

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, xxx 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 | Podstawowe pojęcia teorii grafów | 10 |
| 2.5 | Wybrane pojęcia z teorii mnogości i topologii: relacja porządkująca, moc zbiorów, równoliczność, przestrzenie metryczne, miary odległości | 12 |
| 2.5.1 | Relacja porządkująca | 12 |
| 2.5.2 | Przestrzenie metryczne, miary odległości | 14 |
| 3 | Metody porządkowania | 15 |
| 3.1 | Metody porządkowania liniowego | 15 |
| 3.1.1 | Metody diagramowe | 20 |
| 3.1.2 | Metody oparte na zmiennych syntetycznych | 21 |
| 3.1.3 | Metody iteracyjne | 24 |
| 3.1.4 | Metody gradientowe | 25 |
| 3.2 | Metody porządkowania nieliniowego | 26 |
| 3.2.1 | Metody dendrytowe | 26 |
| 3.2.2 | Metody aglomeracyjne | 27 |
| 4 | Szczegółowe opisy wybranych metod porządkowania | 30 |
| 4.1 | Metody porządkowania liniowego | 30 |
| 4.1.1 | Metody iteracyjne | 33 |
| 4.1.2 | Metody gradientowe | 34 |
| 4.2 | Metody porządkowania nieliniowego | 35 |
| 4.2.1 | Metody dendrytowe | 36 |
| 4.2.2 | Metody aglomeracyjne | 37 |
| 5 | Zbiór danych | 40 |
| 5.1 | Opis zbioru | 40 |

Rozdział 1

Wstęp

Rozdział 2

Preliminaria

2.1 Notacja

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

- \mathbb{R} - zbiór liczb rzeczywistych
- \mathbb{N} - zbiór liczb naturalnych
- $O = \{O_1, O_2, \dots, O_n\}$ - zbiór obiektów przestrzennych, $n \in \mathbb{N}$
- $X = \{X_1, X_2, \dots, X_n\}$ - zbiór zmiennych (cech), gdzie $n \in \mathbb{N}$
- Ω - przestrzeń zdarzeń elementarnych
- ω - zdarzenie elementarne
- \mathcal{F} - rodzina podzbiorów zbioru Ω
- σ - sigma ciało zbiorów \mathcal{F}
- \mathcal{B} - rodzina zbiorów borelowskich
- \mathfrak{B} -rodzina wszystkich zbiorów otwartych
- $\text{med}(X_j)$ - mediana cechy X_j
- ρ - relacja porządkująca
- G - graf prosty
- $V(G)$ - zbiór wierzchołków grafu G
- $E(G)$ - krawędzie grafu G

2.2 Słownik użytych pojęć

W pracy zostały wykorzystane następujące pojęcia, których wytłumaczenie znajduje się poniżej.

- Statystyka matematyczna [5, w oparciu o rozdział 1]
Statystyka matematyczna zajmuje się metodami wnioskowania o całej zbiorowości statystycznej na podstawie zbadania pewnej jej części zwanej próbką lub próbą.
- Model statystyczny [2, w oparciu o rozdział 2]
Modelem statystycznym nazywamy przestrzeń próby doświadczenia tj. wartości zmiennych losowych o jednakowym rozkładzie, rodzinę podzbiorów zbioru zmiennych oraz prawdopodobieństwo występowania danej zmiennej. Można tutaj wskazać analogię do rachunku prawdopodobieństwa, tj. uporządkowanej trójki (Ω, \mathcal{F}, P) .
- Cecha statystyczna [9, Rozdział 1]
Cecha statystyczna jest to liczbowy opis przedmiotu dociekań tj. konkretnej dziedziny życia społeczno-gospodarczego. Służy ona do scharakteryzowania podmiotu badania.

