

# SPIS TREŚCI

|  |    |
|--|----|
| WSTĘP .....  | 7  |
| 1. PODSTAWOWE ZAGADNIENIA STATYSTYCZNEJ ANALIZY WIELOWYMIAROWEJ .....  | 11 |
| 1.1. Zagadnienia wstępne .....   | 11 |
| 1.2. Typy skal pomiarowych i ich charakterystyka .....   | 15 |
| 1.3. Transformacja normalizacyjna i ujednocianie zmiennych .....   | 17 |
| 1.4. Pomiar podobieństwa obiektów w świetle skal pomiaru i wag zmiennych .....   | 26 |
| 1.5. Strategie postępowania w pomiarze odległości dla danych porządkowych .....  | 35 |
| 2. UOGÓLNIONA MIARA ODLEGŁOŚCI GDM .....   | 40 |
| 2.1. Wprowadzenie .....  | 40 |
| 2.2. Uogólniony współczynnik korelacji .....   | 40 |
| 2.3. Charakterystyka uogólnionej miary odległości .....  | 42 |
| 2.4. Silne i słabe strony uogólnionej miary odległości .....   | 48 |
| 2.5. Postać uogólnionej miary odległości dla zmiennych z różnych skal pomiaru .....  | 50 |
| 2.6. Postać uogólnionej miary odległości dla zróżnicowanych wag zmiennych .....  | 51 |
| 2.7. Wzmacnianie skali pomiaru dla danych porządkowych w statystycznej analizie wielowymiarowej z wykorzystaniem odległości GDM2 ..... | 52 |
| 2.8. Kwadrat odległości euklidesowej a współczynnik korelacji liniowej Pearsona i cosinus kąta między wektorami .....                  | 57 |
| 2.9. GDM a współczynnik korelacji liniowej Pearsona i cosinus kąta między wektorami .....  | 59 |
| 3. OBSZARY ZASTOSOWAŃ UOGÓLNIONEJ MIARY ODLEGŁOŚCI GDM W STATYSTYCZNEJ ANALIZIE WIELOWYMIAROWEJ .....                                  | 64 |
| 3.1. Wyznaczanie macierzy odległości w procesie klasyfikacji obiektów ...  | 64 |
| 3.2. Ocena podobieństwa wyników klasyfikacji zbioru obiektów w czasie  | 78 |
| 3.3. Uogólniona miara odległości GDM jako syntetyczny miernik rozwoju w metodach porządkowania liniowego .....                         | 83 |
| 3.4. Ocena podobieństwa wyników porządkowania liniowego zbioru obiektów w czasie .....   | 88 |

|  |     |
|--|-----|
| 4. WYBÓR GRUP METOD NORMALIZACJI WARTOŚCI ZMIENNYCH W STATYSTYCZNEJ ANALIZIE WIELOWYMIAROWEJ DLA DANYCH METRYCZNYCH .....  | 92  |
| 4.1. Wyniki porządkowania liniowego zbioru obiektów z wykorzystaniem miar syntetycznych a wybór grup metod normalizacji wartości zmiennych .....                               | 92  |
| 4.2. Wybór grup metod normalizacji wartości zmiennych w skalowaniu wielowymiarowym .....   | 103 |
| 4.3. Wybór grup metod normalizacji wartości zmiennych w analizie skupień .....   | 110 |
| 5. UOGÓLNIONA MIARA ODLEGŁOŚCI GDM W ŚWIETLE WYBRANYCH EKSPERYMENTÓW SYMULACYJNYCH .....   | 115 |
| 5.1. Losowe generowanie danych o znanej strukturze klas w pakiecie <code>clusterSim</code> .....   | 115 |
| 5.2. Analiza porównawcza metod klasyfikacji dla danych o znanej strukturze klas .....  | 120 |
| 5.3. Ocena wybranych procedur analizy skupień dla danych porządkowych .....  | 125 |
| 6. WYBRANE ZASTOSOWANIA UOGÓLNIONEJ MIARY ODLEGŁOŚCI GDM Z WYKORZYSTANIEM PROGRAMU <b>R</b> .....  | 131 |
| 6.1. Porządkowanie liniowe zbioru obiektów na podstawie danych porządkowych z rynku nieruchomości .....  | 131 |
| 6.2. Porządkowanie liniowe zbioru obiektów na podstawie danych metrycznych dotyczących warunków zamieszkiwania ludności w miastach ...   | 135 |
| 6.3. Ocena podobieństwa wyników porządkowania liniowego zbioru obiektów w czasie na podstawie danych metrycznych dotyczących warunków zamieszkiwania ludności w miastach ..... | 138 |
| 6.4. Analiza skupień zbioru obiektów opisanych danymi porządkowymi z rynku nieruchomości .....   | 141 |
| 6.5. Analiza skupień zbioru obiektów opisanych danymi metrycznymi dotyczącymi zanieczyszczenia powietrza .....   | 145 |
| LITERATURA .....   | 151 |
| ANEKS .....  | 159 |
| SKOROWIDZ .....  | 165 |

