

POLITECHNIKA ŁÓDZKA

WYDZIAŁ FIZYKI TECHNICZNEJ, INFORMATYKI I MATEMATYKI
STOSOWANEJ

Kierunek: Matematyka

Specjalność: Matematyczne Metody Analizy Danych Biznesowych

WYBRANE ZASTOSOWANIE STATYSTYCZNYCH METOD
PORZĄDKOWANIA DANYCH WIELOWYMIAROWYCH

Kamila Choja
Nr albumu: 204052

Praca licencjacka
napisana w Instytucie Matematyki Politechniki Łódzkiej

Promotor: dr, mgr inż. Piotr Kowalski

ŁÓDŹ, LIPIEC 2018

Spis treści

1	Wstęp	2
2	Preliminaria	3
2.1	Notacja	3
2.2	Słownik użytych pojęć	4
2.3	Podstawowe pojęcia rachunku prawdopodobieństwa oraz statystyki	6
2.4	Wybrane operacje statystyczne dla zmiennych	9
2.5	Podstawowe pojęcia teorii grafów	11
2.6	Wybrane pojęcia z teorii mnogości, topologii i algebry liniowej	13
2.6.1	Relacja porządkująca	13
2.6.2	Przestrzenie metryczne, miary odległości	15
3	Metody porządkowania	16
3.1	Metody porządkowania liniowego	16
3.1.1	Metody diagramowe	20
3.1.2	Metody oparte na zmiennych syntetycznych	21
3.1.3	Metody iteracyjne	24
3.2	Metody porządkowania nieliniowego	24
3.2.1	Metody dendrytowe	25
3.2.2	Metody aglomeracyjne	26
4	Zastosowanie wybranych metod porządkowania danych wielowymiarowych	29
4.1	Opis zbioru	29
4.2	Użyte programy	30
4.3	Implementacje wybranych metod	30
4.3.1	Stymulacja zmiennych	30
4.3.2	Transformacje normalizacyjne	31
4.3.3	Metody porządkowania nieliniowego	32
4.3.4	Metody porządkowania liniowego	36
4.3.5	Metoda sum	38
4.3.6	Metoda rang	39
4.3.7	Metoda Hellwiga	40
4.3.8	Porównanie wyników metod porządkowania liniowego dla całego zbioru	42
4.3.9	Zastosowanie funkcji odpowiedzialnych za porządkowania	43
4.3.10	Porównanie wyników	43
4.3.11	Podsumowanie	46
5	Podsumowanie	47

Rozdział 1

Wstęp

Wielowymiarowa analiza danych jest istotnym pojęciem we współczesnej analizie. Pośród wielu zadań z tej dziedziny, w tej pracy chcemy skupić się na zadaniu porządkowania danych, tj. wskazywaniu uporządkowania obiektów, reprezentowanych przez wielowymiarowe dane. Przedstawione w pracy rozwiązania mogą mieć zastosowanie do wielu problemów napotykanym w codziennej pracy analityków. Do poprawnego zrozumienia zagadnienia statystycznego porządkowania danych należy jednak zgromadzić wiedzę i teorię z wielu obszarów matematyki i statystyki oraz zaprezentować ich zastosowanie w praktycznych przykładach. Takie zgrupowanie wyżej wskazanych zagadnień jest głównym zadaniem niniejszej pracy.

Praca została uporządkowana w 5 rozdziałach. Rozdział 1 stanowi bieżący wstęp. Z kolei rozdział 2 zawiera potrzebne teorie różnych działów matematyki i statystyki, które są wykorzystywane w kolejnych częściach pracy. Omówione są w nim także zagadnienia dotyczące podstaw rachunku prawdopodobieństwa oraz statystyki (sekcja 2.3), niezbędne do rozumienia wielowymiarowych danych jako losowej próby prostej pewnej wielowymiarowej zmiennej losowej oraz celem wprowadzenia jednolitych oznaczeń. Znaczną część tego rozdziału poświęcamy matematycznej teorii porządków - tak liniowych jak i częściowych, albowiem algorytmy prezentowane w dalszych rozdziałach nawiązują do teorii porządków, definiowanej w ramach współczesnej teorii mnogości. Z uwagi na fakt, iż nie wszystkie porządki wydobywane z danych są porządkami liniowymi, omawiamy również podstawy oraz wybrane elementy z teorii grafów. Porządki częściowe mogą być bowiem prezentowane na strukturach grafowych ze znakomitą korzyścią dla przejrzystości. Oprócz tego w rozdziale 2 zawarliśmy też elementy teorii przestrzeni metrycznych, gdyż odległości pomiędzy wektorami danych są istotnym elementem w prawie każdym z omawianych algorytmów. Tutaj zawarte zostały również opisy podstawowych przekształceń na zbiorach danych, wykorzystywanych na etapie wstępnego ich przetwarzania.

Kluczowe teorie dotyczące samego porządkowania danych statystycznych zgromadzone są w rozdziale 3. W tej części pracy omawiane są również najistotniejsze sposoby wydobywania porządków z danych statystycznych, prezentowana jest ich systematyka oraz omawiane są właściwości oraz odmiany. Podstawowy podział na: metody porządkowania liniowego - pozwalające określić ukryte porządki liniowe, oraz metody nieliniowe - służące do wskazywania grafowych reprezentacji odkrytych porządków częściowych, rozdziela dwie główne sekcje tego rozdziału. Z kolei w rozdziale 4 przedstawiamy eksperymenty przeprowadzone celem lepszego zaprezentowania treści rozdziału 3. Dla potrzeb tej pracy został wytworzony zbiór wielowymiarowych danych - reprezentujący pewien podzbiór ogłoszeń znaczącego portalu z ofertami sprzedaży pojazdów. Ponadto w rozdziale tym opisujemy własne implementacje wybranych algorytmów z rozdziału 3 opracowane w języku R i prezentujemy uzyskane porządki. W rozdziale 5 zawarte jest treściwe podsumowanie, zarówno z zakresu opisu statystycznych algorytmów porządkowania danych jak i dla wyników przeprowadzonych eksperymentów.

Rozdział 2

Preliminaria

2.1 Notacja

Poniżej znajduje się lista pojęć powszechnie używanych w pracy wraz z symbolami, które się im przypisuje.

- \mathbb{R} - zbiór liczb rzeczywistych,
- \mathbb{N} - zbiór liczb naturalnych,
- K - oznaczenie dowolnego ciała zbioru,
- $O = \{O_1, O_2, \dots, O_n\}$ - zbiór obiektów przestrzennych, tj. opisywanych przez wiele atrybutów, $n \in \mathbb{N}$,
- $X = [x_{ij}]$ - macierz surowych danych, gdzie x_{ij} -oznacza wartość j -tej zmiennej dla i -tego obiektu, gdzie: $i = 1, \dots, n$, $j = 1, \dots, m$, $n, m \in \mathbb{N}$. W rozdziale drugim, przez X najczęściej będziemy oznaczać dowolny zbiór,
- $N = [n_{ij}]$ - macierz znormalizowanych danych, gdzie n_{ij} oznacza wartość j -tej cechy i -tego obiektu,
- x^S - oznaczenie zmiennej mającej charakter stymulacyjny,
- x^D - oznaczenie zmiennej mającej charakter destymulacyjny,
- x^N - oznaczenie zmiennej mającej charakter nominacyjny,
- $D = [d_{ik}]$ - oznaczenie macierzy odległości, gdzie d_{ik} oznacza odległość między i -tym i k -tym obiektem, gdzie: $i, k = 1, \dots, n$, $n \in \mathbb{N}$,
- s_i - oznaczenie zmiennej syntetycznej i -tego obiektu,
- $P_0 = [n_{0j}]$ - oznaczenie obiektu wzorcowego, gdzie n_{0j} - znormalizowana j -ta współrzędna obiektu wzorcowego,
- Y - w rozdziale drugim używana jest najczęściej do oznaczenia zmiennej losowej,
- Ω - oznaczenie dowolnej przestrzeni zdarzeń elementarnych ω , z rodziną podzbiorów \mathcal{F} ,
- \mathcal{B} - rodzina zbiorów borelowskich,
- \mathfrak{B} - rodzina wszystkich zbiorów otwartych,

- \leq - oznaczenie relacji częściowego porządku, przy dołożeniu warunku spójności oznaczać będzie relację liniowego porządku,
- $\text{med}(\cdot)$ - oznaczenie mediany zbioru,
- $\max(\cdot)$ - oznaczenie maksymalnej wartości zbioru,
- $\min(\cdot)$ - oznaczenie minimalnej wartości zbioru,
- $\overline{(\cdot)}$ - oznaczenie mocy zbioru,
- G - oznaczenie ogólnego grafu prostego dla którego $V(G)$ jest zbiorem wierzchołków grafu, a $E(G)$ zbiorem jego krawędzi,

2.2 Słownik użytych pojęć

W pracy zostały wykorzystane następujące pojęcia:

- Statystyka matematyczna [6, w oparciu o rozdział 1]
Statystyka matematyczna jest nauką zajmującą się opisywaniem i analizą zjawisk przy użyciu metod rachunku prawdopodobieństwa.
- Cecha statystyczna [6, Rozdział 1]
Cecha statystyczna jest to właściwość wspólna dla danego zbioru obserwacji. Jej wartości pozwalają rozróżnić elementy zbioru między sobą. cechy statystyczne można podzielić na te mierzalne, tj. ilościowe (np. długość, ciężar), oraz niemierzalne tj. jakościowe (np. kolor, płeć, zawód, województwo).

W celu prezentacji dużych ilości danych, w analizie danych korzysta się z pojęcia macierzy. Poniżej zostanie przedstawiona formalna definicja macierzy oraz definicja macierzy obserwacji, czyli zbiorze obiektów opisywanych przez zmienne.

Definicja 2.2.1. *Macierz [1, Rozdział 1]*

Niech K będzie ciałem i $m, n \in \mathbb{N}$. Macierz o m -wierszach, n -kolumnach i o wyrazach z K nazywamy każdą funkcję postaci $A: \{1, \dots, m\} \times \{1, \dots, n\} \rightarrow K$

Przykład 1. *Macierz A o m -wierszach i n -kolumnach najczęściej zapisuje się postaci $A = [a_{ij}]_{1 \leq i \leq m, 1 \leq j \leq n}$, tj.*

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

Uwaga 1. *W statystyce koncepcja matematyczna macierzy jest rozszerzana, gdyż niektóre kolumny mają wartości z poza ciała zbioru liczb \mathbb{R} (mogą być np. tekstem).*

Definicja 2.2.2. *Macierz obserwacji [8, Rozdział 2]*

Niech $m > 1$ oraz $n > 1$ będą liczbami naturalnymi. Macierz obserwacji nazywamy macierz rozmiaru $n \times m$ postaci

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}$$

gdzie: x_{ij} - zaobserwowana wartość j -tej cechy dla i -tego obiektu.

Definicja 2.2.3. *Macierz odległości zmiennych [9, Rozdział 1.6]*

Macierzą odległości cech zmiennych nazywamy macierz, której elementami są odległości między parami badanych obiektów:

$$D = [d_{ik}].$$

gdzie:

d_{ik} -odległość między i -tym a k -tym obiektem, dla $i, k = 1, 2, \dots, n$

Uwaga 2. *Zauważmy, że macierz obserwacji nie musi być symetryczna, z kolei w przypadku macierzy odległości jest to wymagane.*

W statystyce posługujemy się pojęciami skal do opisu różnych typów danych, które przyjmowane przez nas mogą podlegać analizie. W związku z powyższym zdefiniujemy następujące rodzaje skal:

- Skala porządkowa [9, Rozdział 1.2]
Zmienna opisana jest na skali porządkowej jeśli jej zbiór wartości jest zbiorem, w którym wprowadzony jest porządek np. porządek liczb. Nie zawsze porządek ten jest ustalony w sposób matematyczny. W przypadku rozpatrywania zmiennych jakościowych, porządek ustalamy na podstawie opinii ekspertów lub ogólnie przyjętych poglądów np. poziom wykształcenia, oceny w systemie szkolnym.
- Skala przedziałowa [9, Rozdział 1.2]
Zmienna jest opisana na skali przedziałowej gdy podobnie jak na skali porządkowej jej zbiór wartości jest zbiorem uporządkowanym z wprowadzoną funkcją odległości. Dodatkowo na skali tej możliwe jest wyznaczenie umownego punktu - zera. Przykład zmiennych przedstawianych na skali przedziałowej to: temperatura, czas.
- Skala ilorazowa [9, Rozdział 1.2]
Zmienna jest opisana na skali ilorazowej, jeśli jej zbiór wartości jest zbiorem postaci $[0, \infty)$ będący podzbiorem zbioru liczb \mathbb{R} , lub też taki zbiór wartości który można utożsamić z podzbiorem liczb \mathbb{R} . Przykłady zmiennych opisywanych na skali ilorazowej: napięcie elektryczne, bezrobocie.

Uwaga 3. *Jeżeli zmiennej opisującej dany obiekt nie da się odnieść do żadnej z powyższych skal, to zmienna ta nazywana jest nominalną.*

Ze względu na to, że zmienne opisujące obiekty mogą mieć różny charakter, poniżej zostały wprowadzone definicje trzech różnych typów zmiennych, którymi posługujemy się w statystyce.

Definicja 2.2.4. *Stymulanta [9, Rozdział 1.5]*

Stymulantami nazywane są te zmienne (cechy), dla których pożądane są wysokie wartości w badanych obiektach, ze względu na rozpatrywane zjawisko.

Definicja 2.2.5. *Destymulanta [9, Rozdział 1.5]*

Destymulantami nazywane są te zmienne, dla których niepożądane są wysokie wartości w badanych obiektach, ze względu na rozpatrywane zjawisko.

Definicja 2.2.6. *Nominanta [9, Rozdział 1.5]*

Nominantami nazywane są te zmienne, które mają określoną najkorzystniejszą wartość. Odchylenia od tej wartości są niepożądane, ze względu na rozpatrywane zjawisko.

2.3 Podstawowe pojęcia rachunku prawdopodobieństwa oraz statystyki

Na potrzeby pracy, zostały wykorzystane pojęcia rachunku prawdopodobieństwa oraz statystyki, konieczne do zrozumienia danych jako próby losowej. W tym celu niezbędne było wprowadzenie definicji prawdopodobieństwa, zmiennej losowej, a także pojęć powiązanych z tymi definicjami tj. ciała zbiorów, σ -ciała zbiorów, przestrzeni zdarzeń elementarnych, zdarzenia losowego.

Definicja 2.3.1. *Ciało zbiorów [10, Rozdział 8.1]*

Rodzinę \mathcal{F} podzbiorów, niepustego zbioru X nazywamy ciałem zbiorów, jeżeli spełnia ona następujące warunki:

1. $\emptyset \in \mathcal{F}$,
2. jeżeli $A \in \mathcal{F}$, to $X \setminus A \in \mathcal{F}$,
3. jeżeli $A \in \mathcal{F}$, to $A \cup B \in \mathcal{F}$.

Definicja 2.3.2. σ -algebra/ciało zbiorów [10, Rozdział 8.1]

Ciało zbiorów \mathcal{F} nazywamy σ -ciałem zbiorów, jeżeli spełnia ona warunek dla dowolnych zbiorów $A_n \in \mathcal{F}, n \in \mathbb{N}$, mamy $\bigcup_{i=1}^{\infty} A_n \in \mathcal{F}$.

Najważniejszym σ -ciałem zbiorów w matematyce są σ -ciała zbiorów Borelowskich, dlatego też wprowadzimy definicję zbiorów borelowskich.

Definicja 2.3.3. *Zbiory borelowskie [4, w opraciu o rozdział 2]*

Zbiorami borelowskimi względem danej przestrzeni X , nazywamy zbiory należące do σ -ciała X generowanego przez rodzinę $\mathfrak{B}(X)$ - wszystkich zbiorów otwartych w X . Rodzinę wszystkich zbiorów borelowskich względem X , oznaczamy $\mathcal{B}(X)$.

σ -ciała zbiorów często mogą być po prostu utożsamiane ze zbiorami, które można zmierzyć, w związku z tym wprowadzimy definicję miary zbioru.

Definicja 2.3.4. *Miara zbioru [4, Rozdział 2.10]*

Funkcję μ określoną na ciele \mathcal{F} podzbiorów zbioru Ω nazywamy miarą, jeśli spełnia następujące warunki:

1. $\forall A \in \mathcal{F} \quad \mu(A) \in [0, \infty]$,
2. $\mu(\emptyset) = 0$,
3. jeśli A_1, A_2, \dots jest ciągiem rozłącznych zbiorów \mathcal{F} -mierzalnych takich, że $\bigcup_{k=1}^{\infty} A_k \in \mathcal{F}$, to

$$\mu\left(\bigcup_{k=1}^{\infty} A_k\right) = \sum_{k=1}^{\infty} \mu(A_k).$$

Definicja 2.3.5. *Przestrzeń mierzalna [4, Rozdział 2.10]*

Przestrzeni mierzalną nazywamy parę (X, \mathcal{F}) , gdzie \mathcal{F} jest σ -ciałem podzbiorów zbioru X

Definicja 2.3.6. *Funkcja mierzalna [10, w oparciu o rozdział 8.2]*

Niech X będzie niepustym zbiorem, \mathcal{F} σ -ciałem na X i $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, \infty\}$. Funkcję $f : X \rightarrow \overline{\mathbb{R}}$ nazywamy mierzalną, jeżeli zbiór $\{x \in X : f(x) > a\}$ jest mierzalny przy dowolnym $a \in \mathbb{R}$.

Mając powyższe definicje, wprowadźmy możemy wprowadzić poszukiwane definicje rachunku prawdopodobieństwa.