Definicja 2.2.1. *Macierz [1, w oparciu o rozdział 1.1]*

Niech X będzie skończonym podzbiorem liczb \mathbb{R} . Macierzą wymiaru $m \times n$ (tzn. o m wierszach i n kolumnach) nazywamy prostokątną tablicę utworzoną z elementów zbioru X , postaci:

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

gdzie:

$a_{ij} \in X$, dla $1 \leq i \leq m$, $1 \leq j \leq n$

- Macierz obserwacji [9, Rozdział 2]
Niech $m > 1$ oraz $n > 1$ będą liczbami naturalnymi. Macierzą 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 2.2.2. *Macierz odległości zmiennych [10, Rozdział 1.6]*

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

$$D = [d_{ii'}].$$

gdzie:

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

- Skala [10, w oparciu o rozdział 1.2]
Skalą nazywamy pewien skończony zbiór, umożliwiający porównywanie obiektów na podstawie wartości lub własności wybranych zmiennych.
- Skala porządkowa [10, Rozdział 1.2]
Skala ta pozwala na stwierdzeniu o identyczności lub różnicy porównywanych obiektów, a także na porównywanie wariantów cech zaobserwowanych w obiektach. Nie pozwala ona określić odległości między obiektami. Umożliwia zliczanie obiektów uporządkowanych (liczby relacji równości, nierówności, większości i mniejszości). Przykład zmiennych przedstawianych na skali porządkowej: wykształcenie, kolejność zawodników na podium, oceny w systemie szkolnym.
- Skala przedziałowa [10, Rozdział 1.2]
Jest to skala, która w stosunku do skali porządkowej, pozwala obliczyć odległość między obiektami, dokonując pomiaru cech za pomocą liczb rzeczywistych. Dla skali tej możliwe jest korzystanie z operacji dodawania oraz odejmowania. Dla skali tej istnieje charakterystyczna wartość - punkt zerowy. Jest on wyznaczany w sposób umowny, umożliwia on do zachowania różnic między wartościami cechy, przy zmianie jednostek miary. Przykład zmiennych przedstawianych na skali przedziałowej: temperatura, rok urodzenia.
- Skala ilorazowa [10, Rozdział 1.2]
Skala ta, podobna jest do skali przedziałowej, z tym że występuje w niej zero bezwzględne - punkt, który mówi o tym, że dana zmienna nie występuje, oraz ogranicza lewostronnie zakres skali ilorazowej. Powoduje to, że można na tej skali obok operacji dodawania i odejmowania, dokonywać także dzielenia i mnożenia, a tym samym przedstawiać dowolną wartość cechy danego obiektu jako wielokrotność wartości cechy dla innego obiektu. Przykład zmiennych przedstawianych na skali ilorazowej: napięcie elektryczne, bezrobocie, inflacja.
- Stymulanta [10, Rozdział 1.5]
Stymulantami nazywane są zmienne, których wysokie wartości badany w badanych obiektach są pożądane z punktu widzenia rozpatrywanego zjawiska.
- Destymulanta [10, Rozdział 1.5]
Destymulantami nazywane są zmienne, których wysokie wartości badany w badanych obiektach są niepożądane z punktu widzenia rozpatrywanego zjawiska.
- Nominanta [10, Rozdział 1.5]
Nominantami nazywane są zmienne, których odchylenia wartości w badanym obiekcie od wartości (lub przedziału wartości) uznawanych za najkorzystniejsze są niepożądane z punktu widzenia rozpatrywanego zjawiska.

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 [11, 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 [11, 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

$$\bigcup_{i=1}^{\infty} A_n \in \mathcal{F}.$$

Definicja 2.3.3. Zbiory borelowskie [4, w oparciu 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 $\mathcal{B}(X)$.

Definicja 2.3.4. Przestrzeń zdarzeń elementarnych [7, 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 2.3.5. Miara zbioru [4, Rozdział 2.10]

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

1. $\mu(A) \in [0, \infty]$ dla każdego zbioru $A \in \mathcal{F}$,
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 2.3.6. Przestrzeń mierzalna [4, Rozdział 2.10]

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

Definicja 2.3.7. Funkcja mierzalna [11, 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 $a \in \mathbb{R}$.

Definicja 2.3.8. Zdarzenie losowe [7, 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 [7, 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. $P(A) \geq 0 \quad \forall A \in \mathcal{F}$
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 [7, 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} .

