>(0,0 - 0,0)</td></tr> <tr><td>Bezr.: 0,0</td><td>(0,0 - 0,0)</td></tr> <tr><td>Bier.: 1</td><td>(0,0 - 1)</td></tr> </table>	Prawd.	Przedziały skum. prawd.	Prac.: 0,0	(0,0 - 0,0)	Bezr.: 0,0	(0,0 - 0,0)	Bier.: 1	(0,0 - 1)
Prawd.	Przedziały skum. prawd.																											
Prac.: 0,4	(0,0 - 0,4)																											
Bezr.: 0,3	(0,4 - 0,7)																											
Bier.: 0,3	(0,7 - 1)																											
Prawd.	Przedziały skum. prawd.																											
Prac.: 0,6	(0,0 - 0,6)																											
Bezr.: 0,3	(0,6 - 0,9)																											
Bier.: 0,1	(0,9 - 1)																											
Prawd.	Przedziały skum. prawd.																											
Prac.: 0,0	(0,0 - 0,0)																											
Bezr.: 0,0	(0,0 - 0,0)																											
Bier.: 1	(0,0 - 1)																											
Liczba losowa	0,281	0,709	...	0,481																								
Imputowana aktywność ekonomiczna	Pracujący	Bezrobotna	...	Bierny																								

Źródło: opracowanie własne na podstawie [Williamson 2013].

Dołączanie informacji z informacji z innych źródeł w literaturze (m.in. [Rahman 2008; Williamson 2013]) przedstawia się jako problem imputacji. Jednym z podejść dołączenia informacji z badania reprezentacyjnego do syntetycznej bazy spisowej jest syntetyczna rekonstrukcja (*synthetic reconstruction*). Polega ono na losowaniu wariantów cech z warunkowych prawdopodobieństw (dla domeny) oszacowanych na podstawie badania reprezentacyjnego (por. tab. 2).

3. Badanie symulacyjne

Badanie symulacyjne przeprowadzono w celu zaprezentowania idei tworzenia syntetycznych, jednostkowych zbiorów danych o pełnym pokryciu, jak również empirycznej weryfikacji hipotezy o możliwości dołączenia informacji z badań reprezentacyjnych w celu estymacji dla małych domen. Badanie przeprowadzono na podstawie rezultatów Narodowego Spisu Powszechnego Ludności i Mieszkań 2011 dla ludności oraz wybranych badań reprezentacyjnych².

Moment referencyjny spisu został ustalony na 31 marca 2011. W celu zachowania zgodności okresów referencyjnych zbiory danych zostały ograniczone do jednostek, które podlegały pomiarowi w I i II kwartale 2011 roku³. Wybrano sześć zmiennych zgodnych co do definicji ze zmiennymi spisowymi: powiat zamieszkania, płeć, wiek, stan cywilny, wykształcenie, klasa miejscowości zamieszkania (miasto-wieś). Przyjęto założenie, że zbiory danych pochodzące z badań reprezentacyjnych są rozłączne⁴. Jednocześnie w celu zwiększenia liczebności próby oraz zmniejszenia prawdopodobieństwa wystąpienia „zerowych” domen⁵ dokonano konkatenacji zbiorów danych (por. tab. 3). Wagi analityczne⁶ zharmonizowano, korzystając ze wzoru:

$$w'_i = \frac{w_i}{\sum_{i=1}^n w_i} N, \quad (1)$$

gdzie: w'_i – zharmonizowana waga analityczna dla i -tej jednostki w zintegrowanym zbiorze,

w_i – oryginalna waga analityczna,

N – liczebność populacji generalnej.

Tabela 3. Konkatenacja zbiorów danych pochodzących z badań reprezentacyjnych

Zbiór danych	n
BAEL	139 694
BBGD	44 754
DS	31 401
Ogółem	215 849

Źródło: opracowanie własne.

² Ze względu na dostępność danych wybrano Badanie Aktywności Ekonomicznej Ludności, Badanie Budżetów Gospodarstw Domowych oraz Diagnozę Społeczną.

³ Pomiar w Diagnozie Społecznej przeprowadzony został w marcu i kwietniu 2011 roku. BAEL i BBGD charakteryzują się pomiarem kwartalnym. Jako populację wybrano osoby w wieku 15 lat i więcej.