Definicja 2.3.7. *Przestrzeń zdarzeń elementarnych [6, w oparciu o rozdział 1.1]*

Zbiór wszystkich możliwych wyników doświadczenia losowego nazywamy przestrzenią zdarzeń elementarnych i oznaczamy przez Ω . Elementy zbioru Ω nazywamy zdarzeniami elementarnymi i oznaczamy ω .

Definicja 2.3.8. *Zdarzenie losowe [6, w oparciu o rozdział 1.1]*

Zdarzeniem losowym (zdarzeniem) nazywamy każdy podzbiór A zbioru Ω , taki że $A \in \mathcal{F}$, gdzie \mathcal{F} jest rodziną podzbiorów Ω spełniającą następujące warunki:

1. $\Omega \in \mathcal{F}$,
2. Jeśli $A \in \mathcal{F}$, to $A' \in \mathcal{F}$, gdzie $A' = \Omega \setminus A$ jest zdarzeniem przeciwnym do zdarzenia A ,
3. Jeśli $A_i \in \mathcal{F}$, $i = 1, 2, \dots$, to $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$.

Rodzinę \mathcal{F} spełniającą warunki 1 - 3 nazywamy σ -ciałem podzbiorów zbioru Ω

Definicja 2.3.9. *Prawdopodobieństwo [6, w oparciu o rozdział 1.1]*

Prawdopodobieństwem nazywamy dowolną funkcję P o wartościach rzeczywistych, określoną na σ -ciele zdarzeń $\mathcal{F} \subset 2^{\Omega}$, spełniającą warunki:

1. $\forall A \in \mathcal{F} \quad P(A) \geq 0$,
2. $P(\Omega) = 1$,
3. Jeśli $A_i \in \mathcal{F}$, $i = 1, 2, \dots$ oraz $A_i \cap A_j = \emptyset$ dla $i \neq j$, to

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Definicja 2.3.10. *Przestrzeń probabilistyczna [6, w oparciu o rozdział 1.2]*

Przestrzenią probabilistyczną nazywamy uporządkowaną trójkę (Ω, \mathcal{F}, P) , gdzie Ω jest zbiorem zdarzeń elementarnych, \mathcal{F} jest σ -ciałem podzbiorów Ω , zaś P jest prawdopodobieństwem określonym na \mathcal{F} .

Dla tak podanej definicji prawdopodobieństwa, definiujemy:

Definicja 2.3.11. *Zmienna losowa [6, Rozdział 2.1]*

Niech (Ω, \mathcal{F}, P) będzie dowolną przestrzenią probabilistyczną. Dowolną funkcję $Y: \Omega \rightarrow \mathbb{R}$ nazywamy zmienną losową jednowymiarową, jeśli dla dowolnej liczby rzeczywistej y zbiór zdarzeń elementarnych ω , dla których spełniona jest nierówność $Y(\omega) < y$ jest zdarzeniem, czyli

$$\{\omega : Y(\omega) < y\} \in \mathcal{F} \quad \forall y \in \mathbb{R}.$$

Definicja 2.3.12. *Wektor losowy [5, Rozdział 5.1]*

Wektorem losowym nazywamy odwzorowanie $Y: \Omega \rightarrow \mathbb{R}^n$, spełniające następujący warunek: dla każdego układu liczb $t_1, t_2, \dots, t_n \in \mathbb{R}$ zbiór $Y^{-1}((-\infty, t_1] \times \dots \times (-\infty, t_n])$ należy do \mathcal{F} .

Definicja 2.3.13. *Rozkład prawdopodobieństwa zmiennej losowej [5, Rozdział 5.1]*

Rozkładem prawdopodobieństwa zmiennej losowej Y o wartościach w \mathbb{R} nazywamy funkcję μ_Y określoną na $\mathcal{B}(\mathbb{R})$ zależnością:

$$\mu_Y(B) = P_Y(B) = P(Y^{-1}(B)), \quad B \in \mathcal{B}(\mathbb{R}).$$

Definicja 2.3.14. Rozkład dyskretny [5, Rozdział 5.1]

Mówimy, że zmienna losowa jednowymiarowa Y ma rozkład dyskretny, jeśli istnieje przeliczalny zbiór $S \subset \mathbb{R}$, taki że $\mu_Y(S) = 1$.

Definicja 2.3.15. Gęstość i rozkład ciągły [5, Rozdział 5.1]

Jeśli μ jest rozkładem prawdopodobieństwa na \mathbb{R} i istnieje całkowalna funkcja $f : \mathbb{R} \rightarrow \mathbb{R}$ taka, że:

$$\mu(A) = \int_A f(y)dy, A \in \mathcal{B}(\mathbb{R})$$

to funkcję f nazywamy gęstością rozkładu μ . Rozkład, który ma gęstość, nazywamy rozkładem ciągłym.

Definicja 2.3.16. Wartość oczekiwana [6, Rozdział 2.6]

Niech X będzie zmienną losową typu dyskretnego lub ciągłego. Wartością oczekiwaną zmiennej losowej Y nazywamy:

$$E(Y) = \begin{cases} \sum_{i=1}^n y_i p_i, & \text{jeśli zmienna ma rozkład dyskretny i przyjmuje dokładnie } n \text{ wartości,} \\ \sum_{i=1}^{\infty} y_i p_i, & \text{jeśli zmienna przyjmuje nieskończenie, ale przeliczalnie wiele wartości,} \\ \int_{-\infty}^{\infty} y f(y) dy, & \text{jeśli zmienna ma rozkład ciągły.} \end{cases}$$

Definicja 2.3.17. Wariancja [5, Rozdział 5.6]

Niech Y będzie zmienną losową. Jeśli $E(Y - EY)^2 < \infty$, to liczbę tę nazywamy wariancją zmiennej losowej Y o wartościach rzeczywistych i oznaczamy:

$$\text{Var } Y = E(Y - EY)^2.$$

Definicja 2.3.18. Odchylenie standardowe [5, Rozdział 5.6]

Niech Y będzie zmienną losową. Odchyleniem standardowym zmiennej losowej Y nazywamy pierwiastek z wariancji:

$$\sigma_Y = \sqrt{\text{Var } Y}.$$

Definicja 2.3.19. Próba losowa n -elementowa [2, Rozdział 2]

Niech (Ω, \mathcal{F}, P) będzie przestrzenią probabilistyczną, zaś $Y_i : \Omega^n \rightarrow \mathbb{R}$ będzie zmienną losową. Próbką losową n -elementową nazywamy n -elementowy ciąg niezależnych zmiennych losowych o jednakowym prawdopodobieństwie, tzn. ciąg postaci $(\omega_1, \omega_2, \dots, \omega_n) \in \Omega^n$, taki że $Y_i(\omega_1, \omega_2, \dots, \omega_n) = Y(\omega_i)$.

Uwaga 4. Powyższa definicja próby losowej jest wyrażona w języku rachunku prawdopodobieństwa. W języku statystyki próbę losową rozumiemy jako wartość pochodząca z realizacji takiego doświadczenia.

Definicja 2.3.20. Rozkład normalny (Gaussa) [5, Rozdział 5.10]

Jeśli zmienna losowa Y ma gęstość postaci

$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu_Y)^2}{2\sigma^2}}$$

dla $y \in \mathbb{R}$ i pewnych $\mu_Y \in \mathbb{R}$ i $\sigma^2 > 0$. To mówimy, że zmienna losowa ma rozkład normalny z parametrami μ i σ^2 , co zapisujemy $\mathcal{N}(\mu, \sigma^2)$.

W przypadku, gdy $\mu = 0$ i $\sigma^2 = 1$, to rozkład ten nazywamy standardowym rozkładem normalnym i oznaczamy $\mathcal{N}(0, 1)$, a gęstość jest postaci

$$f(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}}.$$

2.4 Wybrane operacje statystyczne dla zmiennych

Normalizacja zmiennych

Ważnym krokiem przed rozpoczęciem pracy na zbiorze danych jest ujednolicenie ich charakteru, tj. przekształcenie zmiennych (mierzonych na skali przedziałowej lub ilorazowej) opisujących obiekty w zbiorze, w celu pozbycia się dysproporcji między nimi czy też dominacji jednych zmiennych nad drugimi. W tym celu stosuje się transformację normalizacyjną. Wyróżniamy trzy podstawowe typy przekształceń normalizacyjnych:

- standaryzacja,
- unitaryzacja,
- przekształcenie ilorazowe,
- rangowanie zmiennych.

W dalszej części pracy j -ta zmienna znormalizowana, i -tego obiektu jest oznaczona jako n_{ij} . Dodatkowo przyjmujemy oznaczenia: \bar{x} - średnia z próby, $\sigma(x)$ - odchylenie standardowe z próby.

1. W wyniku standaryzacji zmienne uzyskują odchylenie standardowe równe 1 i wartości oczekiwanej równą 0. W tym celu dla każdej zmiennej, będącej cechą obiektu oblicza się odchylenie standardowe na podstawie wartości tej zmiennej dla wszystkich obiektów, a także wartość oczekiwaną na tej samej zasadzie. W kolejnym kroku dla każdego obiektu liczymy jego znormalizowaną wartość tj. od wartości zmiennej odejmujemy średnią wartość danej cechy, a otrzymaną różnicę dzielimy przez odchylenie standardowe dla tej cechy. Powyższy opis można zapisać w postaci:

$$n_{ij} = \frac{x_{ij} - \bar{x}_j}{\sigma(x_j)}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

2. W pracy korzystamy również z unitaryzacji, stosowanej w celu uzyskania zmiennych o ujednoliconym zakresie zmienności, najczęściej jest to przedział $[0, 1]$. W tym celu od wartości zmiennej w obiekcie, odejmowana jest minimalna wartość występująca dla tej cechy, a następnie różnica ta dzielona jest przez różnicę między maksymalną a minimalną wartością zmiennej, która została poddana normalizacji. Znormalizowana zmienna jest postaci:

$$n_{ij} = \frac{x_{ij} - \min_i(x_{ij})}{\max_i(x_{ij}) - \min_i(x_{ij})}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

3. Kolejna metoda normalizacyjna, wykorzystana w pracy to przekształcenie ilorazowe. Stosuje się je aby odnieść wartości zmiennej do ustalonej wartości - może to być wartość oczekiwana danej zmiennej na tle analizowanych obiektów, wartość minimalna lub maksymalna tej cechy. W pracy za tę wartość przyjęliśmy średnią wartość zmiennej. Każda wartość zmiennej dla danego obiektu jest dzielona przez wartość oczekiwaną tej zmiennej, a postać znormalizowanej zmiennej to:

$$n_{ij} = \frac{x_{ij}}{\bar{x}_j}, \quad \bar{x}_j \neq 0 \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

4. Przy zastosowaniu metody rang wykorzystuje się normalizację rangową. Przekształcenie to, najczęściej stosowane jest, gdy zmienne opisujące obiekty są wyrażone na skali porządkowej. W pierwszym kroku wartości zmiennych opisujących obiekty, zostają uporządkowane ze względu na ich wartości po procesie normalizacji. W kolejnym kroku wartościom

zmiennej przyporządkowywane są rangi - czyli wartości liczbowe, będące najczęściej numerami miejsc zajmowanych przez obiekty w uporządkowanym zbiorze. Postać zmiennej znormalizowanej rangowo:

$$n_{ij} = r, \quad \text{dla} \quad x_{hj} = x_{ij}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

gdzie:

r -ranga nadana i -temu obiektowi znajdującemu się na r -tym miejscu w uporządkowanym zbiorze, ze względu na wartość j -tej zmiennej.

Stymulacja zmiennych

W celu ujednolicenia charakteru zmiennych należy poddać je pewnym przekształceniom, polegającym na zamianie destymulant i nominant na stymulanty. Tego typu transformacje nazywamy stymulacją. Można wyróżnić dwie najczęściej stosowane metody, tj. przekształcenie ilorazowe oraz przekształcenie różnicowe. W zależności od skali, na której mierzone są zmienne, należy stosować odpowiednie przekształcenie stymulacyjne.

Przekształcenie ilorazowe można stosować tylko dla zmiennych mierzonych na skali ilorazowej. Poniżej zaprezentujemy jego postać dla zmiennych o charakterze destymulant oraz nominant.

- Niech x^D oznacza dane o charakterze destymulant, wtedy dane x^S wyznaczone wg. poniższego wzoru, mają charakter stymulant:

$$x_{ij}^S = [x_{ij}^D]^{-1}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m,$$

- Analogicznie do powyższego, niech x^N oznaczają dane o charakterze nominant, wtedy dane x^S wyznaczone wg. poniższego wzoru, mają charakter stymulant:

$$x_{ij}^S = \frac{\min\{x_j^N, x_{ij}^N\}}{\max\{x_j^N, x_{ij}^N\}}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

gdzie:

x_j^N - oznacza pożądaną wartość j -tej zmiennej,

x_{ij}^N - oznacza wartość j -tej zmiennej w i -tym obiekcie.

Dla zmiennych mierzonych na skali ilorazowej czy też przedziałowej stosuje się przekształcenie różnicowe. Postać tego przekształcenia dla zmiennych o charakterze destymulant oraz nominant jest następująca.

- Ponownie niech x^D oznacza dane o charakterze destymulant, wtedy dane x^S wyznaczone wg. poniższego wzoru, mają charakter stymulant:

$$x_{ij}^S = \max_i \{x_{ij}^D\} - x_{ij}^D, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m,$$

- Niech teraz x^N oznaczają dane o charakterze nominant, wtedy dane x^S wyznaczone wg. poniższego wzoru, mają charakter stymulant:

$$x_{ij}^S = -|x_{ij}^N - x_j^N|, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

2.5 Podstawowe pojęcia teorii grafów

W pracy zostaną opisane zarówno metody porządkowania liniowego jak i nieliniowego. W tym celu należy wprowadzić definicje związane z teorią grafów, niezbędne przy opisywaniu metod porządkowania nieliniowego.

W celu wprowadzenia kluczowych definicji, należy wcześniej podać podstawowe pojęcia dotyczące grafów. Zaczniemy od wprowadzenia definicji pary uporządkowanej oraz nieuporządkowanej, gdyż pojęcia te zostały wykorzystane w definicji grafu.

Definicja 2.5.1. Para uporządkowana [7, w oparciu o rozdział 3]

Niech X będzie dowolnym, niepustym zbiorem oraz niech dane będą dwa elementy $a, b \in X$. Parę uporządkowaną nazywamy parę postaci $\langle a, b \rangle$, gdzie element a jest poprzednikiem, zaś element b jest następnikiem: $\langle a, b \rangle = \{\{a\}, \{a, b\}\}$.

Definicja 2.5.2. Para nieuporządkowana [7, w oparciu o rozdział 3]

Niech X będzie dowolnym, niepustym zbiorem oraz niech dane będą dwa elementy $a, b \in X$. Parę nieuporządkowaną nazywamy zbiór postaci $\{a, b\}$, tj. zawierający elementy a i b i nie zawierający żadnego innego elementu. W przypadku, gdy $a = b$, to para nieuporządkowana $\{a, b\}$, składa się dokładnie z jednego elementu.

W analogiczny sposób jak parę dwóch punktów można wprowadzić parę dwóch zbiorów. Istotne jest aby podkreślić różnicę pomiędzy parą dwóch wierzchołków, które tworzą krawędź, a parą zbiorów definiującą graf.

Definicja 2.5.3. Graf [12, w oparciu o rozdział 2]

Grafem nazywamy parę $G = (V, E) = (V(G), E(G))$, gdzie V jest niepustym, skończonym zbiorem wierzchołków grafu G , zaś E jest skończonym podzbiorem zbioru nieuporządkowanych par elementów zbioru V .

Definicja 2.5.4. Pętle [12, Rozdział 2]

Niech $G = (V(G), E(G))$ będzie grafem oraz niech $a \in V(G)$. Pętlami w grafie nazywamy krawędzie reprezentowane przez a , tj. łączące wierzchołek z samym sobą. Innymi słowy jest to para nieuporządkowana składająca się z jednego elementu.

Definicja 2.5.5. Wierzchołki sąsiednie [12, Rozdział 2]

Mówimy, że dwa wierzchołki v i w w grafu G są sąsiednie, jeśli istnieje krawędź vw która je łączy.

$$v \text{ ————— } w$$

Analogicznie definiuje się krawędzie sąsiednie.

Definicja 2.5.6. Krawędzie sąsiednie [12, Rozdział 2]

Dwie krawędzie e i f grafu G są sąsiednie, jeśli mają wspólny wierzchołek, tj

$$\exists d \in V(G) \quad d \in e \wedge d \in f$$

$$\text{---} \overset{e}{\quad} d \text{---} \overset{f}{\quad} \text{---}$$

Aby dowiedzieć się o połączeniu dwóch wierzchołków w grafie, wprowadzimy pojęcie trasy.

Definicja 2.5.7. *Trasa/marszruta [12, Rozdział 3]*

Trasę (lub marszrutę) w danym grafie G nazywamy skończony ciąg krawędzi postaci $v_0v_1, v_1v_2, \dots, v_{m-1}v_m$, zapisywany również w postaci $v_0 \rightarrow v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_m$, w którym każde dwie kolejne krawędzie są albo sąsiednie, albo identyczne. Taka trasa wyznacza ciąg wierzchołków v_0, v_1, \dots, v_m . Wierzchołek v_0 nazywamy wierzchołkiem początkowym, a wierzchołek v_m wierzchołkiem końcowym trasy, mówimy też wtedy, o trasie o d wierzchołka v_0 do wierzchołka v_m . Liczbę krawędzi na trasie nazywamy długością trasy.

Definicja 2.5.8. *Ścieżka [12, Rozdział 3]*

Trasę, w której wszystkie krawędzie są różne, nazywamy ścieżką.

Definicja 2.5.9. *Droga [12, Rozdział 3]*

Ścieżkę, w której wierzchołki o znaczone kolejno: v_0, v_1, \dots, v_m są różne (z wyjątkiem, być może, równości $v_0 = v_m$, nazywamy drogą.