Definicja 2.3.11. Zmienna losowa [7, Rozdział 2.1]

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

$$\{\omega : X(\omega) < x\} \in \mathcal{F} \text{ dla każdego } x \in \mathbb{R}$$

Definicja 2.3.12. *Wektor losowy [6, Rozdział 5.1]*

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

Definicja 2.3.13. *Rozkład prawdopodobieństwa zmiennej losowej X [6, Rozdział 5.1]*

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

$$\mu_X(B) = P_X(B) = P(X^{-1}(B)), \quad B \in \mathcal{B}(\mathbb{R})$$

Definicja 2.3.14. *Rozkład dyskretny [6, Rozdział 5.1]*

Mówimy, że zmienna losowa X ma rozkład dyskretny, jeśli istnieje przeliczalny zbiór $S \in \mathbb{R}$, taki że $\mu_X(S) = 1$.

Definicja 2.3.15. *Gęstość i rozkład ciągły [6, Rozdział 5.1]*

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

$$\mu(A) = \int_A f(x)dx, \quad A \in \mathcal{B}(\mathbb{R})$$

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

Definicja 2.3.16. *Wartość oczekiwana [7, Rozdział 2.6]*

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

$$E(X) = \mu_X = \begin{cases} \sum_{i=1}^n x_i p_i, & \text{jeśli zmienna ma rozkład dyskretny} \\ \int_{-\infty}^{\infty} x f(x) dx, & \text{jeśli zmienna ma rozkład ciągły} \end{cases}$$

Definicja 2.3.17. *Wariancja [6, Rozdział 5.6]*

Niech X będzie zmienną losową o skończonej wartości oczekiwanej tj. $E|X| < \infty$. Wariancję zmiennej losowej X nazywamy liczbę

$$\text{Var} X = \mathcal{D}^2 X = E(X - EX)^2.$$

W przypadku zmiennej losowej o rozkładzie dyskretnym lub ciągłym mamy

$$\text{Var} X = \begin{cases} \sum_{i=1}^n (x_i - \mu_X) p_i = \sum_{i=1}^n x_i^2 p_i - (\mu_X)^2, & \text{jeśli zmienna ma rozkład dyskretny} \\ \int_{-\infty}^{\infty} (x - \mu_X)^2 f(x) dx = \int_{-\infty}^{\infty} x^2 f(x) dx - (\mu_X)^2, & \text{jeśli zmienna ma rozkład ciągły} \end{cases}$$

Definicja 2.3.18. *Odchylenie standardowe [6, Rozdział 5.6]*

Niech X będzie zmienną losową. Odchyleniem standardowym nazywamy pierwiastek z wariancji.

$$\sigma_X = \sqrt{\mathcal{D}^2 X}$$

Definicja 2.3.19. *Rozkład normalny (Gaussa) [6, Rozdział 5.10]*

Jeśli zmienna losowa X ma gęstość postaci

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

dla $x \in \mathbb{R}$ i pewnych $\mu_X \in \mathbb{R}$ 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

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

2.4 Podstawowe pojęcia teorii grafów

W pracy zostaną opisane zarówno metody porządkowania liniowego jak i nieliniowego tzn. w ujęciu matematycznym - porządku częściowego, w tym celu należy wprowadzić definicje związane z teorią grafów, niezbędne przy opisywaniu metod porządkowania nieliniowego.

W celu wprowadzenia złożonych definicji, należy wcześniej podać podstawowe pojęcia dotyczące grafów. Upřednio zostanie jeszcze wprowadzona definicja pary uporządkowanej oraz nieuporządkowanej, gdyż pojęcia te zostały wykorzystane w definicji grafu.

Definicja 2.4.1. Para uporządkowana [8, w oparciu o rozdział 3]

Niech dane będą dwa elementy a i b . Parą 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 2.4.2. Para nieuporządkowana [8, w oparciu o rozdział 3]

Niech dane będą dwa elementy a i b . Parą nieuporządkowaną nazywamy zbiór postaci $\{a, b\}$, 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.

Definicja 2.4.3. 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 2.4.4. Pętle [12, Rozdział 2]

Pętlami nazywamy krawędzie wielokrotne, łączące wierzchołek z samym sobą.