⁴ Liczebność badań częściowych stanowi zwykle bardzo niewielki odsetek liczebności całej populacji, więc prawdopodobieństwo wylosowania jednej jednostki do dwóch badań jest zbliżone do zera.

⁵ Domeny, które nie dostały się do próby.

⁶ Wagi wynikające ze schematu losowania w poszczególnych badaniach.

Uzyskano zbiór danych o liczebności 215 849 rzeczywistych jednostek. W celu zapewnienia zgodności rozkładów brzegowych analizowanych cech z ograniczeniami spisowymi w połączonym zbiorze danych dokonano przekształcenia wag analitycznych z pomocą metody iteracyjnego dopasowania proporcjonalnego (*Iterational Proportional Fitting*, IPF)⁷. Liczebności cząstkowe⁸ zostały rozszacowane z wykorzystaniem modelu logliniowego [Peck 2011]:

$$N_{ij} = a_i b_i n_{ij} \quad (2)$$

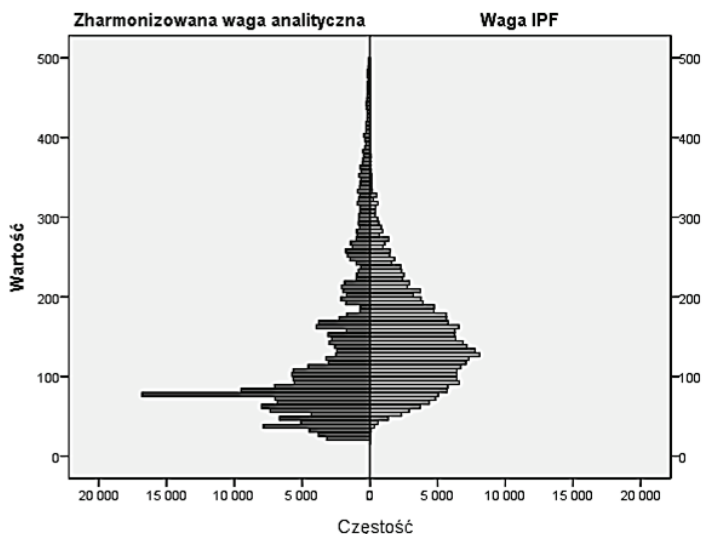
zapisanego jako prawdopodobieństwa:

$$\pi_{ij} = a_i b_i p_{ij}, \quad (3)$$

gdzie π_{ij} i p_{ij} to, odpowiednio, prawdopodobieństwa oszacowane z próby i populacji (spisu):

$$\log\left(\frac{\pi_{ij}}{p_{ij}}\right) = \log(a_i) + \log(b_i) + \epsilon_{ij}. \quad (4)$$

Zakłada się, że liczebności empiryczne są zmiennymi niezależnymi o rozkładzie Poissona. Dopasowanie modelu przeprowadzane jest metodą największej wiarygodności przy użyciu algorytmu Newtona-Raphsona. Rozkład zmodyfikowanych wag analitycznych w_i^{IPF} ilustruje rys. 1.



Rys. 1. Rozkład wag: analitycznych i IPF w zbiorze danych z badań reprezentacyjnych

Źródło: opracowanie własne.

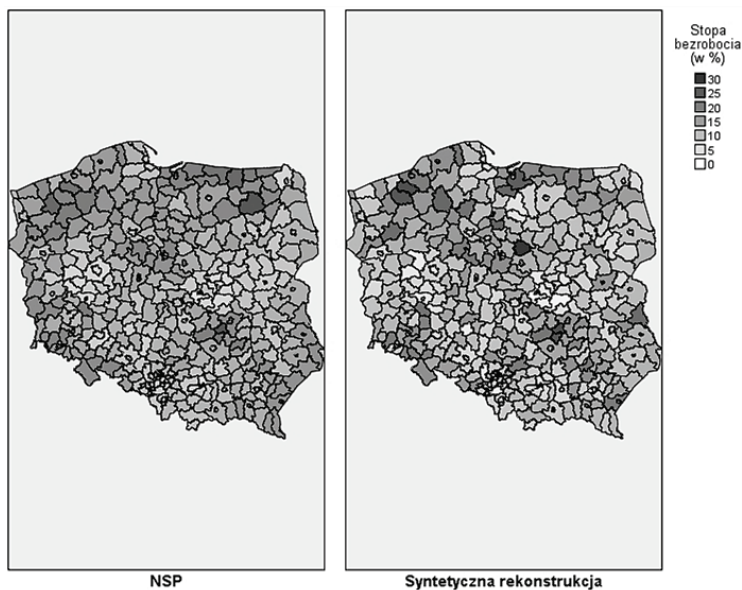
⁷ Inną metodą dostosowania wag do określonych wartości globalnych jest m.in. kalibracja [Józefowski, Szymkowiak 2012].