Definicja 2.5.10. *Droga zamknięta/ścieżka zamknięta [12, Rozdział 3]*

Droga jest zamknięta gdy rozpoczyna się i kończy w tym samym punkcie, tj. według przyjętej notacji $v_0 = v_m$.

Definicja 2.5.11. *Cykl [12, Rozdział 3]*

Cyklem nazywamy drogę zamkniętą.

Definicja 2.5.12. *Graf spójny [12, Rozdział 3]*

Graf jest spójny wtedy i tylko wtedy, gdy każda para wierzchołków jest połączona drogą.

Definicja 2.5.13. *Drzewo [12, Rozdział 4]*

Drzewem nazywamy graf spójny, nie zawierający cykli.

Definicja 2.5.14. *Graf skierowany(digraf albo graf zorientowany) [12, Rozdział 7]*

Graf skierowany lub digraf D , składa się z niepustego zbioru skończonego $V(D)$ elementów nazywanych wierzchołkami i skończonej rodziny $E(D)$ par uporządkowanych elementów zbioru $V(D)$, nazywanych łukami. Zbiór $V(D)$ nazywamy zbiorem wierzchołków, a rodzinę $E(D)$ rodziną łuków digrafu D (krawędzi grafu skierowanego). Łuk (v, w) zwykle zapisujemy jako vw . Graf skierowany oznaczamy zwykle w postaci pary uporządkowanej $G = \langle V, E \rangle$.

Uwaga 5. Każdy graf jednoznacznie wyznacza pewną relację dwuargumentową (binarną) w skończonym zbiorze V . Można również powiedzieć odwrotnie, że każda relacja dwuargumentowa (binarna) r w skończonym zbiorze V , wyznacza jednoznacznie graf zorientowany, którego węzłami są elementy skończonego zbioru V , z kolei krawędziami są uporządkowane pary $\langle v, v' \rangle$, należące do r .

Uwaga 6. Niech dany będzie digraf D składający się z niepustego zbioru skończonego wierzchołków $V(D)$ i skończonej rodziny krawędzi $E(D)$. W momencie gdy

$$\forall a, b \in V(D) \quad \langle a, b \rangle \in E(D) \quad \Rightarrow \quad \langle b, a \rangle \in E(D)$$

to taki graf skierowany, może być utożsamiony z grafem niezorientowanym.



2.6 Wybrane pojęcia z teorii mnogości, topologii i algebry liniowej

2.6.1 Relacja porządkująca

W niniejszej pracy skupiamy się na zagadnieniu porządkowania danych wielowymiarowych. Konieczne jest zatem przywołanie odpowiednich sformułowań dotyczących matematycznej definicji porządku. Najbardziej podstawowym pojęciem jest relacja porządku, która zostanie zdefiniowana poniżej. W sekcji tej zostaną również podane pojęcia równoliczności zbiorów, mocy zbioru, które są potrzebne przy definiowaniu własności porządkowania liniowego zbioru obiektów. Dodatkowo przytoczona zostanie definicja zbioru skończonego, ze względu na zastosowanie metod porządkowania na zbiorze skończonym.

Definicja 2.6.1. *Relacja [7, Rozdział 3]*

Niech dane będą zbiory X i Y . Relacją (dwuargumentową) między elementami zbiorów X i Y nazywamy dowolny podzbiór $\rho \subset X \times Y$. Jeśli $X = Y$ to mówimy, że ρ jest relacją na zbiorze X .

Definicja 2.6.2. *Relacja porządkująca (częściowego porządku) [3, Rozdział 2]*

Niech dana relacja ρ , którą oznaczać będziemy przez \leq , będzie określona dla elementów ustalonego zbioru X . Mówimy, że relacja \leq jest relacją częściowego porządku, jeśli spełnione są warunki:

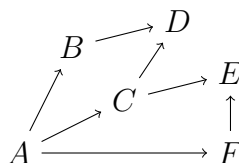
1. $x \leq x$ dla każdego x (zwrotność),
2. jeśli $x \leq y$ i $y \leq x$, to $x = y$ (słaba antysymetryczność),
3. jeśli $x \leq y$ i $y \leq z$, to $x \leq z$ (przechodność).

Przykład 2. *Częściowego porządku na zbiorze*

Wykorzystanie częściowego porządku, obrazuje diagram Hassego, będący grafem skierowanym, którego wierzchołki zostały poddane relacji porządkowania i reprezentują elementy skończonego zbioru X . Aby go skonstruować, należy postępować według poniższych kroków:

- Punkty obrazujące elementy zbioru X , umieszcza się na płaszczyźnie.
- Punkt $x \in X$ łączony jest odcinkiem z punktem $y \in X$ (gdzie $y \neq x$), jeśli x jest następnikiem y , czyli gdy $y \leq x$ oraz nie istnieje taki punkt $z \in X$, że $y \leq z \leq x$ oraz $x \neq z \neq y$.

Zobrazowanie relacji częściowego porządku, dla punktów $A, B, C, D, E, F \in X$.



Definicja 2.6.3. *Relacja liniowo porządkująca (liniowy porządek) [3, Rozdział 2]*

Niech dany będzie niepusty zbiór X . Relację \leq porządkującą zbiór X , nazywamy relacją liniowo porządkującą lub porządkiem liniowym, gdy dla dowolnych $x, y \in X$ spełnia ona następujący warunek spójności tzn. $x \leq y$ lub $y \leq x$. Parę (X, \leq) nazywamy zbiorem liniowo uporządkowanym lub łańcuchem.

Definicja 2.6.4. *Elementy wyróżnione [3, Rozdział 2]*

Niech X będzie zbiorem częściowo uporządkowanym przez relację \leq oraz niech $a \in X$. Mówimy, że:

1. a jest elementem najmniejszym w X , gdy dla każdego $x \in X$ $a \leq x$
2. a jest elementem minimalnym w X , gdy jest jedynym elementem najmniejszym w X
3. a jest elementem największym w X , gdy dla każdego $x \in X$ $x \leq a$
4. a jest elementem maksymalnym w X , gdy jest jedynym elementem największym w X

Definicja 2.6.5. *Ograniczenie górne [3, Rozdział 2]*

Niech $A \subseteq X$, gdzie (X, \leq) jest zbiorem uporządkowanym. Element $x \in X$ nazywamy ograniczeniem górnym zbioru A względem relacji \leq , gdy dla każdego $a \in A$, $a \leq x$.

Definicja 2.6.6. *Ograniczenie dolne [3, Rozdział 2]*

Niech $A \subseteq X$, gdzie (X, \leq) jest zbiorem uporządkowanym. Element $y \in X$ nazywamy ograniczeniem dolnym zbioru A względem relacji \leq , gdy dla każdego $a \in A$, $y \leq a$.

Definicja 2.6.7. *Zbiór ograniczony [3, Rozdział 2]*

Niech $A \subseteq X$, gdzie (X, \leq) jest zbiorem uporządkowanym. Zbiór nazywamy ograniczonym z góry (ograniczonym z dołu), jeśli ma on ograniczenie górne (dolne). Zbiór ograniczony z dołu i z góry nazywamy ograniczonym.

Definicja 2.6.8. *Kres górny [3, Rozdział 2]*

Niech $A \subseteq X$, gdzie (X, \leq) jest zbiorem uporządkowanym. Jeśli zbiór A jest ograniczony z góry i wśród ograniczeń górnych zbioru A istnieje element najmniejszy x_0 , to element ten nazywamy kresem górnym zbioru A i oznaczamy symbolem $\sup A$. Tak więc $x_0 = \sup A$, gdy spełnione są następujące warunki:

1. dla każdego $a \in A$ $a \leq x_0$,
2. $\forall x \in X$ $(\forall a \in A \ a \leq x) \Rightarrow x_0 \leq x$

Definicja 2.6.9. *Kres dolny [3, Rozdział 2]*

Niech $A \subseteq X$, gdzie (X, \leq) jest zbiorem uporządkowanym. Jeśli zbiór A jest ograniczony z dołu i wśród ograniczeń dolnych zbioru A istnieje element największy x_0 , to element ten nazywamy kresem dolnym zbioru A i oznaczamy symbolem $\inf A$. Tak więc $y_0 = \inf A$, gdy spełnione są następujące warunki:

1. $y_0 \leq a$ dla każdego $a \in A$,
2. $\forall y \in X$ $(\forall a \in A \ y \leq a) \Rightarrow y \leq y_0$

Definicja 2.6.10. *Zbiory równoliczne [3, Rozdział 5]*

Mówimy, że zbiory A i B są równoliczne (tej samej mocy), gdy istnieje bijekcja, tj. funkcja f różnowartościowa, przekształcająca zbiór A na zbiór B , tzn. $f : A \rightarrow B$. Piszemy wtedy: $\overline{A} = \overline{B}$.

Definicja 2.6.11. *Zbiór skończony [3, Rozdział 5]*

Mówimy, że zbiór A jest skończony, gdy jest pusty lub równoliczny ze zbiorem $\{1, \dots, n\}$, dla pewnego $n \in \mathbb{N}$. Gdy zbiór jest równoliczny ze zbiorem $\{1, \dots, n\}$, to mówimy że jest on n -elementowy, tj. mocy równej n .

Definicja 2.6.12. *Zbiór przeliczalny [3, Rozdział 5]*

Mówimy, że zbiór X jest przeliczalny, gdy jest skończony lub jest równoliczny z \mathbb{N} .

2.6.2 Przestrzenie metryczne, miary odległości

Niezbędnym jest również wprowadzenie podstawowych pojęć z topologii, ze względu na stosowanie funkcji odległości w celu uporządkowania obiektów.

Definicja 2.6.13. *Metryka [7, Rozdział 9]*

Niech X będzie niepustym zbiorem, wtedy funkcję $d : X \times X \rightarrow [0, \infty)$, nazywamy metryką jeśli spełnione są warunki:

1. $\forall x, y \in X \quad (d(x, y) = 0 \iff x = y),$
2. $\forall x, y \in X \quad d(x, y) = d(y, x),$
3. $\forall x, y, z \in X \quad d(x, y) \leq d(x, z) + d(z, y).$

Definicja 2.6.14. *Przestrzeń metryczna [7, Rozdział 9]*

Niech X będzie niepustym zbiorem, d metryką, wówczas parę (X, d) nazywamy przestrzenią metryczną.

Przykład 3. *Metryka euklidesowa w \mathbb{R}^2*

Niech $d_e : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ będzie metryką euklidesową, wówczas

$$\forall (x_1, x_2), (y_1, y_2) \in \mathbb{R}^2 \quad d_e((x_1, y_1), (x_2, y_2)) := \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}.$$

Przykład 4. *Metryka miejska (Manhattan) w \mathbb{R}^2*

Niech $d_m : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ będzie metryką miejską, wówczas

$$\forall (x_1, x_2), (y_1, y_2) \in \mathbb{R}^2 \quad d_m((x_1, y_1), (x_2, y_2)) := |x_1 - x_2| + |y_1 - y_2|.$$

Przykład 5. *Przestrzeń euklidesowa n -wymiarowa \mathbb{R}^n*

Niech $d_e : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ będzie metryką euklidesową, wówczas

$$\forall (x_1, x_2, \dots, x_n), (y_1, y_2, \dots, y_n) \in \mathbb{R}^n \quad d_e(x, y) := \sqrt{\sum_{i=1}^n |x_i - y_i|^2}.$$

Rozdział 3

Metody porządkowania

Rozdział ten został opracowany w oparciu o pozycję [9, Rozdział 2]. Omówimy w nim wybrane metody porządkowania zarówno liniowego jak i nieliniowego. W kolejnym rozdziale szczegółowo przyjrzymy się wybranym metodom wraz z przedstawieniem ich algorytmów oraz dokładnych opisów matematycznych.

Metody porządkowania liniowego umożliwiają utworzenie uporządkowanej listy obiektów na podstawie określonego kryterium (np. wartości zmiennych). Z kolei metody porządkowania nieliniowego zwracają graf połączeń podobnych obiektów, ze względu na opisujące je zmienne.

3.1 Metody porządkowania liniowego

W wielowymiarowej przestrzeni zmiennych, porządkowanie liniowe obiektów sprowadza się do rzutowania na prostą punktów, które reprezentują obiekty poddane porządkowaniu. Taka operacja pozwala na ustalenie hierarchii obiektów.

Poniżej zostaną przedstawione własności uporządkowania liniowego obiektów, wraz z podaniem ich matematycznej interpretacji.

Obiekty uporządkowane liniowo charakteryzują się tym, że:

- każdy obiekt ma przynajmniej jednego sąsiada i nie więcej niż dwóch sąsiadów,
- jeżeli sąsiadem i -tego obiektu jest k -ty obiekt, to jednocześnie sąsiadem k -tego obiektu jest i -ty obiekt,
- dokładnie dwa obiekty mają tylko jednego sąsiada.

Powyżej wymienione własności są wynikiem posiadania jedynie skończonej ilości obiektów, które poddane są uporządkowaniu. W następnej części chcielibyśmy:

- sformalizować rozumienie powyższych własności,
- udowodnić ich poprawność,
- rozważyć dostateczność tych własności w zbiorach o skończonej ilości obiektów.

Na początku zaczniemy od sprecyzowania takich pojęć jak sąsiad względem relacji.

Definicja 3.1.1. *Sąsiad względem relacji \leq*

Niech X będzie niepustym zbiorem, a x, y będą dwoma różnymi elementami należącymi do tego zbioru. Mówimy, że $y \in X$ jest sąsiadem $x \in X$, co zapisujemy ySx , jeśli

$$(y \leq x \vee x \leq y) \quad \wedge \quad (\neg \exists_{z \in X} \quad x \neq z \neq y \Rightarrow y \leq z \leq x \vee x \leq z \leq y).$$

Twierdzenie 3.1.2. *Własności porządku liniowego w zbiorach skończonych*

Niech \leq będzie relacją porządku liniowego zdefiniowaną w X , gdzie X jest zbiorem ze skończoną liczbą obiektów, złożonym co najmniej z dwóch elementów. Wtedy:

1. $\forall x \in X \quad \overline{\{y \in X, ySx\}} \in \{1, 2\},$
2. $\forall x, y \in X \quad ySx \Rightarrow xSy,$
3. $\overline{\{x \in X, \overline{\{y \in X, \overline{ySx} = 1\}} = 1\}} = 2,$

gdzie S oznacza sąsiada względem relacji \leq .

Dowód. Poniżej zostaną udowodnione powyższe własności.

1. Niech $x \in X$. Przypuśćmy na początek, że $\overline{\{y \in X, ySx\}} = 0$, tzn. że obiekt x nie posiada sąsiadów w tej relacji. Nasz zbiór X jest jednak co najmniej dwuelementowy, zatem istnieje element $y \in X$ i $x \neq y$. Wobec spójności linowego porządku z Definicji 2.6.3 zachodzi wtedy

$$x \leq y \vee y \leq x.$$

Jednak wiemy, że y nie może być sąsiadem x gdyż ten nie posiada sąsiadów. Zatem z definicji sąsiada musi istnieć $z \in X$ różny od obu $x \neq z \neq y$, spełniający warunek

$$z \leq x \vee x \leq z.$$

Powyższe rozumowanie dla y można by dalej zastosować do z , uzyskując kolejne z_1 a później z_2, z_3, \dots dowolną ilość różnych elementów, z których każdy występuje w relacji liniowego porządku z x , ale żaden z nich nie jest sąsiadem. Jednak nasz zbiór X jest zbiorem skończonym, więc nigdy nie uda nam się utworzyć dowolnej ilości różnych elementów ze zbioru X (elementy się wyczerpią). Zatem nasze przypuszczenie, że $\overline{\{y \in X, ySx\}} = 0$ jest fałszywe.

Przypuśćmy dalej, że $\overline{\{y \in X, ySx\}} \geq 3$. Niech a, b, c będą trzema różnymi elementami z X będącymi sąsiadami dla x . Wtedy bez straty o gólności możemy przyjąć, że $a \leq x, b \leq x$ lub $x \leq a, x \leq b$. Istotnie mając 3 elementy w relacji wtedy co najmniej dwa muszą znajdować się po zgodnej stronie, a z dokładnością do o znaczeń możemy przyjąć, że będą nimi a o raz b . Ustalmy zatem, że $a \leq x, b \leq x$. Wobec definicji 2.6.3 wiemy, że $a \leq b$ lub $b \leq a$. Jeśli $a \leq b$ to $a \leq b \leq x$. Co przeczy temu, że a jest sąsiadem x . Jeśli $b \leq a$ to $b \leq a \leq x$ co przeczy temu, że b jest sąsiadem. Analogicznie postępujemy dla przypadku $x \leq a, x \leq b$. Uzyskujemy zatem sprzeczność, będącą efektem przypuszczenia, że mogą istnieć takie 3 elementy a, b, c . Zatem ostatecznie $\overline{\{y \in X, ySx\}} \in \{1, 2\}$.

2. Niech $x, y \in X$ o raz niech ySx . Korzystając z definicji sąsiada 3.1.1 mamy, że skoro ySx to

$$(y \leq x \vee x \leq y) \quad \wedge \quad (\neg \exists z \in X \quad x \neq z \neq y \Rightarrow y \leq z \leq x \vee x \leq z \leq y).$$

Natomiast xSy oznacza, że

$$(x \leq y \vee y \leq x) \quad \wedge \quad (\neg \exists z \in X \quad y \neq z \neq x \Rightarrow x \leq z \leq y \vee y \leq z \leq x).$$

Wobec powyższego widać, że te dwa zdania znaczą to samo, stąd widać że $ySx \Rightarrow xSy$.