Definicja 2.4.5. Trasa/marszruta [12, Rozdział 3]

Trasą (lub marszrutą) 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 od wierzchołka v_0 do wierzchołka v_m . Liczbę krawędzi na trasie nazywamy długością trasy.

Definicja 2.4.6. Ścieżka [12, Rozdział 3]

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

Definicja 2.4.7. Droga [12, Rozdział 3]

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

Definicja 2.4.8. Droga zamknięta/ścieżka zamknięta [12, Rozdział 3]

Droga lub ścieżka jest zamknięta, jeśli $v_0 = v_m$.

Definicja 2.4.9. Cykl [12, Rozdział 3]

Ścieżką zamkniętą zawierającą co najmniej jedną krawędź nazywamy cyklem.

Definicja 2.4.10. *Graf spójny [12, Rozdział 3]*

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

Definicja 2.4.11. *Drzewo [12, Rozdział 4]*

Drzewem nazywamy graf spójny, nie zawierający cykli.

Definicja 2.4.12. *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 $A(D)$ par uporządkowanych elementów zbioru $V(D)$, nazywanych łukami. Zbiór $V(D)$ nazywamy zbiorem wierzchołków, a rodzinę $A(D)$ rodziną łuków digrafu D . Łuk (v, w) zwykle zapisujemy jako vw . Graf skierowany oznaczamy zwykle w postaci pary uporządkowanej $G = \langle V, E \rangle$

Uwaga 1. *Każdy graf jednoznacznie wyznacza pewną relację dwuargumentową (binarną) w zbiorze V . Można również powiedzieć odwrotnie, że każda relacja dwuargumentowa (binarna) r w zbiorze V , wyznacza jednoznacznie graf zorientowany, którego węzłami są elementy zbioru V , z kolei krawędziami są uporządkowane pary (v, v') , należące do r .*

Definicja 2.4.13. *Graf niezorientowany [10, Rozdział 2.3]*

Grafem niezorientowanym, nazywamy graf $G = \langle V, E \rangle$, jeżeli relacja binarna tego grafu jest symetryczna, tj. dla dowolnych wierzchołków $v, v' \in V$, $(v, v') \in E$ wtedy i tylko wtedy $(v', v) \in E$.

2.5 Wybrane pojęcia z teorii mnogości i topologii: relacja porządkująca, moc zbiorów, równoliczność, przestrzenie metryczne, miary odległości

2.5.1 Relacja porządkująca

W niniejszej pracy skupiam 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 ze względu na korzystanie z tych pojęć, przy definiowaniu właściwości porządkowania liniowego zbioru obiektów. Dodatkowo przytoczona została definicja zbioru skończonego, ze względu na zastosowanie metod porządkowania na zbiorze skończonym.

Definicja 2.5.1. *Relacja [8, 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 2.5.2. *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 1. *Częściowego porządku na zbiorze*

Wykorzystanie częściowego porządku na płaszczyźnie \mathbb{R}^2 , obrazuje diagram Hassego, będący grafem skierowanym, którego wierzchołki zostały poddane relacji porządkowania i reprezentują elementy skończonego zbioru $X \subset \mathbb{R}$. 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 $y \in X$, jeśli x jest następnikiem y , czyli gdy $y < x$ oraz nie istnieje taki punkt $z \in X$, że $y < z < x$.

Definicja 2.5.3. *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.

$$x \leq y \text{ lub } y \leq x$$

Parę (X, \leq) nazywamy zbiorem liniowo uporządkowanym lub łańcuchem.

Definicja 2.5.4. *Dobry porządek [3, Rozdział 2]*

Niech dany będzie zbiór X . Relację \leq porządkującą zbiór X , nazywamy dobrym porządkiem na zbiorze X , gdy w każdym niepustym podzbiorze zbioru X istnieje element najmniejszy względem relacji \leq . Jeśli relacja \leq na zbiorze X jest dobrym porządkiem, to mówimy, że para (X, \leq) jest zbiorem dobrze uporządkowanym.

Definicja 2.5.5. *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 $a \leq x$ dla każdego $a \in A$.

Definicja 2.5.6. *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 $y \leq a$ dla każdego $a \in A$.