## WSTĘP

Prezentowana książka stanowi podsumowanie rozważań autora zawartych w wielu opracowaniach dotyczących miary odległości, która została w pierwotnej wersji zaproponowana dla zmiennych porządkowych [Walesiak 1993a, s. 44-45], a następnie dla danych metrycznych [Walesiak 2002a] i nominalnych [Walesiak 2003c]. Podstawowe części książki zostały opublikowane m.in. w „Argumenta Oeconomica”, „Przeglądzie Statystycznym”, „Badaniach Operacyjnych i Decyzjach”, Pracach Naukowych Akademii Ekonomicznej we Wrocławiu (Uniwersytetu Ekonomicznego we Wrocławiu) oraz były referowane na konferencjach naukowych, w tym na konferencji Sekcji Klasyfikacji i Analizy Danych PTS (zob. [Walesiak, Bąk, Jajuga 2002; Walesiak 2003b; 2004b; 2011b; 2013; Walesiak, Dudek 2009a; 2010b]), konferencji Międzynarodowej Federacji Towarzystw Klasyfikacyjnych IFCS (zob. [Walesiak, Dziechciarz, Bąk 1998; Walesiak, Dudek 2010a]) oraz Niemieckiego Towarzystwa Klasyfikacyjnego (zob. [Jajuga, Walesiak, Bąk 2003]).

Dotychczas uogólniona miara odległości została zaprezentowana w zwartej postaci w trzech wydaniach książkowych Wydawnictwa Akademii Ekonomicznej (Uniwersytetu Ekonomicznego) we Wrocławiu (zob. [Walesiak 2002b; 2006; 2011d]). Obecna monografia zawiera istotne zmiany i uzupełnienia wynikające w znacznej mierze z nowych badań. Całkowicie nowe są podrozdział 2.7 oraz rozdział 4. Wprowadzono istotne zmiany w podrozdziale 1.3.

Praca składa się z sześciu rozdziałów.

W rozdziale pierwszym przedstawiono podstawowe zagadnienia statystycznej analizy wielowymiarowej. Wyjaśniono w nim takie podstawowe pojęcia, jak obiekt, zmienna, macierz i kostka danych. Scharakteryzowano typy skal pomiarowych oraz zagadnienie transformacji normalizacyjnej i ujednociania zmiennych z punktu widzenia skal pomiarowych. Ponadto zaprezentowano szeroką klasyfikację miar podobieństwa obiektów z uwzględnieniem problematyki ważenia zmiennych oraz skal ich pomiaru. Rozdział kończą rozważania dotyczące strategii postępowania w pomiarze odległości dla danych porządkowych.

W rozdziale drugim przedstawiono szczegółową charakterystykę uogólnionej miary odległości GDM (*Generalised Distance Measure*). W konstrukcji miary odległości GDM wykorzystano ideę uogólnionego współczynnika korelacji, który obejmuje współczynnik korelacji liniowej Pearsona i współczynnik korelacji zmiennych porządkowych tau Kendalla. W związku z tym w części pierwszej tego rozdziału zaprezentowano uogólniony współczynnik korelacji. W dalszej części scharakteryzowano uogólnioną miarę odległości GDM dla jednakowych i zróżnicowanych wag zmiennych. Następnie wskazano silne i słabe strony uogólnionej miary odległości. Rozważania teoretyczne zilustrowano licznymi przykładami poglądowymi. Zapre-

zentowano postać uogólnionej miary odległości GDM uwzględniającą zmienne mierzone na skali metrycznej, porządkowej, nominalnej oraz zmienne z różnych skal pomiaru. Zaproponowano metodę wzmacniania skali pomiaru dla danych porządkowych w statystycznej analizie wielowymiarowej z wykorzystaniem odległości GDM2. Ponadto przedstawiono analizę związków między kwadratem odległości euklidesowej a współczynnikiem korelacji liniowej Pearsona i cosinusem kąta między wektorami oraz między uogólnioną miarą odległości GDM a współczynnikiem korelacji liniowej Pearsona i cosinusem kąta między wektorami.