3. Intuicyjnie te dwa elementy posiadające po jednym sąsiadzie są elementami maksymalnym i minimalnym w tym zbiorze. Udowodnimy kolejno:

- Element minimalny w zbiorze ma pojedynczego sąsiada. Zauważmy, że zbiór musi posiadać dokładnie 1 element minimalny, tzn. $x_m \in X$ takie, że

$$\forall x \in X \quad x_m \leq x.$$

Istotnie przypuśćmy, że nie istnieje element minimalny. Niech x_1 będzie dowolnym elementem z X . Skoro nie istnieje element minimalny, to istnieje $x_2 \in X$ takie, że $x_2 \leq x_1$ i $x_2 \neq x_1$. Dla x_2 z braku elementu minimalnego, musi istnieć z kolei $x_3 \leq x_2$ takie, że $x_3 \neq x_2$. Itd. Co nie jest możliwe, gdyż zbiór X jest przecież skończonym zbiorem. Rozważmy dalej przypuszczenie gdyby były dwa lub więcej takich elementów. Wtedy to, z antysymetryczności, oczywiście musiałyby być sobie równe. Jeśli x_m, y_m są jednocześnie minimalne to

$$\forall x \in X \quad x_m \leq x,$$

oraz

$$\forall x \in X \quad y_m \leq x.$$

Skąd natychmiast mamy, że $x_m \leq y_m$ oraz $y_m \leq x_m$. Wobec antysymetryczności z definicji 2.6.2 mamy, że $x_m = y_m$ wbrew naszemu przypuszczeniu, że są od siebie różne. Pozostaje pokazać, że element minimalny ma pojedynczego sąsiada. Przypuśćmy, że $y, z \in X$ są dwoma różnymi sąsiadami dla x_m . Wtedy $x_m \leq y \vee y \leq x_m$ oraz $x_m \leq z \vee z \leq x_m$. Skoro x_m jest minimalny to musi to oznaczać, że

$$x_m \leq y \wedge x_m \leq z.$$

Wobec spójności z definicji 2.6.3 zachodzi $y \leq z$ lub $z \leq y$. Sprzeczność, gdyż wtedy któryś z nich nie mógłby być sąsiadem dla x_m .

- Element maksymalny x_M w zbiorze ma pojedynczego sąsiada. Analogicznie do powyższego punktu, zbiór musi posiadać dokładnie 1 element maksymalny, tzn. $x_M \in X$ takie, że

$$\forall x \in X \quad x \leq x_M.$$

Istotnie przypuśćmy, że nie istnieje element maksymalny. Niech x_1 będzie dowolnym elementem z X . Skoro nie istnieje element maksymalny, to istnieje $x_2 \in X$ takie, że $x_1 \leq x_2$ i $x_1 \neq x_2$. Dla x_2 z braku elementu maksymalnego, musi istnieć taki element $x_3 \in X$ i $x_3 \neq x_2$, że $x_2 \leq x_3$. Itd. Co nie jest możliwe, gdyż z założenia zbiór X jest skończonym zbiorem. Rozważmy dalej przypuszczenie gdyby były dwa lub więcej takich elementów. Wtedy to z antysymetryczności, musiałyby być sobie równe. Jeśli x_M, y_M są jednocześnie maksymalne, to

$$\forall x \in X \quad x \leq x_M,$$

oraz

$$\forall x \in X \quad x \leq y_M.$$

Stąd natychmiast mamy, że $x_M \leq y_M$ oraz $y_M \leq x_M$. Wobec antysymetryczności z definicji 2.6.2, mamy że $x_M = y_M$, co wbrew naszemu przypuszczeniu daje, że elementy te nie są od siebie różne. Pozostaje pokazać, że element maksymalny ma pojedynczego sąsiada. Przypuśćmy, że $y, z \in X$ są dwoma różnymi sąsiadami dla x_M . Wtedy $x_M \leq y \vee y \leq x_M$ oraz $x_M \leq z \vee z \leq x_M$. Skoro x_M jest elementem maksymalny to musi to zatem oznaczać

$$y \leq x_M \wedge z \leq x_M.$$

Wobec spójności z definicji 2.6.3 zachodzi $y \leq z$ lub $z \leq y$. Sprzeczność, gdyż wtedy któryś z nich nie mógłby być sąsiadem dla x_M .

- Żaden inny element nie może mieć pojedynczego sąsiada. Przypuśćmy, że $x \in X$ nie będąc ani elementem minimalnym ani maksymalnym ma pojedynczego sąsiada. Wobec definicji elementu minimalnego i maksymalnego oraz spójności zachodzi

$$x_m \leq x \leq x_M.$$

Zatem albo x_m jest sąsiadem x albo istnieje $y_1 \in X$ taki, że $y_1 \leq x$. Tworzy to kilka możliwych przypadków. W pierwszym x_m będzie tym jedynym sąsiadem, w drugim x_M nim będzie, w ostatnim natomiast, ani x_m , ani x_M nie będą sąsiadami.

Zajmiemy się najpierw pierwszym z nich, tj. x_m jest sąsiadem x . Zauważmy teraz, że z faktu, iż x_M jest elementem maksymalnym zbioru X , wynika że $x \leq x_M$. Nie jest jednak sąsiadem elementu x . Zatem istnieje takie $x_1 \in X$, że $x \leq x_1 \leq x_M$. Jednak x_1 również nie może być sąsiadem x co powoduje, że istnieje taki element $x_2 \in X$, że $x \leq x_2 \leq x_1 \leq x_M$. Itd. Jednakże, skoro zbiór X jest zbiorem skończonym, to musi istnieć taki element $x_j \in X$, że $x \leq x_j$ i $x_j \neq x_m$, który będzie sąsiadem z x , zatem xSx_j . Zatem ostatecznie xSx_m i xSx_j , a to przeczy założeniu, że x ma pojedynczego sąsiada.

Przejdźmy teraz do drugiego przypadku, tj. gdy x_M jest sąsiadem dla x . Wtedy x_m nie jest sąsiadem dla x jednak wiedząc, że $x_m \leq x$ musi istnieć $y_1 \in X$, taki że $y_1 \leq x$. Jednak wiedząc iż y_1 nie jest sąsiadem dla x wnioskujemy, że istnieje taki $y_2 \in X$, że $x_m \leq y_1 \leq y_2 \leq x$. Itd. Jednakże, skoro zbiór X jest zbiorem skończonym, to musi istnieć taki element $y_i \in X$, że $y_i \leq x$, który będzie sąsiadem z x . Zatem ostatecznie y_iSx i xSx_M , co przeczy założeniu że x ma pojedynczego sąsiada.

Zajmijmy się teraz trzecim przypadkiem, tj. gdy ani x_m oraz x_M nie są sąsiadami elementu x . Z faktu, iż zbiór X posiada element minimalny x_m , który nie jest sąsiadem elementu x , wynika że istnieje taki element $x_1 \in X$, że $x_m \leq x_1 \leq x$. Co więcej istnieje taki $x_2 \in X$, że $x_m \leq x_1 \leq x_2 \leq x$. Itd. I znów skoro zbiór X jest skończony, to istnieje taki element $x_j \in X$, że $x_j \leq x$ i x_jSx . Z drugiej strony, skoro zbiór X posiada element maksymalny x_M , który nie jest sąsiadem elementu x , wynika że istnieje taki element $y_1 \in X$, że $x \leq y_1 \leq x_M$. Analogicznie do wcześniejszych kroków, istnieje taki element $y_2 \in X$, że $x \leq y_2 \leq y_1 \leq x_M$. Itd. Zbiór X jest zbiorem skończonym, zatem musi istnieć taki element $y_i \in X$, że $x \leq y_i$ i xSy_i . Łącząc te dwa warunki, wynika że x musi mieć dwóch sąsiadów. Co kończy dowód własności.

□

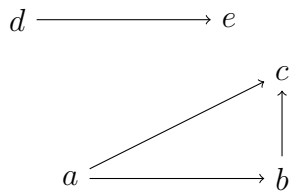
Własności, które wykazaliśmy powyżej są często podawane niemal na równi z definicją takiego uporządkowania. Poniżej zostaną jednak podane przykłady takich relacji, które mimo, że posiadają powyższy zestaw własności, to nie opisują relacji będących porządkami liniowymi.

Przykład 6. Rozważmy zbiór dwuelementowy $X = \{a, b\}$ gdzie $a \leq b$ jest jedynym punktem tej relacji. Tak zdefiniowana relacja spełnia wszystkie własności, ale nie spełnia założenia o zwrotności - zatem relacja ta nie jest liniowym porządkiem. Diagram Hassego prezentujący tę relację, jest postaci:



Przykład 7. Rozważmy zbiór $X = \{a, b, c, d, e\}$ oraz relację definiującą następujące sąsiedztwa (wypisaną bez par symetrycznych) aSb, bSc, aSc, dSe . Ponadto dołożymy warunek zwrotności, tj. $a \leq a, b \leq b, c \leq c, d \leq d, e \leq e$. Diagram Hassego prezentujący relację porządku tego zbioru,

jest postaci:



Z diagramu widać, że taka relacja spełnia wszystkie omawiane wcześniej własności - jednak nie jest spójna. I tak np. nie możemy porównać elementów a i d , bowiem nie możemy określić czy $d \leq a$ lub $a \leq d$.

W podsumowaniu tej sekcji należy podkreślić, że by uporządkować liniowo obiekty z macierzy obserwacji, charakteryzujące je zmienne muszą być mierzone przynajmniej na skali porządkowej. Istotne jest również aby miały jednakowy charakter. Na potrzeby pracy zakładamy, że zmienne opisujące obiekty powinny być stymulantami. Dlatego też gdy nimi nie są, należy poddać je np. stymulacji. Operacja ta umożliwi w dalszym kroku przejście do transformacji normalizacyjnej, która konieczna jest gdy zmienne opisujące obiekty mierzone są na skali przedziałowej lub ilorazowej, a chcemy uzyskać ich porównywalność.

Metody porządkowania liniowego można podzielić na metody diagramowe, procedury oparte na zmiennej syntetycznej oraz procedury iteracyjne bazujące na funkcji kryterium dobroci uporządkowania. W kolejnej sekcji zostaną pokrótce przedstawione różne metody, z wyszczególnieniem najważniejszych założeń o każdej z nich.

3.1.1 Metody diagramowe

W metodach diagramowych stosuje się graficzną reprezentację macierzy odległości zwanej diagramem. Macierz konstruowana jest w oparciu o odległości między obiektami, wyznaczone za pomocą dowolnej metryki. Porządkowanie obiektów polega na porządkowaniu diagramu, tzn. przestawieniu wierszy i odpowiadających im kolumn, aby wzdłuż przekątnej skupiały się najmniejsze odległości zaś im dalej od głównej przekątnej tym większe odległości między zmiennymi opisującymi porządkowane obiekty. Narzędzie pomocnicze w porządkowaniu danych, może stanowić kryterium postaci:

$$F^1 = \sum_{i=1}^n \sum_{k>1}^n d_{ik} w_{ik}$$

gdzie:

$D = [d_{ik}]$ - macierz odległości między i -tym i k -tym obiektem,

$W = [w_{ik}]$ $i, k = 1, 2, \dots, n$ - macierz wag elementów macierzy odległości, wymiaru równego wymiarowi macierzy odległości.

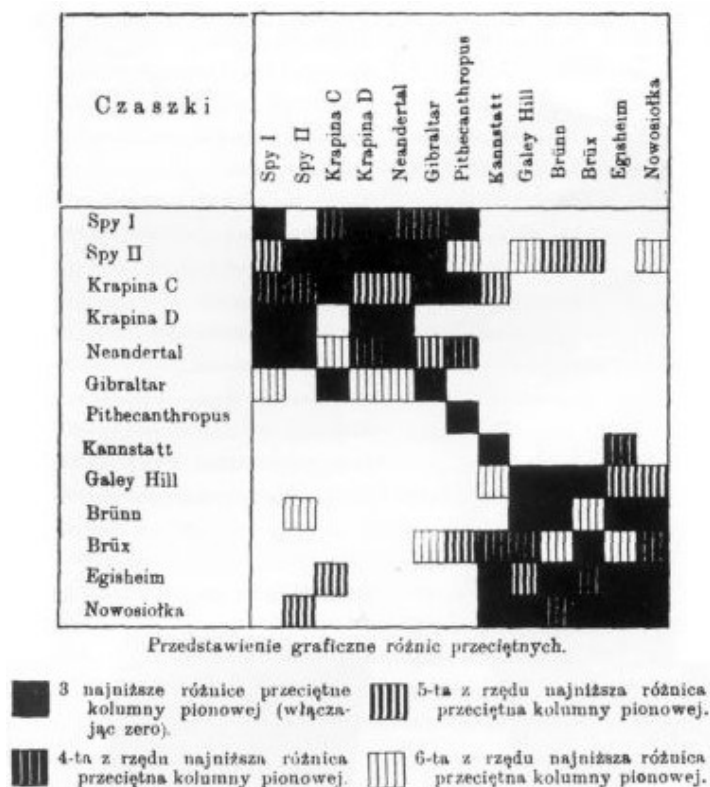
Elementy macierzy wag wyznaczone są za pomocą jednego z poniższych wzorów:

(a) $w_{ik} = \frac{|i-k|}{n-1},$

(b) $w_{ik} = \frac{1}{n(n-1)}[2n|i-k-1| + i+k - (i-i)^2],$

(c) $w_{ik} = \frac{1}{n(n-1)}[2n|i-k| + 2 - i - k - (i-i)^2].$

Zaprezentujmy teraz przykład uporządkowanego diagramu, przedstawiającego wynik badań Jana Czekanowskiego, dotyczących metod badania różnic między kopalnymi czaszkami ludzkimi. Analizując diagram zauważamy że wzdłuż głównej przekątnej skupiają się najmniejsze odległości między zmiennymi, a im dalej od niej tym odległości zwiększają się. Przykład został znaleziony na stronie [11].



Rysunek 3.1: Diagram opublikowany w podręczniku „Statystyka dla antropologów” Jana Czekanowskiego, 1913r. [11]

Metody diagramowe nie leżą w zakresie zainteresowań tej pracy, z uwagi na mały formalizm matematyczny. Są jednak istotnym narzędziem do wizualizacji porządku.

3.1.2 Metody oparte na zmiennych syntetycznych

W tym podrozdziale zostaną opisane metody porządkowania oparte na zmiennych syntetycznych, tj. funkcji wyznaczonej na podstawie wartości zmiennych opisujących obiekty, której wartości będą służyć do porządkowania zbioru. Metody oparte na zmiennych syntetycznych dzielimy na wzorcowe i bezwzorcowe. Poniżej zostaną one opisane szczegółowo, jednak wcześniej zostaną przedstawione wzory wyznaczające zmienną syntetyczną.

Sposoby wyznaczania zmiennej syntetycznej

W pracy dla zachowania ogólności, przyjmujemy że wszystkie zmienne opisujące obiekty mają jednakowe wagi, w związku z tym wzory służące do wyznaczenia zmiennej syntetycznej są postaci:

1. średniej arytmetycznej:

$$s_i = \frac{1}{m} \sum_{j=1}^m n_{ij}, \quad i = 1, 2, \dots, n,$$

2. średniej geometrycznej:

$$s_i = \prod_{j=1}^m (n_{ij})^{\frac{1}{m}}, \quad i = 1, 2, \dots, n,$$

3. średniej harmonicznej

$$s_i = \left[\sum_{j=1}^m \frac{1}{n_{ij}} \right]^{-1} \cdot m, \quad i = 1, 2, \dots, n,$$

gdzie: s_i - wartość zmiennej syntetycznej w i -tym obiekcie,

Metoda wzorcowa

W metodach tych zakłada się istnienie obiektu wzorcowego $P_0 = [n_{0j}]$, $j = 1, 2, \dots, m$. Zmienne tego obiektu są znormalizowane. Przyjmują one optymalne wartości, które są ustalane na podstawie ogólnie przyjętych norm, subiektywnej opinii dotyczącej obserwowanego obiektu, lub też opinii ekspertów. Wtedy obiekty, porządkowane są na podstawie odległości od obiektu wzorcowego. Poszczególne metody mogą różnić się sposobem wyznaczania obiektu wzorcowego, odległości oraz miary syntetycznej, na podstawie której dokonywane jest porządkowanie.

Metoda Hellwiga

Metoda Hellwiga jest jedną z najstarszych metod wzorcowych. W metodzie tej, obiekt wzorcowy wyznaczony jest na podstawie wystandaryzowanych zmiennych wejściowych. Współrzędnym obiektu wzorcowego przyporządkowuje się maksimum, gdy zmienne wejściowe są stymulantami lub minimum gdy zmienne są destymulantami. Obiekty są uporządkowywane na podstawie odległości od obiektu wzorcowego, przy wykorzystaniu odległości euklidesowej. Do porządkowania używana jest natomiast miara syntetyczna, postaci:

$$s_i = 1 - \frac{d_{i0}}{d_0}, \quad i = 1, 2, \dots, n,$$

gdzie:

d_{i0} - odległość i -tego obiektu, od obiektu wzorcowego

d_0 - wartość wyznaczana dla każdej zmiennej, będąca sumą średniej odległości od obiektu oraz podwójonej wartości odchylenia standardowego dla tej zmiennej. Doprecyzujemy

Współrzędne obiektu wzorcowego są obliczane na podstawie wzoru:

$$n_{0j} = \begin{cases} \max_i(n_{ij}) & \text{gdy } n_j \text{ jest stymulantą, } j = 1, 2, \dots, m, \quad i = 1, 2, \dots, n, \\ \min_i(n_{ij}) & \text{gdy } n_j \text{ jest destymulantą, } j = 1, 2, \dots, m, \quad i = 1, 2, \dots, n. \end{cases}$$

Dalej wyznaczamy miarę syntetyczną, licząc kolejno:

- Dla każdego obiektu, wyznaczana jest odległość od obiektu wzorcowego:

$$d_{i0} = \left[\sum_{j=1}^m (n_{ij} - n_{0j})^2 \right]^{\frac{1}{2}}$$

- Następnie dla każdej zmiennej wyznaczana jest średnia odległość od obiektu wzorcowego wg. wzoru:

$$\bar{d}_0 = \frac{1}{n} \sum_{i=1}^n d_{i0}$$

- W kolejnym kroku wyznaczamy odchylenie standardowe dla każdej zmiennej za pomocą wzoru:

$$\sigma(d_0) = \left[\frac{1}{n} \sum_{i=1}^n (d_{i0} - \bar{d}_0)^2 \right]^{\frac{1}{2}}$$

- Mając powyższe, możemy wyznaczyć wartość d_0 jako sumę średniej odległości oraz podwojonej wartości odchylenia standardowego:

$$d_0 = \overline{d_0} + 2\sigma(d_0)$$

Wartości miary s_i zazwyczaj są z przedziału $[0, 1]$, jednakże mogą zdarzyć się wartości ujemne. Należy tu zaznaczyć, że wartości miary są tym wyższe, im mniej jest oddalony obiekt od obiektu wzorcowego.