Definicja 2.5.7. *Zbiór ograniczony z góry, zbiór ograniczony z dołu, 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 2.5.8. *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. $a \leq x_0$ dla każdego $a \in A$,
2. jeśli $a \leq x$ dla każdego $a \in A$, to $x_0 \leq x$.

Definicja 2.5.9. *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 $x_0 = \inf A$, gdy spełnione są następujące warunki:

1. $y_0 \leq a$ dla każdego $a \in A$,
2. jeśli $y \leq a$ dla każdego $a \in A$, to $y \leq y_0$.

Definicja 2.5.10. *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 2.5.11. *Zbiór skończony [3, Rozdział 5]*

Mówimy, że zbiór A jest skończony, gdy istnieje taka liczba naturalna n , że zbiór A jest równoliczny z przedziałem w zbiorze liczb naturalnych

$$[0, n) = \{k \in \mathbb{N} : k < n\}.$$

Zatem, gdy zbiór A jest równoliczny ze zbiorem $\{1, \dots, k\}$, to mówimy że jest on k -elementowy, tj. mocy równej k .

2.5.2 Przestrzenie metryczne, miary odległości

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

Definicja 2.5.12. *Metryka [8, 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 2.5.13. *Przestrzeń metryczna [8, Rozdział 9]*

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

Przykład 2. *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 3. *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 4. *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 [10, Rozdział 2], omówiono w nim ogólnie metody porządkowania zarówno liniowego jak i nieliniowego, po to by w kolejnym rozdziale szczegółowo przyjrzeć się wybranym metodą wraz z przedstawieniem ich algorytmów oraz dokładnych opisów matematycznych, z uwzględnieniem ich wad oraz zalet.

Metody porządkowania liniowego pozwalają na ustalenie hierarchii obiektów ze względu na określone kryterium. Z kolei metody porządkowania nieliniowego nie pozwalają na ustalenie hierarchii, natomiast w wyniku uporządkowania możliwe jest wskazanie dla każdego z obiektów poddanych porządkowaniu, wskazaniu obiektów podobnych ze względu na opisujące je zmienne.

3.1 Metody porządkowania liniowego

Porządkowanie liniowe obiektów polega, w ujęciu geometrycznym, na rzutowaniu na prostą punktów reprezentujących obiekty, umieszczonych w wielowymiarowej przestrzeni zmiennych. Takie postępowanie pozwala na ustalenie hierarchii obiektów, czyli uporządkowanie ich od obiektu stojącego najwyżej w tej hierarchii do obiektu znajdującego się najniżej. Poniżej zostaną przedstawione własności uporządkowania liniowego obiektów, wraz z podaniem ich matematycznej interpretacji.

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

Powyżej wymienione własności są wynikiem posiadania jedynie skończonej ilości obiektów, które podane są uporządkowaniu. W następnej części chciałabym

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

Na początku zacznę 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. Mówi się, 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 $X \times X$, 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. $\overline{\{x \in X, \overline{\{y \in X, \overline{\{y \in X, ySx\}} = 1\}} = 1\}} = 2$

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 2.5.3 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 ogólności 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 oznaczeń możemy przyjąć, że będą nimi a oraz b . Ustalmy zatem, że $a \leq x, b \leq x$. Wobec definicji 2.5.3 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. Zupełnie analogicznie postępujemy dla przypadku $x \leq a, x \leq b$. Zatem uzyskujemy 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 $x, y \in X$ oraz niech ySx . Korzystając z definicji sąsiada 3.1.1 mamy, że skoro ySx to

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

dotatkowo jeśli xSy , to

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

Zatem łącząc te dwa podejścia otrzymalibyśmy, że $ySx \wedge xSy \Rightarrow xSy \wedge ySx$, stąd ostatecznie, 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. Istotnie 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 2.5.2 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 zatem oznaczać

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

Wobec spójności z definicji 2.5.3 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 2.5.2 mamy, że $x_M = y_M$, co wbrew naszemu przypuszczeniu daje, że elementy te 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 minimalny to musi to zatem oznaczać

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

Wobec spójności z definicji 2.5.3 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 $y_1 \leq x$.