⁸ Algorytm metody IPF został opisany w [Roszka 2011].

Zbiór danych przeważono według wag IPF. Następnie, zgodnie z metodą przedstawioną w tab. 1, utworzono syntetyczny zbiór danych o liczebności $N = 32\,679\,635$ zawierający łączną obserwację wszystkich analizowanych cech.

W kolejnym kroku do syntetycznego, pełnego zbioru danych dołączono informację o aktywności ekonomicznej ludności z BAEL. Zbiór danych z badań reprezentacyjnych podzielono na 55 170 warstw (domen) na podstawie wariantów analizowanych cech i w każdej warstwie obliczono warunkowe prawdopodobieństwa dla każdego rodzaju aktywności zawodowej. Następnie w syntetycznym zbiorze danych dla każdej jednostki losowano liczbę pseudolosową z rozkładu jednostajnego z przedziału $(0,1)$ i w zależności od przedziału wartości wynikającego ze skumulowanych warunkowych prawdopodobieństw dla danej domeny imputowano określony rodzaj aktywności.

Ocenę jakości uzyskanych szacunków przeprowadzono poprzez porównanie stopy bezrobocia uzyskanej z syntetycznego zbioru z rzeczywistymi wynikami spisowymi. W pierwszej kolejności porównano przestrzenny rozkład prawdziwej i teoretycznej (wynikającej z modelu imputacji) stopy bezrobocia (por. rys. 2).

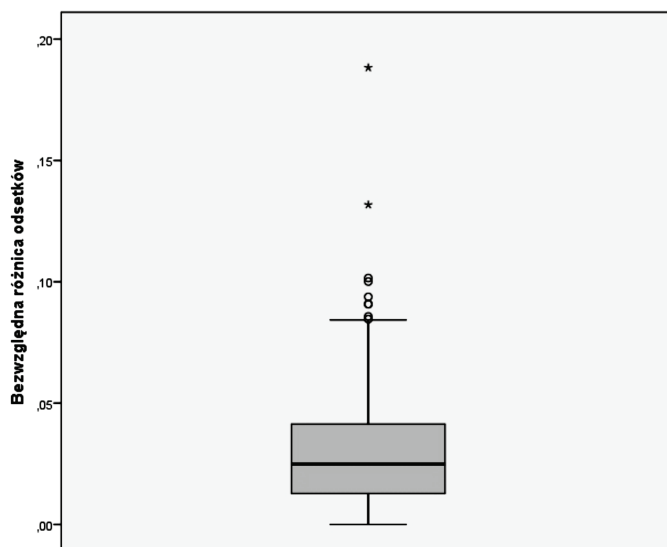


Rys. 2. Rozkład stopy bezrobocia w ujęciu powiatów na podstawie NSP 2011 i syntetycznego zbioru danych

Źródło: opracowanie własne.

Porównując rozkład stopy bezrobocia w ujęciu powiatów, stwierdzono pewne podobieństwo, sugerujące, że imputacja metodą syntetycznej rekonstrukcji w dobry sposób odtwarza rozkład imputowanej cechy.

Następnie obliczono bezwzględne różnice między stopą bezrobocia obliczoną na podstawie NSP a stopą oszacowaną na podstawie syntetycznego zbioru danych. Średnia bezwzględna różnica wynosiła 2,98% przy odchyleniu standardowym równym 2,31%. Najmniejsza różnica wynosiła zaledwie 0,002%, największa zaś aż 18,82% (por. rys. 3).