Metoda dystansowa

Podobnie jak we wcześniejsze metodach, na początku zmienne należy poddać stymulacji, oraz przekształceniu normalizacyjnemu, wybranemu na podstawie skal do których należą zmienne opisujące obiekty. W kolejnym kroku wyznaczane są współrzędne obiektu wzorcowego, a następnie macierz odległości każdego obiektu od obiektu wzorcowego. Odległość od obiektu wzorca jest wyznaczana przy zastosowaniu dowolnej metryki, np. metryki euklidesowej. Dla metody tej, miara syntetyczna jest wyznaczana za pomocą przekształcenia unitaryzacyjnego postaci:

$$s_i = \left(\frac{d_{i0} - \min_i(d_{i0})}{\max_i(d_{i0}) - \min_i(d_{i0})} \right)^p, \quad i = 1, 2, \dots, n, \quad p \in \mathbb{N}.$$

Miara syntetyczna uzyskana tą metodą jest unormowana i przyjmuje wartości z przedziału: $[0, 1]$. Czym niższa wartość miary, tym bliżej obiektu wzorcowego leży dany obiekt.

Metody bezwzorcowe

W metodach tych zakładamy, że nie istnieje obiekt wzorcowy. Porządkowanie dokonywane jest na podstawie wartości zmiennej syntetycznej, wyznaczonej dla każdego obiektu. Poniżej zostaną omówione wybrane metody porządkowania bezwzorcowego.

Metoda rang

Definicja 3.1.3. Ranga [9, Rozdział 1.5]

Rangą nazywamy zmienną o wartościach ze zbioru $\mathbb{R}_+ \setminus \{0\}$, będącą najczęściej liczbą całkowitą. Wartości reprezentują numery miejsc obiektów po uporządkowaniu malejąco.

Metoda ta opiera się na normalizacji rangowej, w związku z tym zmienne poddane porządkowaniu, powinny być mierzone na skali porządkowej. Dla każdego obiektu wyznacza się sumę przyporządkowanych mu rang ze względu na wszystkie zmienne. Na końcu obliczana jest wartość zmiennej syntetycznej, jako średniej wartości rang. W oparciu o tę wartość następuje porządkowanie obiektów, tj. im wartość zmiennej syntetycznej jest mniejsza tym wyżej w hierarchii znajduje się uporządkowany obiekt. Wzór na obliczenie wartości zmiennej syntetycznej:

$$s_i = \frac{1}{m} \sum_{j=1}^m n_{ij}, \quad i = 1, 2, \dots, n,$$

gdzie:

n_{ij} -zmienna znormalizowana rangowo, tj. $n_{ij} = r$ dla $x_{rj} = x_{ij}$, $r, i = 1, 2, \dots, n$,

r -ranga nadana i -temu obiektowi znajdującemu się na r -tym miejscu w uporządkowanym szeregu obiektów ze względu na j -tą zmienną.

Metoda sum

Metoda ta używana jest w momencie, gdy zmienne mierzone są na skali ilorazowej lub przedziałowej. W związku z tym tuż po stymulacji zmiennych, należy dokonać przekształcenia normalizacyjnego, za pomocą unitaryzacji. W kolejnym kroku, dla każdego obiektu wyznaczana jest zmienna syntetyczna, jako średnia arytmetyczna wartości zmiennych przy przyjęciu jednakowych wag dla każdej zmiennej. Następnie muszą zostać wyeliminowane ujemne wartości zmiennej syntetycznej, do czego służy poniższe przekształcenie:

$$s'_i = s_i - \min_i(s_i), \quad i = 1, 2, \dots, n.$$

Końcowa postać zmiennej syntetycznej otrzymywana jest, przy wykorzystaniu normalizacji postaci:

$$s''_i = \frac{s'_i}{\max_i(s'_i)}, \quad i = 1, 2, \dots, n.$$

Powyższe przekształcenia ujednolicają zakres miary syntetycznej do przedziału $[0, 1]$. Im wyższa wartość zmiennej syntetycznej, tym wyżej w hierarchii znajduje się obiekt.

3.1.3 Metody iteracyjne

W metodach tych przyjmowania jest funkcja kryterium dobroci porządkowania, dla której w kolejnych iteracjach poszukiwane jest takie uporządkowanie liniowe obiektów, które optymalizuje wartość funkcji kryterium, aż do osiągnięcia przez nią wartości optymalnej tj. maksymalnej lub minimalnej.

Metoda Szczotki

Metoda ta polega na znalezieniu takiego uporządkowania liniowego obiektów, dla którego funkcja kryterium dobroci uporządkowania osiąga maksimum:

$$F^2 = \sum_{k=1}^{n-1} k \sum_{i=1}^{n-k} d_{ik} \rightarrow \max$$

gdzie:

d_{ik} - odległość euklidesowa między i -tym i k -tym obiektem.

W pierwszym kroku działania tej metody, przeprowadzane jest dowolne liniowe uporządkowanie obiektów. Następnie, dla tego uporządkowania obliczana jest wartość funkcji kryterium dobroci uporządkowania, według powyższego wzoru. W kolejnych etapach wyznaczana jest wartość tej funkcji, dla każdej transpozycji pary obiektów. Powyższe kroki wykonywane są do momentu, gdy dowolna transpozycja pary obiektów nie spowoduje zwiększenia wartości funkcji kryterium dobroci uporządkowania.

3.2 Metody porządkowania nieliniowego

Metody porządkowania nieliniowego w odróżnieniu od metod porządkowania liniowego, polegają nie na uporządkowaniu obiektów w sposób hierarchiczny, a na określeniu dla każdego z nich, stopnia podobieństwa z innymi obiektami, na podstawie opisujących je zmiennych.

Aby zastosować metody porządkowania nieliniowego, zmienne opisujące obiekty, powinny być mierzone na skali przedziałowej lub ilorazowej. Gdy zmienne te mierzone są na skali przedziałowej lub ilorazowej, należy dokonać ich normalizacji.

Metody porządkowania nieliniowego dzielimy na metody dendrytowe oraz aglomeracyjne. Metody dendrytowe prowadzą do powstania dendrytu, prezentującego położenie obiektów ze względu na ich podobieństwo między sobą. Z kolei metody aglomeracyjne sprowadzają się do powstania dendrogramu lub łańcucha połączeń, prezentując sposób łączenia obiektów do siebie podobnych.

3.2.1 Metody dendrytowe

Definicja 3.2.1. *Dendryt [9, Rozdział 2.3]*

Dendrytem nazywamy acykliczny graf spójny, bez pętli.

Metody dendrytowe opierają się na pojęciach teorii grafów. Metody te polegają na stworzeniu dendrytu, którego wierzchołki odpowiadają odpowiednim obiektom poddanym porządkowaniu. Krawędzie łączące wierzchołki odpowiadają zaś odległością między obiektami. Przykładem metod dendrytowych jest taksonomia wrocławska, metoda Prima. Poniżej zostanie jednak opisana jedynie metoda taksonomii wrocławskiej.

Taksonomia wrocławska

W metodzie tej obiekty dzielone są na grupy obiektów najbardziej do siebie podobnych, tj. takich, dla których odległość między sobą jest jak najmniejsza. W związku z tym w pierwszej kolejności należy wyznaczyć macierz odległości obiektów D np. przy użyciu metryki euklidesowej, a następnie w każdym wierszu (kolumnie) macierzy, wyznaczamy jest element najmniejszy:

$$d_{ik} = \min_k d_{ik}, \quad i, k = 1, 2, \dots, n, i \neq k.$$

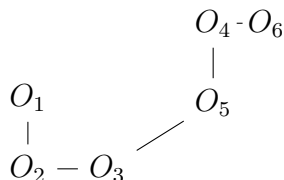
Za pomocą grafu prezentowane są pary obiektów, najbardziej do siebie podobnych. W grafie tym, długość krawędzi łączących wierzchołki (czyli obiekty poddane porządkowaniu) odpowiadają odległości między parą obiektów. Jeżeli wśród połączonych par obiektów, pojawiają się krawędzie dwukrotne, należy je wyeliminować, ze względu na to, że kolejność połączeń w dendrycie nie jest istotna. Obiekty w dendrycie nie mogą się powtarzać, jeżeli natomiast niektóre obiekty w łączeniu wystąpią wielokrotnie, to obiekty te zostaną połączone w zespoły zwane skupieniami. Metoda ta kończy swoje działanie, w momencie uzyskania grafu spójnego.

Zaprezentujemy teraz tę metodę na przykładzie.

Przykład 8. Niech dana będzie macierz odległości $D = [d_{ik}]_{i \leq 6, k \leq 6}$ sześciu różnych obiektów $\{O_1, O_2, \dots, O_6\}$, w której to w kolejnych wierszach znajdują odległości i -tego obiektu od obiektu k -tego:

$$D = \begin{bmatrix} 0.00 & \underline{0.35} & 0.70 & 0.95 & 2.36 & 2.99 \\ \underline{0.35} & 0.00 & 0.45 & 1.45 & 2.00 & 0.36 \\ 0.70 & \underline{0.45} & 0.00 & 1.05 & 0.90 & 0.75 \\ 0.95 & 1.45 & 1.05 & 0.00 & 0.50 & \underline{0.16} \\ 2.36 & 2.00 & 0.90 & \underline{0.50} & 0.00 & 0.52 \\ 2.99 & 0.36 & 0.75 & \underline{0.16} & 0.52 & 0.00 \end{bmatrix}$$

Mając wyznaczoną macierz odległości, postępujemy według wyżej wskazanej reguły, która prowadzi do uzyskania dendrytu postaci:



W powstałym dendrycie wierzchołkami są obiekty przechodzące zbioru C , z kolei krawędzie łączące te wierzchołki są współrzędne wektora c , które powstały przez wybór najmniejszej odległości między dołączanymi obiektami, w kolejnych etapach dołączania obiektów do dendrytu.

3.2.2 Metody aglomeracyjne

Przed opisem działania metod aglomeracyjnych, wprowadzimy definicję łańcucha połączeń.

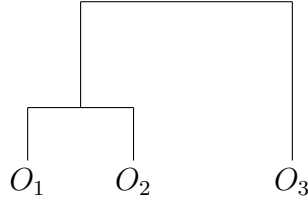
Definicja 3.2.2. *Łańcuch połączeń*

Niech $O = \{O_1, O_2, \dots, O_n\}$, $n \in \mathbb{N}$ oznacza zbiór obiektów. Łańcuchem połączeń nazwiemy skończony ciąg $(C_i)_{i=1}^N$, $N \in \mathbb{N}$, $n \leq N \leq 2^n$, taki że:

1. $C_i = \{O_i\}$ dla $i = 1, 2, \dots, n$
2. $\forall_{i > n} \quad \exists_{j, k \in \mathbb{N}, j \leq k < i} \quad C_i = C_j \cup C_k$
3. $C_N = O = \{O_1, O_2, \dots, O_n\}$
4. $\forall_{i < j, i, j \in \mathbb{N}} \quad C_i \cap C_j \neq \emptyset \Rightarrow C_i \subset C_j$

Uwaga 7. Graficzną reprezentację łańcucha połączeń nazywamy dendrogramem.

Przykład 9. Zaprezentujemy teraz przykład tworzenia łańcucha połączeń. Niech dany będzie zbiór obiektów $O = \{O_1, O_2, O_3\}$. Wtedy przykładowy łańcuch połączeń jest postaci $C = \{\{O_1\}, \{O_2\}, \{O_3\}, \{O_1, O_2\}, \{O_1, O_2, O_3\}\}$. Zaprezentujemy go teraz graficznie, w postaci dendrogramu:



Opis działania metod aglomeracyjnych:

Istotą metod aglomeracyjnych jest utworzenie łańcucha połączeń i dendrogramu. W ten sposób zobrazowana jest kolejność łączenia obiektów, na podstawie zmniejszającego się podobieństwa między obiektami włączonymi do dendrogramu, a tymi wcześniej do niego należącymi. Położenie obiektów oraz grup obiektów, które powstały w kolejnych etapach tworzenia łańcucha, jest przeprowadzone na podstawie kolejności połączeń tych obiektów i grup. Każde ogniwo w łańcuchu oznacza grupy obiektów podobnych do siebie. Wyjściowym założeniem metod aglomeracyjnych jest to, że każdy obiekt stanowi odrębną, jednoelementową grupę $(C_i, i = 1, 2, \dots, n)$.

W kolejnych krokach następuje łączenie ze sobą grup obiektów najbardziej podobnych do siebie ze względu na ich zmienne. Podobieństwo weryfikowane jest na podstawie odległości między grupami.

Na początku odległości między jednoelementowymi grupami C_1, \dots, C_n wyznacza wyjściowa macierz odległości D . W macierzy D poszukiwane są najmniejsze odległości pomiędzy grupami obiektów:

$$d_{ii'} = \min_{ik} d_{ik}, \quad i = 1, 2, \dots, n_i, \quad k = 1, 2, \dots, n_{i'}, \quad i, i' = 1, 2, \dots, n, i \neq i'.$$

gdzie:

$d_{ii'}$ - odległość i -tej od i' -tej grupy.

W kolejnym kroku, obiekty o najmniejszej odległości między sobą łączone są w jedną grupę, dzięki czemu liczba grup zmniejsza się o jeden. Zostaje rozpoczęty proces tworzenia łańcucha połączeń. Ponownie badane są odległości między nowo stworzoną grupą, a pozostałymi grupami. Proces trwa do momentu stworzenia pełnego łańcucha połączeń, tj. jednej grupy.

Ogólna postać wzoru służącego do wyznaczenia odległości nowo powstałej grupy $C_{i''}$, (powstałej dzięki połączeniu grup C_i i $C_{i'}$), od grup które zostały $C_{i'''}$ to:

$$d_{i''i'''} = \alpha_i d_{ii'} + \alpha_{i'} d_{i'i'''} + \beta d_{ii'''} + \gamma |d_{ii'''} - d_{i'i'''}|$$

gdzie: $\alpha_i, \alpha_{i'}, \beta, \gamma$ - współczynniki przekształceń, różne dla poszczególnych metod aglomeracyjnych

Możemy wyróżnić sześć różnych metod aglomeracyjnych, różniących się sposobem wyznaczenia odległości między grupami obiektów. Poniżej zostaną podane współczynniki przekształceń dla każdej z nich, a w dalszej części pracy wybrane z nich zostaną szczegółowo omówione.

- metoda najbliższego sąsiedztwa (metoda pojedynczego wiązania):
parametry przekształceń $\alpha_i = 0,5 \quad \alpha_{i'} = 0,5 \quad \beta = 0 \quad \gamma = 0,5$.
- metoda najdalszego sąsiedztwa (metoda pełnego wiązania):
parametry przekształceń $\alpha_i = 0,5 \quad \alpha_{i'} = 0,5 \quad \beta = 0 \quad \gamma = -0,5$.
- metoda średniej międzygrupowej (metoda średnich połączeń):
parametry przekształceń $\alpha_i = \frac{n_i}{n_i + n_{i'}} \quad \alpha_{i'} = \frac{n_{i'}}{n_i + n_{i'}} \quad \beta = 0 \quad \gamma = 0$.
- metoda mediany:
parametry przekształceń $\alpha_i = 0,5 \quad \alpha_{i'} = 0,5 \quad \beta = -0,25 \quad \gamma = 0$.
- metoda środka ciężkości:
parametry przekształceń $\alpha_i = \frac{n_i}{n_i + n_{i'}}; \alpha_{i'} = \frac{n_{i'}}{n_i + n_{i'}} \quad \beta = \frac{-n_i n_{i'}}{(n_i + n_{i'})^2} \quad \gamma = 0$.
- metoda Warda:
parametry przekształceń $\alpha_i = \frac{n_i + n_{i'''}}{n_i + n_{i'} + n_{i'''}} \quad \alpha_{i'} = \frac{n_{i'} + n_{i'''}}{n_i + n_{i'} + n_{i'''}} \quad \beta = \frac{-n_{i'''}}{n_i + n_{i'} + n_{i'''}} \quad \gamma = 0$.