Zajmiemy się najpierw pierwszą częścią, tj. x_m jest sąsiadem x . Na mocy definicji sąsiada 3.1.1, wynika, że skoro x_m jest sąsiadem elementu x oraz element x ma pojedynczego sąsiada i nie jest ani elementem maksymalnym, ani minimalnym, to

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

Ale z faktu, że x_m jest elementem minimalnym, wynika że $x_m \leq x$, zatem x nie może zachodzić $x \leq x_m$. Zauważmy teraz, że z faktu, iż x_M jest elementem maksymalnym zbioru X , wynika że $x \leq x_M$, oraz dodatkowo istnieje takie $x_1 \in X$, że $x \leq x_1 \leq x_M$. Dalej, istnienie x_1 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 $x \leq x_j$, 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 drugiej części, tj. istnieje $y_1 \in X$, taki że $y_1 \leq x$. Zaczniemy najpierw od założenia, że x_M jest sąsiadem x . Wówczas analogicznie do powyższego i definicji sąsiada 3.1.1, mamy że skoro x_M jest sąsiadem x i element ten nie jest ani maksymalny ani minimalny to, z definicji sąsiada mamy

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

Ale z faktu, że x_M jest elementem maksymalnym, wynika że $x \leq x_M$. Dokładając również istnienie elementu minimalnego w zbiorze X , mamy że $x_m \leq x$. Co więcej istnieje taki element $y_1 \in X$, że $x_m \leq y_1 \leq x$. Dalej istnienie y_1 powoduje, 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.

□

Powyżej wykazane własności są często podawane, niemal na równi z definicją takiego uporządkowania, w pozycjach książkowych omawiających praktyczny aspekt porządkowania danych. Poniżej podane zostaną przykłady takich relacji, które też posiadają powyższy zestaw własności, ale nie opisują relacji będących porządkami liniowymi.

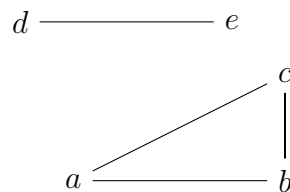
Przykład 5. Rozważmy zbiór dwuelementowy $X = \{a, b\}$ gdzie $a \leq b$ jest jedynym punktem tej relacji. Tak zdefiniowana relacja spełnia wszystkie własności, ale nie spełnia założenia o zwrotności - zatem relacja ta nie jest liniowym porządkiem.

Diagram Hassego prezentując tę relację jest postaci:



Przykład 6. 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, charakteryzujące je zmienne muszą być mierzone przynajmniej na skali porządkowej. Gdy zmienne te mierzone są na skali przedziałowej lub ilorazowej, należy dokonać ich normalizacji, dla zapewnienia ich porównywalności.

Metody porządkowania liniowego można podzielić na metody diagramowe, procedury oparte na zmiennej syntetycznej oraz procedury iteracyjne bazujące funkcji kryterium dobroci uporządkowania tzn. funkcji, którą się przyjmuje, lub też tworzy się ją aby w kolejnych iteracjach szukać takiego uporządkowania, które optymalizuje zbiór wartości tej funkcji. W kolejnej sekcji zostaną pokrótce przedstawione różne metody, z wyszczególnieniem najważniejszych założeniach 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, tj. przestawieniu wierszy i odpowiadających im kolumn diagramu, tak aby symbole graficzne reprezentujące najmniejsze odległości skupiały się wzdłuż głównej przekątnej, zaś w miarę oddalania się od głównej przekątnej znajdowały się symbole graficzne odpowiadające coraz to większym odległością.

Jako narzędzie pomocnicze w porządkowaniu danych, może stanowić kryterium postaci:

$$F^1 = \sum_{i=1}^n \sum_{i'>1}^n d_{ii'} w_{ii'}$$

gdzie:

$d_{ii'}$ - odległość euklidesowa między i -tym i i' -tym obiektem .