W rozdziale trzecim zaprezentowano obszary zastosowań uogólnionej miary odległości w statystycznej analizie wielowymiarowej. Podstawowymi obszarami zastosowań tej miary są wyznaczanie macierzy odległości w procesie klasyfikacji zbioru obiektów, w skalowaniu wielowymiarowym oraz zastosowanie miary GDM jako syntetycznego miernika rozwoju w metodach porządkowania liniowego. Ponadto w rozdziale tym zaprezentowano metody oceny podobieństwa wyników klasyfikacji zbioru obiektów oraz oceny podobieństwa wyników porządkowania liniowego zbioru obiektów w czasie.

Rozdział czwarty poświęcono zagadnieniu wyboru metody normalizacji wartości zmiennych w statystycznej analizie wielowymiarowej dla danych metrycznych. W kolejnych trzech podrozdziałach zaprezentowano zagadnienie wyboru metody normalizacji wartości zmiennych w porządkowaniu liniowym zbioru obiektów z wykorzystaniem miar syntetycznych, w skalowaniu wielowymiarowym oraz w analizie skupień.

Rozdział piąty zawiera rezultaty wybranych eksperymentów symulacyjnych pozwalających ocenić zachowanie się uogólnionej miary odległości GDM przy różnych strukturach danych. W pierwszym podrozdziale scharakteryzowano zagadnienie losowego generowania danych o znanej strukturze klas w pakiecie `clusterSim`. W drugim podrozdziale przedstawiono analizę porównawczą metod klasyfikacji dla danych o znanej strukturze klas dla trzech typów danych. W dwóch pierwszych eksperymentach wykorzystano dane metryczne oraz porządkowe o znanej strukturze klas obiektów wygenerowane z wykorzystaniem z funkcji `cluster.Gen` pakietu `clusterSim`. W eksperymencie trzecim zbiory danych utworzono z wykorzystaniem funkcji pakietu `mlbench` (`spirals`, `smiley`, `cassini`) oraz zbiorów własnych (`worms`, `w3`, `skad`). W podrozdziale trzecim, na podstawie porządkowych danych symulacyjnych wygenerowanych z wykorzystaniem z funkcji `cluster.Gen` pakietu `clusterSim`, przeprowadzono ocenę przydatności wybranych procedur analizy skupień obejmujących miarę odległości GDM, dziewięć metod klasyfikacji oraz osiem indeksów służących ustaleniu liczby klas.

W rozdziale szóstym zaprezentowano wybrane zastosowania uogólnionej miary odległości GDM1 i GDM2 w statystycznej analizie wielowymiarowej z wykorzystaniem programu **R**. Znaczna część skryptów wykorzystuje pakiet `clusterSim`. Zastosowania dotyczyły porządkowania liniowego i analizy skupień zbioru obiektów na podstawie danych porządkowych z rynku nieruchomości oraz porządkowania

liniowego na podstawie danych metrycznych dotyczących warunków zamieszkiwania ludności w miastach i analizy skupień obiektów opisanych danymi metrycznymi dotyczącymi zanieczyszczenia powietrza. Ponadto dokonano oceny podobieństwa wyników porządkowania liniowego zbioru obiektów w czasie na podstawie danych metrycznych dotyczących warunków zamieszkiwania ludności w miastach.

Pracę zamyka zestawienie wykorzystywanej literatury, aneks oraz skorowidz rzeczowy.

Wersję instalacyjną programu **R** oraz dodatkowe pakiety (w tym pakiet `clusterSim` autorstwa Marka Walesiaka i Andrzeja Dudka) można pobrać ze strony: <http://www.r-project.org/>. Wszystkie skrypty zawarte w książce przetestowano, używając wersji 3.3.0 programu **R**.

Na stronie internetowej <http://keii.ue.wroc.pl> znajdują się pliki zawierające wykorzystywane dane oraz skrypty realizujące zastosowania zamieszczone w książce.

Książka jest przeznaczona dla pracowników naukowych zajmujących się zastosowaniem metod statystycznej analizy wielowymiarowej w każdej dziedzinie wiedzy, w tym w badaniach ekonomicznych. Ponadto odbiorcami książki mogą być słuchacze wyższych uczelni studiujący zagadnienia statystycznej analizy wielowymiarowej i jej zastosowań.