Metoda najbliższego sąsiedztwa

W metodzie tej odległość między dwoma grupami obiektów jest równa odległości pomiędzy najbliższymi obiektami (sąsiadami), które należą do dwóch różnych grup. Odległość ta opisana jest wzorem:

$$d_{ii'} = \min_{ik} d_{ik}(\mathbf{O}_i \in \mathbf{C}_i, \mathbf{O}_k \in \mathbf{C}_{i'}),$$

$$i = 1, 2, \dots, n_i, \quad k = 1, 2, \dots, n_{i'}, \quad i, i' = 1, 2, \dots, n, \quad i \neq i',$$

gdzie:

$$\mathbf{O}_i = [n_{ij}], \quad j = 1, 2, \dots, m.$$

Metoda najdalszego sąsiedztwa

W metodzie tej odległość między dwoma grupami obiektów jest równa odległości pomiędzy najdalszymi obiektami (sąsiadami), które należą do dwóch różnych grup. Odległość ta opisana jest wzorem:

$$d_{ii'} = \max_{ik} d_{ik}(\mathbf{O}_i \in \mathbf{C}_i, \mathbf{O}_k \in \mathbf{C}_{i'}),$$

$$i = 1, 2, \dots, n_i, \quad k = 1, 2, \dots, n_{i'}, \quad i, i' = 1, 2, \dots, n, \quad i \neq i',$$

Metoda średniej międzygrupowej

W metodzie tej odległość między dwoma grupami obiektów równa jest średniej arytmetycznej odległości między wszystkimi parami obiektów należących do dwóch różnych grup. Odległość ta opisana jest wzorem:

$$d_{ii'} = \frac{1}{n_i n_{i'}} \sum_{k=1}^{n_{i'}} \sum_{i=1}^{n_i} d_{ik} (\mathbf{O}_i \in \mathbf{C}_i, \mathbf{O}_k \in \mathbf{C}_{i'})$$

$$i, i' = 1, 2, \dots, n, \quad i \neq i'.$$

Metoda mediany

W metodzie tej odległość między grupami obiektów jest równa medianie odległości pomiędzy wszystkimi parami obiektów należących do dwóch grup. Odległość ta opisana jest wzorem:

$$d_{ii'} = \text{med}_{i,k} \{d_{ik} (\mathbf{O}_i \in \mathbf{C}_i, \mathbf{O}_k \in \mathbf{C}_{i'})\},$$

$$i = 1, 2, \dots, n_i, \quad k = 1, 2, \dots, n_{i'}, \quad i, i' = 1, 2, \dots, n, \quad i \neq i'.$$

Rozdział 4

Zastosowanie wybranych metod porządkowania danych wielowymiarowych

4.1 Opis zbioru

Zbiór danych jest opracowaniem własnym, na podstawie ofert sprzedaży samochodów osobowych, zamieszczonych na portalu *www.otomoto.pl* w okresie listopad - grudzień 2017 roku. Zebrane dane dotyczą szczegółowych informacji odnośnie samochodu, tj. jego marki, modelu, wersji, typu, koloru lakieru, pojemności silnika, roku produkcji, przebiegu, liczby drzwi, rodzaju skrzyni biegu, rodzaju paliwa, rodzaju napędu, wyposażenia w: ABS, komputer pokładowy, ESP, klimatyzację. Oprócz danych ściśle związanych z budową i wyposażeniem samochodu, pojawiły się również atrybuty, tj. cechy umieszczone w kolumnach, związane z informacją o tym czy auto jest uszkodzone oraz bezwypadkowe, czy jest sprowadzane, jaki jest kraj aktualnej rejestracji, czy było serwisowane, czy sprzedający jest pierwszym właścicielem. Dodatkowo oprócz powyższych, został dodany atrybut najbardziej interesujący kupującego - czyli cena oraz województwo tj. miejsce skąd wystawiana została oferta. Zbiór został dołączony do pracy na płycie.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q			
1	MARKA	MODEL	WERSJA	TYP	WOJEW.	CENA.NET	CENA.BR	MOC[km]	POJEMN.	CROK.	PRO.	PRZEBIEG	KOLOR	L.DZ.	RZWI.	RODZAJ	F.SKRZYNIA	NAPED	KRAJ	AKT
2	Hyundai	i20	II	kompakt	malopolskie	46500		90	1396	2016		18300	biały		3	diesel	manualna	na przedni	Polska	
3	Hyundai	i20	I	kompakt	mazowieckie	22900		86	1248	2013		319000	niebieski		5	benzyna+L	manualna	na przedni	Polska	
4	Subaru	Legacy	V	kombi	mazowieckie	36900		150	1998	2010		149000	srebrny-metal		5	benzyna	manualna	4x4(stały)	Polska	
5	Ford	Mondeo	Mk4	sedan	dolnoslaskie	30000		146	1999	2008		166290	czarny		5	benzyna+L	manualna	na przedni	Polska	
6	Opel	Astra	G	kompakt	slaskie	3990		136	1998	1998		230000	czarny		5	benzyna	manualna	na przedni	Polska	
7	Mazda	Premacy		minivan	dolnoslaskie	7750		100	1998	2004		210563	sreby		5	diesel	manualna	na przedni	Niemcy	
8	Seat	Leon	II	kompakt	dolnoslaskie	19900		200	1984	2006		198000	czarny		5	benzyna	manualna	na przedni	Szwajcaria	
9	Volkswagen	Passat	B6	sedan	lodzkie	13999		105	1900	2005		202749	czarny		5	diesel	manualna	na przedni	Polska	
10	Opel	Zafira	A	minivan	lodzkie	7900		125	1800	2001		196000	srebrny		5	benzyna	manualna	4x4(stały)	Niemcy	
11	Mercedes-	Klasa A	W168	kompakt	lodzkie	6500		102	1598	2002		200000	niebieski		5	diesel	manualna	na przedni	Polska	
12	Skoda	Octavia	II	kombi	malopolskie	19500		140	1986	2008		220000	czarny		5	diesel	manualna	na przedni	Polska	
13	Seat	Toledo	II	kompakt	wielkopolskie	11300		105	1598	2003		174600	szary		4	benzyna	manualna	na przedni	Niemcy	
14	Peugeot	508		kombi	slaskie	42500		115	1560	2014		156800	biały		5	diesel	manualna	na przedni	Polska	
15	Skoda	Octavia	III	kombi	wielkopolskie	50900		105	1598	2015		91800	biały		5	diesel	manualna	na przedni	Polska	
16	Peugeot	Partner	I	kombi	lodzkie	7990		90	2000	2004		275763	złoty		5	diesel	manualna	na przedni	Polska	
17	Porsche	Cayman		coupe	wielkopolski	95900	117957	300	3400	2006		83000	srebrny		3	benzyna	automatyczna	na tylne koła	Polska	
18	Toyota	Auris	II	kompakt	mazowieckie	46000		132	1598	2014		46000	biały-metal		5	benzyna	manualna	na przedni	Polska	
19	Toyota	Auris	II	kompakt	dolnoslaskie	30300		99	1329	2014		151783	biały		5	benzyna	manualna	na przedni	Polska	
20	Mazda		3 II	kompakt	zachodniopomorskie	25200		105	1598	2009		135800	czarny		5	benzyna	manualna	na przedni	Austria	
21	Mazda		3 II	kombi	malopolskie	30600		150	1999	2009		116000	brazowy		5	benzyna	manualna	na przedni	Niemcy	
22	Mazda		6 I	kombi	dolnoslaskie	15700		146	2000	2007		190000	czarny		5	diesel	manualna	na przedni	Polska	
23	Mazda		6 I	kompakt	lodzkie	17900		147	2000	2006		238482	szary		5	benzyna+L	manualna	na przedni	Polska	
24	Mazda		6 I	kombi	kujawsko-pomorskie	11000		121	1998	2005		204000	srebrny		5	diesel	manualna	na przedni	Polska	
25	Citroen	C4	II	kompakt	wielkopolskie	36200		92	1560	2014		86561	biały		5	diesel	manualna	na przedni	Polska	
26	Citroen	C4	II	kompakt	malopolskie	36997		90	1397	2014		33000	biały		5	benzyna	manualna	na przedni	Polska	
27	Citroen	C4	II	kompakt	lodzkie	16900		90	1397	2014		35000	szary		5	benzyna	manualna	na przedni	Francja	

Rysunek 4.1: Podgląd stworzonego zbioru

Użyte zmienne

W stworzonym zbiorze danych znajduje się 29 atrybutów, opisujących 61 różnych rekordów, tj. obiektów reprezentowanych przez wiersze, którym przypisano pewne wartości atrybutów. Wśród zebranych danych można wyróżnić zarówno zmienne jakościowe, jak i ilościowe.

Zmiennymi jakościowymi są atrybuty (pisownia zgodna z plikiem danych):

- marka,
- model,
- typ,
- województwo,
- kolor,
- rodzaj.paliwa,
- skrzynia.biegow,
- napęd,
- kraj.aktualnej.rejestracji,
- kraj.pochodzenia,
- stan,
- ABS,
- uszkodzony,
- pierwszy.wlasciciel,
- kto.sprzedaje,
- serwisowany,
- komputer.pokladowy,
- ESP,
- klimatyzacja,
- bezwypadkowy,
- status.pojazdu.sprawozdanego

Zmiennymi ilościowymi są atrybuty:

- cena.netto,
- cena.brutto,
- moc,
- pojemnosc.skokowa,
- rok.produkcji,
- przebieg,
- l.drzwi.

4.2 Użyte programy

Zbiór danych został umieszczony w pliku arkusza programu Excel, z kolei implementacje zostały stworzone w języku R, który ma zastosowanie w statystyce i ekonometrii.

4.3 Implementacje wybranych metod

W sekcji tej zaprezentujemy implementacje wybranych metod porządkowania liniowego oraz nieliniowego. W części dotyczącej porządkowania liniowego porównamy również wyniki porządkowania par metod: metody rang i metody sum, metody sum i metody Hellwiga, metody rang i metody Hellwiga. Zostanie to przeprowadzone na całym zbiorze obiektów. Zaczniemy jednak od przedstawienia ogólnych funkcji odpowiedzialnych za stymulacje oraz transformacje normalizacyjne. Po to by móc następnie je wykorzystać przy prezentacji wybranych metod porządkowania nieliniowego oraz liniowego.

4.3.1 Stymulacja zmiennych

Poniżej zostaną przedstawione metody stymulacji zmiennych opisane wcześniej w sekcji 2.4. Ograniczamy się tu do przypadku, że dana zmienna jest destymulantą i należy ją poddać stymulacji. W tym celu stworzyliśmy dwie funkcje:

- `stymulacja_przekształcenie_ilorazowe(x,y)`,
- `stymulacja_przekształcenie_roznicowe(x,y)`.

W miejscu argumentu 'x' należy wpisać nazwę zbioru na którym dokonywane jest porządkowanie, z kolei w miejscu argumentu 'y' należy podać nazwę kolumny poddanej stymulacji, z tym że nazwa kolumny musi zostać podana w cudzysłowie.

```
stymulacja_przekształcenie_ilorazowe<-function(x,y)
{
  for (i in 1:nrow(x))
  {
    x[i,which(colnames(x)==y)]=1/x[i,which(colnames(x)==y)]
  }
  return(x)
}
```

```
stymulacja_przekształcenie_roznicowe<-function(x,y)
{
  max_wartosc=max(x[which(colnames(x)==y)])
  for (i in 1:nrow(x))
  {
    x[i,which(colnames(x)==y)]=max_wartosc-x[i,which(colnames(x)==y)]
  }
  return(x)
}
```

Uwaga 8. *Stymulacja zmiennych dokonywana jest pojedynczo, tj. jeżeli w naszym zbiorze jest wiele zmiennych mających charakter destymulanty, dla każdej z nich musimy użyć funkcji, na sam koniec nadpisać nasz zbiór, tym nowym wystymulowanym, dzięki czemu przekształcenia zostaną zapisane.*

4.3.2 Transformacje normalizacyjne

Poniżej zostaną przedstawione ogólne funkcje transformacji normalizacyjnej - standaryzacja, unitaryzacja, przekształcenie ilorazowe. Podobnie jak metody stymulacji zmiennych, zostały one wcześniej omówione w sekcji 2.4. Aby wywołać którąś z poniższych funkcji, należy w miejsce 'x' wpisać nazwę zbioru.

```
standaryzacja<-function(x)
{
  suma=0
  srednia=0
  odchylenie=0
  for (j in 2:ncol(x))
  {
    suma[j]=sum(x[j])
    srednia[j]=suma[j]/nrow(x)
    suma_kwadratow=0
    kwadrat=0
    for(i in 1:nrow(x))
    {
      kwadrat=(x[i,j]-srednia[j])^2
      suma_kwadratow=suma_kwadratow+kwadrat
    }
  }
}
```



```

    }
    odchylenie[j]=sqrt(suma_kwadratow/nrow(x))
    for (i in 1:nrow(x))
    {
        x[i,j]=(x[i,j]-srednia[j])/odchylenie[j]
    }
}
return(x)
}

```

```

unitaryzacja<-function(x)
{
    maksi=0
    minim=0
    for (j in 2:ncol(x))
    {
        maksi[j]=max(x[j])
        minim[j]=min(x[j])
        for (i in 1:nrow(x))
        {
            x[i,j]=(x[i,j]-minim[j])/(maksi[j]-minim[j])
        }
    }
    return(x)
}

```

```

przekształcenie_ilorazowe<-function(x)
{
    suma=0
    srednia=0
    for (j in 2:ncol(x))
    {
        suma[j]=sum(x[j])
        srednia[j]=suma[j]/nrow(x)
        for(i in 1:nrow(x))
        {
            x[i,j]=x[i,j]/srednia[j]
        }
    }
    return(x)
}

```

4.3.3 Metody porządkowania nieliniowego

Zostanie tutaj zaprezentowane zastosowanie nieliniowego porządkowania danych przy pomocy istniejących funkcji biblioteki cluster, w której to funkcja agnes umożliwia uporządkowanie zbioru po wyborze odpowiedniej metody aglomeracyjnej. Mamy tu do wyboru metody: single - metoda najbliższego sąsiedztwa, complete - metoda najdalszego sąsiedztwa, ward - metoda

Warda, average - metoda średniej między grupowej. Poniżej zostanie zaprezentowane zastosowanie metod single oraz complete wraz z porównaniem wyników porządkowania.

Import danych

Na początku należy jednak zaimportować dane, które chcemy poddać porządkowaniu. W tym celu należy zaimportować bibliotekę readxl - gdyż dane pobieramy z excela, a w kolejnym kroku wywołujemy plik, podając jego ścieżkę z rozszerzeniem.xlsx. My użyjemy tutaj zbioru wszystkich obiektów, będących ofertami sprzedaży aut.

```
library(readxl)
zbior_danych <- read_excel("datasets/zbior_danych.xlsx",
                           sheet = "DANE_INNA_WERSJA")
```

Podgląd danych:

```
head(zbior_danych)
```

```
## # A tibble: 6 x 30
##   Nr MARKA  MODEL  WERSJA TYP      WOJEWODZTWO  'CENA.NETTO_[pln]'
##   <dbl> <chr>   <chr>   <chr> <chr>   <chr>          <dbl>
## 1  1.00 Hyundai i20     II     kompakt malopolskie      NA
## 2  2.00 Hyundai i20     I      kompakt mazowieckie      NA
## 3  3.00 Subaru  Legacy V      kombi    mazowieckie      NA
## 4  4.00 Ford    Mondeo Mk4     sedan    dolnoslaskie      NA
## 5  5.00 Opel    Astra  G      kompakt slaskie        NA
## 6  6.00 Mazda   Premacy <NA> minivan dolnoslaskie      NA
## # ... with 23 more variables: 'CENA.BRUTTO_[pln]' <dbl>, 'MOC_[km]' <dbl>,
## # 'POJEMNOSC.SKOKOWA_[cm3]' <dbl>, ROK.PRODUKCJI <dbl>, 'PRZEBIEG_[km]'
## # <dbl>, KOLOR <chr>, L.DZRZWI <dbl>, RODZAJ.PALIWA <chr>,
## # SKRZYNIA.BIEGOW <chr>, NAPED <chr>, KRAJ.AKTUALNEJ.REJESTRACJI <chr>,
## # KRAJ.POCHODZENIA <chr>, STATUS.POJAZDU.SPROWADZONEGO <chr>,
## # PIERWSZY.WLASCICIEL <dbl>, KTO.SPRZEDAJE <chr>, STAN <chr>,
## # SERWISOWANY <dbl>, ABS <dbl>, KOMPUTER.POKLADOWY <dbl>, ESP <dbl>,
## # KLIMATYZAJCA <dbl>, BEZWYPADKOWY <dbl>, USZKODZONY <dbl>
```

Podzbiór danych

W kolejnym, kroku po przyjrzeniu się zbiorowi danych, użytkownik musi zdecydować na których danych ilościowych chce pracować - ważna jest znajomość danych. Dodatkowo pierwszą kolumną musi być kolumna zawierająca numery indeksów obiektów, ze względu na to, że jako wynik zastosowania funkcji odpowiedzialnej za porządkowanie, zostaną zwrócone w kolejności malejącej numery indeksów, mówiące o kolejności uporządkowania. W związku z tym, za pomocą poniższej procedury użytkownik tworzy podzbiór zaimportowanego zbioru, gdzie w cudzysłowie wpisuje nazwy kolumn zawierających zmienne ilościowe, wybrane do porządkowania (zakładamy, że podzbiór będzie nazywał się dane_porzadkowanie)). Wybranymi kolumnami są: cena, moc, pojemność, rok produkcji, przebieg.

```
dane_porzadkowanie<-zbior_danych[c("Nr", "CENA.BRUTTO_[pln]", "MOC_[km]",
                                     "POJEMNOSC.SKOKOWA_[cm3]",
                                     "ROK.PRODUKCJI", "PRZEBIEG_[km]")]
```

Transformacje danych

Przed samym porządkowaniem, wymagane jest aby zmienne miały charakter stymulant oraz by zostały poddane transformacji normalizacyjnej. Aby funkcja dokonująca porządkowania dawała poprawny wynik, użytkownik musi zająć się transformacją przed jej zastosowaniem. Poniżej podaliśmy tego przykład. Dla zmiennych, które stymulantami nie są, należy dokonać stymulacji. Wśród zmiennych poddanych porządkowaniu, do stymulant nie należy zmienna zmienna: przebieg - jest destymulantą, w związku z tym, została przekształcona na stymulantę, za pomocą przekształcenia ilorazowego.

```
dane_porzadkowanie<-stymulacja_przekształcenie_ilorazowe(dane_porzadkowanie,
"PRZEBIEG_[km] ")
```

W celu uzyskania porównywalności między zmiennymi, zostały one poddane transformacji normalizacyjnej - unitaryzacji.

```
dane_porzadkowanie<-unitaryzacja(dane_porzadkowanie)
```

Zastosowanie metod aglomeracyjnych

Chcąc zastosować funkcję `agnes`, należy w pierwszej kolejności wyznaczyć macierz odległości pomiędzy wszystkimi parami obiektów. Do wyznaczenia odległości zostanie użyta metryka euklidesowa - w tym celu zostanie wykorzystana funkcja `dist(x, method="")` - w miejsce 'x' należy wpisać nazwę tabeli zawierających dane do uporządkowania, a w nawiasie[,] na miejscu drugiej współrzędnej należy podać wektor kolumn, na podstawie którego wartości zostanie wyznaczona macierz odległości. W miejsce argumentu `method` należy wpisać nazwę metryki na podstawie której zostanie obliczona odległość - w naszym przypadku będzie to `euclidean` - euklidesowa.

```
odleglosci <- dist(dane_porzadkowanie[,c("CENA.BRUTTO_[p1n]","MOC_[km] ",
"POJEMNOSC.SKOKOWA_[cm3]","ROK.PRODUKCJI","PRZEBIEG_[km] ")] ,
```

Następnie należy zaimportować bibliotekę `cluster`, by móc skorzystać z metod aglomeracyjnych. Jak już zostało wspomniane we wstępie, wywołanie metod aglomeracyjnych odbywa się dzięki funkcji `agnes(x, method="")`. W miejsce argumentu 'x' - zostanie podana wyznaczona macierz odległości z kolei, w kolejnym argumencie - `method` zostanie podana reguła wyznaczania odległości pomiędzy nową grupą a pozostałymi obiektów. Regułą tą może być metoda najbliższego sąsiedztwa, najdalszego, Warda lub średniej między grupowej. My wykorzystamy metodę najbliższego sąsiedztwa oraz najdalszego.