$w_{ii'}$ - wagi elementów macierzy odległości, zdefiniowane w oparciu o jeden z następujących wzorów:

$$w_{ii'} = \frac{|i-i'|}{n-1},$$

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

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

Dodatkowo wagi elementów macierzy odległości tworzą macierz wag postaci:

$$W = [w_{ii'}], \quad i, i' = 1, 2, \dots, n.$$

3.1.2 Metody oparte na zmiennych syntetycznych

W tym podrozdziale zostaną opisane metody porządkowania oparte na zmiennych syntetycznych, tj. funkcji których wartości będą służyć do porządkowania danych. 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

1. dla średniej arytmetycznej:

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

2. dla średniej geometrycznej:

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

3. dla średniej harmonicznnej

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

gdzie:

s_i - wartość zmiennej syntetycznej w i -tym obiekcie,

w_j - waga j -tej zmiennej.

Metody bezwzorcowe

W metodach tych, unormowane wartości podanych zmiennych wejściowych są uśrednianie, przez przypisywanie im odpowiednich wag. Poniżej zostaną omówione wybrane metody porządkowania bezwzorcowego.

Metoda rang

Metoda ta bazuje na normalizacji rangowej, w związku z tym zmienne mierzone są 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. Wzór na obliczenie wartości zmiennej syntetycznej:

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

gdzie:

z_{ij} -zmienna znormalizowana rangowo, tj.

$$z_{ij} = h \text{ dla } x_{hj} = x_{ij}, \quad h, i = 1, 2, \dots, n.$$

gdzie:

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

Metoda sum

Metoda ta bazuje na konstrukcji zmiennej syntetycznej przy pomiarze zmiennych na skali ilorazowej lub przedziałowej. Dla każdego obiektu obliczana jest wartość zmiennej syntetycznej, jako średnia arytmetyczna wartości zmiennych przy przyjęciu jednakowych wag dla każdej zmiennej. Eliminowane są ujemne wartości zmiennej syntetycznej przy wykorzystaniu przekształcenia:

$$s'_i = s_i - \min\{s_i\}, \quad i = 1, 2, \dots, n.$$

Końcowa postać zmiennej syntetycznej otrzymywana jest po przeprowadzeniu normalizacji według wzoru:

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

Powyższe przekształcenia powodują unormowanie miary syntetycznej w przedziale $[0,1]$. Powyższa wielkość wykorzystywana jest do uporządkowania obiektów.

Metoda wzorcowa

W metodach tych zakłada się istnienie obiektu wzorcowego $O_0 = [z_{0j}]$, $j = 1, 2, \dots, m$, w którym zmienne wejściowe z_{0j} , $j = 1, 2, \dots, m$, będące współrzędnymi obiektu wzorcowego, przyjmują optymalne wartości, które to mogą być ustalane na podstawie ogólnie przyjętych norm, subiektywnej opinii dotyczącej obserwowanego obiektu, lub też opinii ekspertów. Poszczególne metody różnią się sposobem wyznaczania obiektu wzorcowego, poniżej zostaną one przedstawione.

Metoda Hellwiga

W metodzie tej, obiekt wzorcowy wyznaczony jest na podstawie wystandaryzowanych zmiennych wejściowych. Współrzednym 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. Miara syntetyczna jest postaci:

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

gdzie:

współrzedne obiektu wzorcowego są obliczane na podstawie wzoru:

$$z_{0j} = \begin{cases} \max_i \{z_{ij}\} & \text{dla } z_j^S, \quad j = 1, 2, \dots, m, \quad i = 1, 2, \dots, n, \\ \min_i \{z_{ij}\} & \text{dla } z_j^D, \quad j = 1, 2, \dots, m, \quad i = 1, 2, \dots, n, \end{cases}$$

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

$$d_0 = \overline{d_0} + 2S(d_0),$$

$$\overline{d_0} = \frac{1}{n} \sum_{i=1}^n d_{i0},$$

$$S(d_0) = \left[\frac{1}{n} \sum_{i=1}^n (d_{i0} - \overline{d_0})^2 \right]^{\frac{1}{2}}.$$

Wartości miary s_i zazwyczaj są z przedziału $[0, 1]^2$. Należy tu zaznaczyć, że wartości miary są tym wyższe, im mniej jest oddalony obiekt od obiektu wzorcowego.