Rys. 3. Rozkład bezwzględnych różnic NSP i syntetycznego zbioru danych

Źródło: opracowanie własne.

Estymacja stopy bezrobocia dla małych obszarów (powiatów) charakteryzowała się umiarkowaną jakością, jednak dla wielu domen uzyskano szacunki zbliżone do wartości rzeczywistych. Jako główny problem w analizowanych przykładzie można wskazać niewielką liczbę zmiennych w syntetycznym zbiorze danych, które dodatkowo są cechami niemierzalnymi. Utworzenie dobrego modelu imputacji na podstawie stosunkowo niewielu informacji jest utrudnione. Należy jednak podkreślić, że pomimo umiarkowanej jakości szacunków, sposób konstrukcji syntetycznego zbioru danych wraz z dołączaniem do niego informacji z innych źródeł może, przy dostępie większej liczby zmiennych, w tym ilościowych, zapewnić tworzenie dobrych jakościowo szacunków dla małych domen.

4. Zakończenie

W artykule pokazano sposób konstrukcji syntetycznych zbiorów danych o pełnym pokryciu na podstawie publikacji spisowych i zbiorów danych pochodzących z badań reprezentacyjnych. Wykazano, że syntetyczne zbiory danych mogą służyć jako

baza do dołączania informacji z innych źródeł, a dzięki pełnemu pokryciu możliwa jest estymacja dla małych domen.

Jako dalsze kierunki badań można wskazać utworzenie zbiorów z większą liczbą zmiennych, w tym zmiennych mierzalnych, jak również utworzenie syntetycznych zbiorów dla okresów międzypisowych, m.in. z wykorzystaniem danych rejestrowych i informacji pochodzących ze sprawozdawczości bieżącej.

Syntetyczne zbiory danych o pełnym pokryciu mogą być alternatywą dla modelowego podejścia do statystyki małych obszarów [Rao 2003], jak również mogą być źródłem informacji dodatkowej dla podejścia modelowego.

Literatura

- Anderson B., 2013, *Estimating Small Area Income Deprivation: An Iterative proportional Fitting Approach*, [w:] Tanton R., Edwards K.L. 2013, *Spatial Microsimulation: A Reference Guide for Users*, J. Bus. Econ. Stat. 4, s. 87-94.
- Dygaszewicz J., 2010, *Integracja rejestrów publicznych*, Główny Urząd Statystyczny, Warszawa.
- Józefowski T., Szymkowiak M., 2009, *Estymatory kalibracyjne w badaniach statystycznych*, Wiadomości Statystyczne, 2009, | nr 1.
- Peck J., 2011, *Extension Commands and Rim Weighting with IBM SPSS Statistics: Theory and Practice*, IBM Corporation.
- Rahman A., 2008, *A Review of Small Area Estimation Problems and Methodological Developments*, Discussion Paper 66, NATSEM, University of Canberra.
- Rao J.N.K., 2003, *Small Area Estimation*, Wiley and Sons.
- Roszka W., 2011, *Iteracyjne dopasowanie proporcjonalne jako metoda poprawiania wyników w badaniach sondażowych*, [w:] Garczarczyk J., Skikiewicz R., *Metody pomiaru i analizy rynku usług. Dylematy badawcze*, Zeszyty Naukowe UEP, 201.
- Williamson P., 2013, *An Evaluation of Two Synthetic Small-Area Microdata Simulation Methodologies: Synthetic Reconstruction and Combinatorial Optimisation*, [w:] Tanton R., Edwards K.L. 2013, *Spatial Microsimulation: A Reference Guide for Users*, J. Bus. Econ. Stat. 4, s. 8794.

CONSTRUCTION OF SYNTHETIC DATA SETS FOR SMALL AREA ESTIMATION

Summary: A growing economy requires the information support in the specific cross-sections. The sample size in sample surveys makes the fulfillment of this demand impossible, while the full study, due to the cost and time of implementation, is rarely conducted. Statistical methods of data integration allow to combine information from multiple sources. Thanks to their application, it is possible to combine sample surveys and full studies, which makes estimation for small domains possible. In this article the process of combining information from the census publication and information from sample surveys to estimate the rate of unemployment in terms of counties will be presented.

Keywords: statistical data integration, small area estimation, sample survey, census.