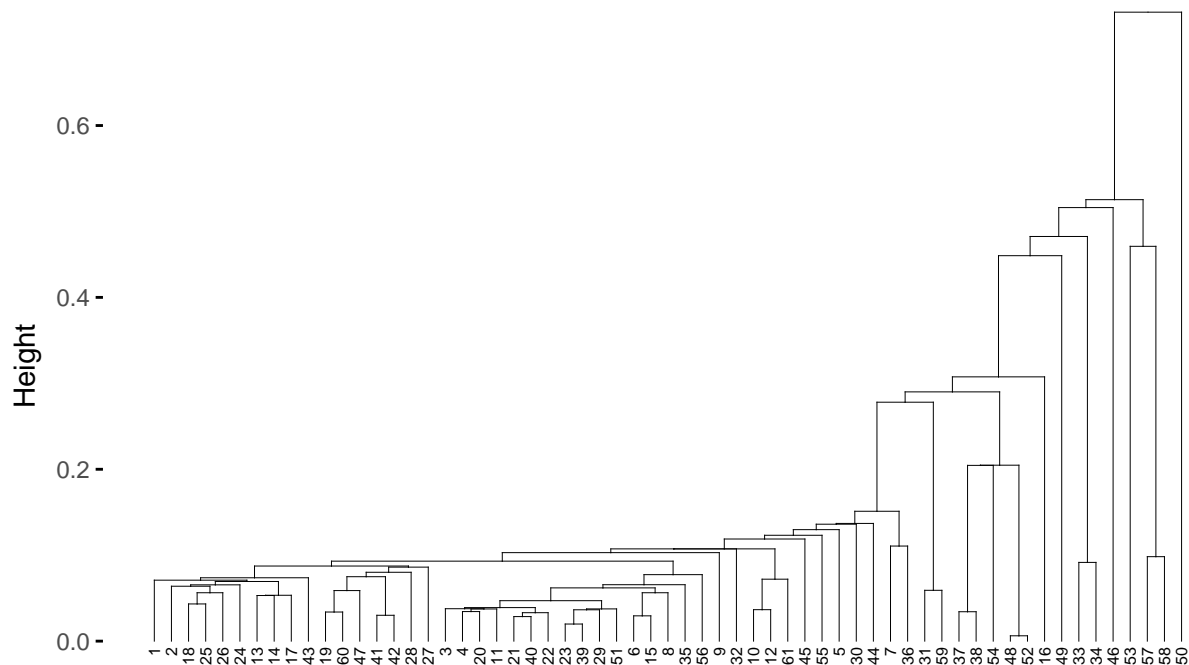
```
library(cluster)
metoda_najblizszego <- agnes(odleglosci, method = "single")
metoda_najdalszego<- agnes(odleglosci, method = "complete")
```

Ostatni etap to graficzne zaprezentowanie wyniku w postaci dendrogramu. W tym celu należy zaimportować bibliotekę `factoextra`, w której to jest funkcja `fviz_dend(x, main = "")`. W miejsce argumentu `x` należy wpisać nazwę obiektu powstałego przy pomocy funkcji `agnes`, `main` to tytuł wykresu. Dodatkowo na osi pionowej zaprezentowane są odległości między obiektami, z kolei na osi poziomej znajdują się numery indeksów obiektów.

Dendrogram dla metody najbliższego sąsiada, prezentuje się następująco:

```
library(factoextra)
fviz_dend(metoda_najblizszego, lwd=0.1, cex=0.45,
main = "Metoda najbliŹszego sąsiedztwa")
```

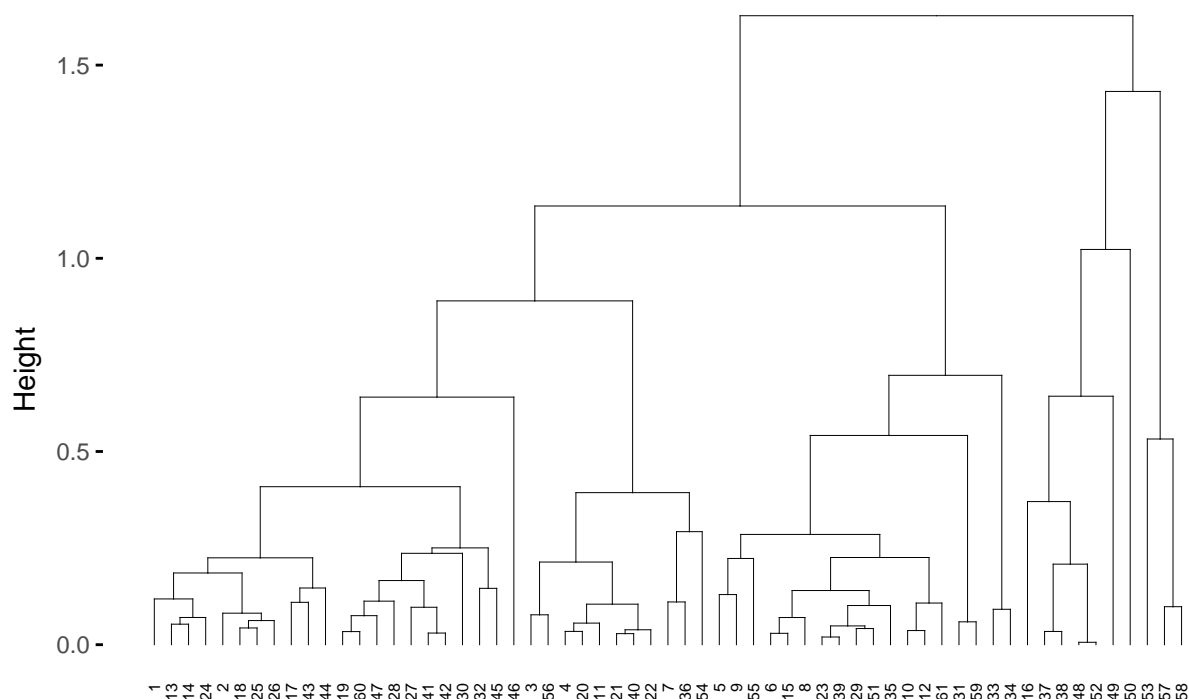
Metoda najbliŹszego sasiedztwa



Dendrogram dla metody najdalszego sąsiada, wygląda natomiast:

```
fviz_dend(metoda_najdalszego, lwd=0.1, cex=0.45,
main = "Metoda najdalszego sąsiedztwa")
```

Metoda najdalszego sasiedztwa



Porównanie wyników uporządkowania

W przypadku dendrogramu uzyskanego metodą najdalszego sąsiedztwa można zauważyć większą odległość wiązań niż dla metody najbliższego sąsiedztwa, które to obrazują odległość między grupami. Dodatkowo obiekty charakteryzujące się najbardziej korzystnym wartością zmiennych, zostały pogrupowane w jedną grupę (tu mowa o obiektach o numerach indeksów 49, 50, 53, 57, 58), natomiast w przypadku metody najbliższego sąsiedztwa obiekty te zostały rozdzielone na pojedyncze grupy. Można również zauważyć, że w przypadku uporządkowania metodą najdalszego sąsiedztwa obiekty tworzą pewnego rodzaju podgrupy - skupiska, z kolei w przypadku uporządkowania metodą najbliższego sąsiada, wynik porządkowania można porównać do wyglądu lawiny, lub góry, u której podnóża znajdują się obiekty o najniższych wartościach opisywanych je zmiennych, które to następnie łączą się i dochodzą do szczytu, który stanowią obiekty o najwyższych wartościach opisywanych je zmiennych.

4.3.4 Metody porządkowania liniowego

Zostaną tu zaprezentowane metody liniowe: metoda sum, rang oraz Hellwiga. Dla wszystkich metod będziemy korzystać z tego samego zbioru danych, który jest podzbiorem wyjściowego zbioru - zawiera 8 najbardziej różniących się obiektów. Różnice dla każdej z metod będą zaczynały się od sekcji związanej z transformacją danych.

Import danych

Wobec powyższego, podobnie jak przy metodach porządkowania nieliniowego, użytkownik musi zaimportować dane, które chce poddać porządkowaniu.

```
library(readxl)
zbior_danych <- read_excel("datasets/8_Rozniacych_sie_obiektow.xlsx",
  sheet = "Arkusz1", col_types = c("numeric","text",
    "text", "text", "text", "text",
    "numeric", "numeric", "text", "numeric",
    "text", "text", "text", "text","text",
    "text", "numeric", "text", "text",
    "text", "numeric", "numeric",
    "numeric", "numeric", "numeric"))
```

Podgląd danych:

```
head(zbior_danych)
```

```
## # A tibble: 6 x 29
##      Nr MARKA  MODEL WERSJA TYP  WOJEWODZTWO 'CENA.BRUTTO_[p~ 'MOC_[km] '
##    <dbl> <chr>  <chr>  <chr> <chr> <chr>          <dbl>      <dbl>
## 1  1.00 Mazda  3      II    komp~ zachodniop~    25200      105
## 2  2.00 Jaguar XF      X260  kombi dolnoslask~    323600      240
## 3  3.00 Subaru B9 Tr~ <NA>  suv  malopolskie    38900      245
## 4  4.00 Volks~ Golf  VII    kombi lodzkie    113900      150
## 5  5.00 Peuge~ 508    <NA>  kombi slaskie    42500      115
## 6  6.00 Opel  Antara <NA>  suv  lodzkie    24000      150
## # ... with 21 more variables: 'POJEMNOSC.SKOKOWA_[cm3]' <dbl>,
## #   ROK.PRODUKCJI <dbl>, 'PRZEBIEG_[km]' <dbl>, KOLOR <chr>, L.DZRZWI
## #   <dbl>, RODZAJ.PALIWA <chr>, SKRZYNIA.BIEGOW <chr>, NAPED <chr>,
## #   KRAJ.AKTUALNEJ.REJESTRACJI <chr>, KRAJ.POCHODZENIA <chr>,
## #   STATUS.POJAZDU.SPROWADZONEGO <chr>, PIERWSZY.WLASCICIEL <dbl>,
## #   KTO.SPRZEDAJE <chr>, STAN <chr>, SERWISOWANY <chr>, ABS <dbl>,
## #   KOMPUTER.POKLADOWY <dbl>, ESP <dbl>, KLIMATYZAJCA <dbl>, BEZWYPADKOWY
## #   <dbl>, USZKODZONY <dbl>
```

Podzbiór danych

W kolejnym, kroku po przyjrzeniu się zbiorowi danych, użytkownik musi zdecydować na których danych ilościowych chce pracować - ważna jest znajomość danych. Dodatkowo pierwszą kolumną musi być kolumna zawierająca numery indeksów obiektów, ze względu na to, że w wyniku zastosowania funkcji odpowiedzialnej za porządkowanie, zostaną zwrócone w kolejności malejącej numery indeksów, mówiące o kolejności uporządkowania. W związku z tym, za pomocą poniższej procedury użytkownik tworzy podzbiór zaimportowanego zbioru, gdzie w cudzysłowie wpisuje nazwy kolumn zawierających zmienne ilościowe, wybrane do porządkowania (przyjmijmy założenie, że podzbiór będzie nazywał się dane_porzadkowanie - będzie to pomocne w dalszej części programu). U nas wybranymi kolumnami są: cena, moc, pojemność, rok produkcji, przebieg.

```
dane_porzadkowanie<-zbior_danych[c("Nr","CENA.BRUTTO_[p1n]","MOC_[km] ",
  "POJEMNOSC.SKOKOWA_[cm3] ",
  "ROK.PRODUKCJI","PRZEBIEG_[km] ")]
```

Transformacje danych

Przed zastosowaniem metod, należy dokonać ich transformacji. W pierwszym kroku należy dokonać stymulacji destymulant - dla metody rang i metody Hellwiga zostanie użyte przekształcenie ilorazowe, dla metody sum przekształcenie ilorazowe. Po stymulacji można przejść do transformacji normalizacyjnej - dla metod: sum i rang użyjemy normalizacji, natomiast dla metody Hellwiga standaryzacji.

4.3.5 Metoda sum

Stymulacja

```
dane_porzadkowanie<-stymulacja_przekształcenie_roznicowe(dane_porzadkowanie,
"PRZEBIEG_[km] ")
```

Transformacja normalizacyjna

```
dane_porzadkowanie<-unitaryzacja(dane_porzadkowanie)
```

Funkcja porządkująca

Chcąc przeprowadzić porządkowanie znormalizowanych danych za pomocą metody sum, należy użyć poniższej funkcji:

```
funkcja_porzadkowanie_metoda_sum<-function(x)
{
  x[,"zmienna_syntetyczna"] <-0 #ostatnia kolumna to zmienna_syntetyczna
  for(i in 1:nrow(x))
  {
    for(j in 2:(ncol(x)-1))
    {
      x[i,ncol(x)]=x[i,ncol(x)]+x[i,j]
    } #liczba kolumn - (kolumna_nr_indeksu+kolumna_zmienna_synt)
    x[i,ncol(x)]=x[i,ncol(x)]/(ncol(x)-2)
  }
  #wyeliminowanie ujemnych wartosci zmiennej syntetycznej
  min_zmienna=min(x$zmienna_syntetyczna)
  for(i in 1:nrow(x))
  {
    x[i,ncol(x)]=x[i,ncol(x)]-min_zmienna
  }
  #normalizacja zm. syntetycznej
  max_zmienna=max(x$zmienna_syntetyczna)
  for(i in 1:nrow(x))
  {
    x[i,ncol(x)]=x[i,ncol(x)]/max_zmienna
  }
  x<-x[order(-x$zmienna_syntetyczna),]
  return(x[1])
}
```

Wywołanie funkcji dla zbioru dane_porzadkowanie

```
funkcja_porzadkowanie_metoda_sum(dane_porzadkowanie)
```

```
## # A tibble: 8 x 1
##   Nr
##   <dbl>
## 1  2.00
## 2  4.00
## 3  3.00
## 4  5.00
## 5  6.00
## 6  1.00
## 7  8.00
## 8  7.00
```

Funkcja zwraca nam indeksy uporządkowanych obiektów, tj. 1 miejsce zajął obiekt z numerem indeksu 2, 2 miejsce obiekt z numerem indeksu równym 4, z kolei miejsce ostatnie zajął obiekt o numerze indeksu równym 7.

4.3.6 Metoda rang

Stymulacja

```
dane_porzadkowanie<-stymulacja_przekształcenie_ilorazowe(dane_porzadkowanie,
"PRZEBIEG_[km]")
```

Transformacja normalizacyjna

```
dane_porzadkowanie<-unitaryzacja(dane_porzadkowanie)
```

Funkcja porządkująca

Chcąc przeprowadzić porządkowanie znormalizowanych danych za pomocą metody rang, należy użyć poniższej funkcji:

```
funkcja_porzadkowanie_metoda_rang<-function(x)
{
  y<-x
  for (i in 2:ncol(x))
  {
    x[ncol(x)+1]=rank(-x[i])
  }
  x[, "zmienna_syntetyczna"] <-0 #ostatnia kolumna to zmienna_syntetyczna
  for(i in 1:nrow(x))
  {
    for(j in (ncol(y)+1):(ncol(x)-1))
    {
      x[i,ncol(x)]=x[i,ncol(x)]+x[i,j]
    }
  }
}
```



```

        j=j+1
    }
    x[i,ncol(x)]=x[i,ncol(x)]/(ncol(x)-7)
}
x<-x[order(x$zmienna_syntetyczna),]
print("Numery indeksów obiektów po uporządkowaniu: ")
return(x[1])
}

```

Wywołanie funkcji dla zbioru dane_porzadkowanie

```
funkcja_porzadkowanie_metoda_rang(dane_porzadkowanie)
```

```
## [1] "Numery indeksów obiektów po uporządkowaniu: "
```

```
## # A tibble: 8 x 1
```

```
##      Nr
```

```
##   <dbl>
```

```
## 1  2.00
```

```
## 2  4.00
```

```
## 3  3.00
```

```
## 4  5.00
```

```
## 5  6.00
```

```
## 6  1.00
```

```
## 7  8.00
```

```
## 8  7.00
```

Na podstawie powyższego wyniku, widać, że 1 miejsce zajęła oferta sprzedaży z numerem indeksu 2, 2 miejsce oferta o numerze indeksu równym 4, a ostatnie oferta z numerem indeksu numer 7. Porównując wyniki porządkowania tej metody, z metodą sum zauważamy, że wyniki są takie same.

4.3.7 Metoda Hellwiga

Stymulacja

```
dane_porzadkowanie<-stymulacja_przekształcenie_ilorazowe(dane_porzadkowanie,
"PRZEBIEG_[km] ")
```

Transformacja normalizacyjna

```
dane_porzadkowanie<-standaryzacja(dane_porzadkowanie)
```

Funkcja porządkująca

Chcąc przeprowadzić porządkowanie znormalizowanych danych za pomocą metody Hellwiga, należy użyć poniższej funkcji:

```

funkcja_porzadkowanie_metoda_Hellwiga<-function(x)
{
  obiekt_wz=0
  for (j in 2:ncol(x))
  {
    obiekt_wz[j]=max(x[j])
  }
  odleg<- x[c("Nr" )]
  for (i in 1:nrow(x))
  {
    SUMKA=0
    for (j in 2:ncol(x))
    {
      SUMKA=SUMKA+(x[i,j]-obiekt_wz[j])^2
    }
    odleg[i,2]=sqrt(SUMKA) #kolumna zawierajaca odleglosci
  }
  d_0=0
  suma=0
  srednia=0
  odchylenie=0
  for (j in 2:ncol(odleg))
  {
    suma[j]=sum(odleg[j])
    srednia[j]=suma[j]/nrow(odleg)
    suma_kwadratow=0
    kwadrat=0
    for(i in 1:nrow(odleg))
    {
      kwadrat=(odleg[i,j]-srednia[j])^2
      suma_kwadratow=suma_kwadratow+kwadrat
    }

    odchylenie[j]=sqrt(suma_kwadratow/nrow(odleg))
    d_0=srednia[j]+2*odchylenie[j] #d_0 to po prostu wartosc
  }
#ostatnia kolumna to jak zawsze zmienna syntetyczna
  x[,"zmienna_syntetyczna"] <-0
  for (i in 1:nrow(x))
  {
    x[i,ncol(x)]=1-(odleg[i,2]/d_0)
  }
  x<-x[order(-x$zmienna_syntetyczna),]
  return(x[1])
}

```

Wywołanie funkcji dla zbioru dane_porzadkowanie

```
metoda_Hellwiga(dane_porzadkowanie)
```

```
## # A tibble: 8 x 1
##   Nr
##   <dbl>
## 1  2.00
## 2  4.00
## 3  3.00
## 4  6.00
## 5  5.00
## 6  1.00
## 7  8.00
## 8  7.00
```

Na podstawie powyższego wyniku, widać, że 1 miejsce zajęła oferta sprzedaży z numerem indeksu 2, a ostatnie oferta z numerem indeksu o numerze 7. Wobec powyższego wyniki porządkowania tą metodą są inne, niż wyniki metody sum.

4.3.8 Porównanie wyników metod porządkowania liniowego dla całego zbioru

Analizując wyniki porządkowania zaprezentowane w poprzedniej sekcji, zauważamy że dla tak małego zbioru uzyskaliśmy zgodność rezultatów dla metody sum oraz rang. Sprawdźmy teraz jak będzie w przypadku zastosowania tych metod dla całego zbioru. W związku z tym należy zaimportować cały zbiór danych. W dalszej kolejności, należy stworzyć trzy kopie zbioru, dla każdej z metod porządkowania. Następnie przeprowadzamy o odpowiednie transformacje o raz stosujemy algorytmy porządkowania.

Import danych

Tak samo jak przy implementacji metod porządkowania nieliniowego, importujemy pełen zbiór obiektów.

```
library(readxl)
zbior_danych <- read_excel("datasets/zbior_danych.xlsx",
                           sheet = "DANE_INNA_WERSJA")
```

Podzbiór danych

Mając zaimportowany zbiór danych, tworzymy podzbiór na którym będziemy pracować.

```
dane_porzadkowanie<-zbior_danych[c("Nr", "CENA.BRUTTO_[p1n]", "MOC_[km]",
                                     "POJEMNOSC.SKOKOWA_[cm3]",
                                     "ROK.PRODUKCJI", "PRZEBIEG_[km]")]
```

Transformacja danych

Ze względu na to, że metody: sum, rang, Hellwiga bazują na różnych przekształceniach, tworzymy 3 podzbiory, które poddamy wcześniej już zaprezenowanym metodom stymulacji.

```
dane_rang<-stymulacja_przekształcenie_ilorazowe(dane_porzadkowanie,"PRZEBIEG_[km]")
dane_hellwig<-stymulacja_przekształcenie_ilorazowe(dane_porzadkowanie,"PRZEBIEG_[km]")
dane_sum<-stymulacja_przekształcenie_roznicowe(dane_porzadkowanie,"PRZEBIEG_[km]")
```

Teraz dokonamy przekształcenia normalizacyjnego wystymulowanych podzbiorów.

```
dane_sum<-unitaryzacja(dane_sum)
dane_rang<-unitaryzacja(dane_rang)
dane_hellwig<-standaryzacja(dane_hellwig)
```

4.3.9 Zastosowanie funkcji odpowiedzialnych za porządkowania

Wykorzystamy tutaj wcześniej już wyjaśnione funkcje odpowiedzialne za porządkowanie.

```
dane_sum<-funkcja_porzadkowanie_metoda_sum(dane_sum)
dane_rang<-funkcja_porzadkowanie_metoda_rang(dane_rang)
dane_hellwig<-funkcja_porzadkowanie_metoda_Hellwiga(dane_hellwig)
```

4.3.10 Porównanie wyników

W celu porównania wyników uporządkowania zostały stworzone trzy tabele pomocnicze. Każda tabela zawiera trzy kolumny, w dwóch pierwszych kolumnach znajdują się indeksy obiektów po uporządkowaniu. Z kolei w kolumnie trzeciej - „porownanie” znajduje się jedna z dwóch wartości: 0 lub 1. Odpowiednio wartość 1 jest przypisywana tym rekordom, dla których zgadza się kolejność uporządkowania przy zastosowaniu dwóch różnych metod porządkowania. Dodatkowo zostały również stworzone trzy tabele o nazwie „podsumowanie”. W nich zliczane są wystąpienia wartości: 0 oraz 1 w kolumnie poprzedniej tabeli - „porownanie”.

Para 1: metoda rang i metoda sum

```
tabela_porownawcza=data.frame(dane_rang,dane_sum)
names(tabela_porownawcza)<-c("dane_rang","dane_sum")
tabela_porownawcza$porownanie=0
#inwersja
for(i in 1:nrow(tabela_porownawcza))
{
  if(tabela_porownawcza$dane_sum[i]==tabela_porownawcza$dane_rang[i])
  {
    tabela_porownawcza$porownanie[i]=1
  }
}
head(tabela_porownawcza,15)
```

```
##      dane_rang dane_sum porownanie
## 1          49      49          1
## 2          50      50          1
## 3          53      53          1
## 4          48      16          0
## 5          52      48          0
```

```
## 6      16      52      0
## 7      57      38      0
## 8      38      57      0
## 9      37      37      1
## 10     58      58      1
## 11     44      44      1
## 12     17      46      0
## 13     56      17      0
## 14     20       1      0
## 15     46      41      0
```

#podsumowanie

```
podsumowanie=as.data.frame(table(tabela_porownawcza$porownanie))
names(podsumowanie)<-c("wartość","ilosc wystąpień")
podsumowanie
```

```
##      wartość ilosc wystąpień
## 1          0          51
## 2          1          10
```

Para 2: metoda rang i metoda Hellwiga

```
tabela_porownawcza=data.frame(dane_rang,dane_hellwig)
names(tabela_porownawcza)<-c("dane_rang","dane_hellwig")
tabela_porownawcza$porownanie=0
#inwersja
for(i in 1:nrow(tabela_porownawcza))
{
  if(tabela_porownawcza$dane_hellwig[i]==tabela_porownawcza$dane_rang[i])
  {
    tabela_porownawcza$porownanie[i]=1
  }
}
head(tabela_porownawcza,15)
```

```
##      dane_rang dane_hellwig porownanie
## 1          49          53          0
## 2          50          49          0
## 3          53          50          0
## 4          48          16          0
## 5          52          57          0
## 6          16          48          0
## 7          57          52          0
## 8          38          58          0
## 9          37          38          0
## 10         58          37          0
## 11         44          46          0
## 12         17          54          0
## 13         56          56          1
## 14         20          36          0
## 15         46           7          0
```

```
#podsumowanie
podsumowanie=as.data.frame(table(tabela_porownawcza$porownanie))
names(podsumowanie)<-c("wartość","ilosc wystąpień")
podsumowanie
```

```
##    wartość ilosc wystąpień
## 1         0          56
## 2         1           5
```

Para 3: metoda sum i metoda Hellwiga

```
tabela_porownawcza=data.frame(dane_sum,dane_hellwig)
names(tabela_porownawcza)<-c("dane_sum","dane_hellwig")
tabela_porownawcza$porownanie=0
#inwersja
for(i in 1:nrow(tabela_porownawcza))
{
  if(tabela_porownawcza$dane_hellwig[i]==tabela_porownawcza$dane_sum[i])
  {
    tabela_porownawcza$porownanie[i]=1
  }
}
head(tabela_porownawcza,15)
```

```
##    dane_sum dane_hellwig porownanie
## 1         49          53           0
## 2         50          49           0
## 3         53          50           0
## 4         16          16           1
## 5         48          57           0
## 6         52          48           0
## 7         38          52           0
## 8         57          58           0
## 9         37          38           0
## 10        58          37           0
## 11        44          46           0
## 12        46          54           0
## 13        17          56           0
## 14         1          36           0
## 15        41           7           0
```

```
#podsumowanie
podsumowanie=as.data.frame(table(tabela_porownawcza$porownanie))
names(podsumowanie)<-c("wartość","ilosc wystąpień")
podsumowanie
```

```
##    wartość ilosc wystąpień
## 1         0          59
## 2         1           2
```

4.3.11 Podsumowanie

Na podstawie powyższych wyników zauważamy, że najwięcej zgodności wyniku porządkowania jest widoczne dla pary pierwszej - metody rang oraz metody sum. W przypadku kolejnych par, zgodność ta jest już niewielka.

Rozdział 5

Podsumowanie

W niniejszej pracy rozważone zostały wybrane zastosowania statystycznych metod porządkowania danych wielowymiarowych. Celem było porównanie tych metod praktycznej statystyki z matematyczną teorią porządków, jak również zaprezentowanie ich zastosowania w praktycznym przykładzie. W pracy dostrzeżono wiele analogii pomiędzy porządkami liniowymi i częściowymi, a rozważanymi metodami porządkowania liniowego oraz nieliniowego. Ponadto, w efekcie przeprowadzonych eksperymentów, zaobserwowano kilka własności wybranych metod porządkowania. Praca zawiera również wszystkie kody źródłowe przeprowadzonych eksperymentów, które w naszej ocenie mogą zostać ponownie wykorzystane w innych zastosowaniach. Podstawowe metody porządkowania danych wielowymiarowych, przypuszczają możliwość występowania w nich porządków liniowych. Przypuszczenie to jest co prawda wysoce niepewne, gdyż dla danych wielowymiarowych, bardziej spodziewane są porządki częściowe. Stosowanie zatem tych algorytmów, jeśli wyniki nie są wyraźnie jednoznaczne, może się okazać wątpliwym. Na korzyść tych algorytmów przemawia natomiast wysoka zgodność z matematyczną teorią takich porządków. W pracy wykazaliśmy, że niesłusznym jest definiowanie porządków liniowych za pomocą jedynie 3 poszukiwanych własności. Jednak obserwujemy, że algorytmy te pracując na danych, które posiadają naturalne uporządkowanie liniowe, istotnie zwracają zgodne i prawdziwe uporządkowania. W dalszej części rozpoznawaliśmy pozostałe algorytmy porządkowania danych, generujących inne matematyczne struktury. Aby w pełni zrozumieć zależności pomiędzy porządkami częściowymi, a algorytmami wyznaczającymi porządki nieliniowe w danych, niezbędna okazała się podstawowa wiedza z zakresu teorii grafów. Matematyczna teoria porządków może być bowiem z łatwością przedstawiana na grafach nazywanych diagramami Hassego. Z drugiej strony algorytmy statystyki okazują się poszukiwać w skończonych zbiorach danych właśnie struktur będących (lub przypominających) grafy. Reprezentowanie porządków na tych strukturach okazuje się mieć dodatkowe atuty w postaci przejrzystej wizualizacji uzyskiwanych wyników. Występują tu jednak pewne istotne różnice. obserwujemy, iż nie każdy porządek częściowy, mógłby za pomocą tychże algorytmów zostać wykryty. Algorytmy te bowiem generują grafy nazywane dendrytami, które nie wyczerpują rodziny wszystkich porządków częściowych. W rozdziale 4 zaprezentowaliśmy implementacje wybranych metod porządkowania liniowego, tj. wybraliśmy 3 różne metody: metodę rang, sum oraz Hellwiga. Metoda Hellwiga zaliczana jest do metod wzorcowych, czyli zakłada istnienie obiektu wzorcowego. Jako, że zmienne opisujące obiekty (zgromadzone oferty sprzedaży aut), są w większości stymulantami, stąd też przyjęliśmy koncepcję, że wartości współrzędnych obiektu wzorcowego przyjmą wartość maksymalną dla każdej zmiennej. Podkreślamy to gdyż zaobserwowaliśmy różnice uzyskiwane przez różne grupy metod. Metody wzorcowe generowały odmienne porządki, w porównaniu z metodą sum czy też metodą rang. Wyniki, o których mówimy, znajdują się w sekcji 4.3.10. Jednakże przyglądając się wynikom porządkowania zawartym w sekcji 4.3.4, widać iż dla tych trzech metod otrzymujemy wyniki wspólne dla obiektów skrajnych, znajdujących się najwyżej lub najniżej w tym

porządkowaniu. Różnice zauważamy głównie dla obiektów ze środka. Rozważone zostały też algorytmy nieliniowe, w szczególności dla metod aglomeracyjnych zaobserwowano, że istotna jest różnica w porządkach przy zastosowaniu odległości liczonej wzorem najbliższego oraz najdalego sąsiada. Umożliwiły nam one zobrazowanie wewnętrznych różnic pomiędzy tymi wersjami algorytmu porządkowania danych. Stąd też można zauważyć, że podobnie jak przy wynikach porządkowania liniowego dla małego zbioru, zaobserwować można zgodność uporządkowania tych metod dla obiektów skrajnych, znajdujących się najwyżej lub najniżej w tym uporządkowaniu. Wybrane wbudowane funkcje, obok samodzielnie stworzonych funkcji porządkowania przedstawionych w sekcji 4.3.4, mogą zostać użyte do innych zbiorów danych, ze względu na ich ogólność i uniwersalność.

Bibliografia

- [1] Grzegorz Banaszak, Wojciech Gajda. *Elementy algebry liniowej (część 1)*. Wydawnictwo Naukowo-Techniczne, Warszawa, 2002.
- [2] Jarosław Bartoszewicz. *Wykłady ze statystyki matematycznej*. Państwowe Wydawnictwo Naukowe, Warszawa, 1996.
- [3] Aleksander Błaszczyk, Sławomir Turek. *Teoria mnogości*. Państwowe Wydawnictwo Naukowe, Warszawa, 2007.
- [4] Patric Billingsley. *Prawdopodobieństwo i miara*. Państwowe Wydawnictwo Naukowe, Warszawa, 1987.
- [5] Jacek Jakubowski, Rafał Sztencel. *Wstęp do teorii prawdopodobieństwa*. SCRIPT, Warszawa, 2004.
- [6] Włodzimierz Kryszicki, Jerzy Bartos, Wacław Dyczka, Krystyna Królikowska, Mariusz Wasilewski. *Rachunek prawdopodobieństwa i statystyka matematyczna w zadaniach: część I. rachunek prawdopodobieństwa*. Państwowe Wydawnictwo Naukowe, Warszawa, 1999.
- [7] Kazimierz Kuratowski. *Wstęp do teorii mnogości i topologii*. Państwowe Wydawnictwo Naukowe, Warszawa, 2004.
- [8] Andrzej Młodak. *Analiza taksonomiczna w statystyce regionalnej*. Centrum Doradztwa i Informacji Difin, Warszawa, 2006.
- [9] Tomasz Panek, Jan Karol Zwierchowski. *Statystyczne metody wielowymiarowej analizy porównawczej: teoria i zastosowania*. Oficyna Wydawnicza, Szkoła Główna Handlowa, Warszawa, 2013.
- [10] Ryszard Rudnicki. *Wykłady z analizy matematycznej*. Państwowe Wydawnictwo Naukowe, Warszawa, 2006.
- [11] Arkadiusz Sołtysiak, Piotr Jaskulski. <http://www.antropologia.uw.edu.pl/MaCzek/maczek.html>. dostęp: 02.07.2018.
- [12] Robin J. Wilson. *Wprowadzenie do teorii grafów*. Państwowe Wydawnictwo Naukowe, Warszawa, 2008.