Metoda Walesiaka

Metoda ta bazuje na konstrukcji zmiennej syntetycznej w oparciu o badanie odległości obiektów od obiektu wzorcowego, przy wykorzystaniu uogólnionej miary odległości. Umożliwia ona porządkowanie obiektów, jeżeli opisujące je charakterystyki są mierzone przynajmniej na skali porządkowej. W takim przypadku, zmienne wejściowe o postaci nominant muszą zostać podane stymulacji. Z kolei gdy zmienne są mierzone na skali przedziałowej lub ilorazowej, należy je znormalizować. Miara syntetyczna oparta na uogólnionej mierze odległości przyjmuje postać:

$$s_i = \frac{1}{2} - \frac{\sum_{j=1}^m w_j a_{i0j} b_{0ij} + \sum_{j=1}^m \sum_{i''=1}^n w_j a_{ii''j} b_{0i''j}}{2 \left[\left(\sum_{j=1}^m \sum_{i''=1}^n w_j a_{ii''j}^2 \right) \cdot \left(\sum_{j=1}^m \sum_{i''=1}^n w_j b_{0i''j}^2 \right) \right]^{\frac{1}{2}}} \quad (3.1)$$

gdzie:

$$w_j \in [0, m]$$

$$\sum_{j=1}^m w_j = m$$

Ostateczna postać zmiennej syntetycznej zależy od skali pomiaru zmiennych.

Jeśli zmienne charakteryzujące obiekty mierzone są na skali ilorazowej lub przedziałowej, stosowane jest następujące podstawienie:

$$a_{ii^*j} = z_{ij} - z_{i^*j} \text{ dla } i^* = 0, i'',$$

$$b_{0i^*j} = z_{0j} - z_{i^*j} \text{ dla } i^* = i, i''.$$

gdzie:

z_{0j} -wystandaryzowana wartość j -tej zmiennej dla obiektu wzorcowego

Z kolei, gdy zmienne charakteryzujące obiekty mierzone są na skali porządkowej to stosowne jest podstawienie:

$$a_{ii^*j} = \begin{cases} 1 & \text{dla } z_{ij} > z_{i^*j}, \\ 0 & \text{dla } z_{ij} = z_{i^*j}, \quad i^* = 0, i', \\ -1 & \text{dla } z_{ij} < z_{i^*j}, \end{cases} \quad (3.2)$$

$$b_{0i^*j} = \begin{cases} 1 & \text{dla } z_{0j} > z_{i^*j}, \quad i^* = i, i', \\ 0 & \text{dla } z_{0j} = z_{i^*j}, \quad i^* = i, i', \\ -1 & \text{dla } z_{0j} < z_{i^*j}. \end{cases} \quad (3.3)$$

Zmienna syntetyczna przyjmuje wartości z przedziału $[0, 1]$. Czym niższa wartość zmiennej syntetycznej, tym bliżej wzorca leży dany obiekt.

Metoda dystansowa

W metodzie tej zmienna syntetyczna wyznaczana jest na podstawie odległości każdego obiektu od obiektu wzorca, przy wykorzystaniu np. metryki euklidesowej. Miara syntetyczna, przy wykorzystaniu przekształcenia unitaryzacyjnego, jest postaci:

$$s_i = \left(\frac{d_{i0} - \min_i \{d_{i0}\}}{\max_i \{d_{i0}\} - \min_i \{d_{i0}\}} \right)^p, \quad i = 1, 2, \dots, m \quad (3.4)$$

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

3.1.3 Metody iteracyjne

W metodach tych przyjmowania jest funkcja kryterium dobroci porządkowania i 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

W metodzie tej poszukiwane jest takie liniowe uporządkowanie obiektów, dla którego funkcja kryterium dobroci uporządkowania osiąga maksimum:

$$F^2 = \sum_{i'=1}^{n-1} i' \sum_{i=1}^{n-i'} d_{ii'} \rightarrow \max \quad (3.5)$$

gdzie:

$d_{ii'}$ - odległość euklidesowa między i -tym i i' -tym obiektem.

Sposób postępowania:

W pierwszym kroku dokonywane jest dowolne liniowe uporządkowanie obiektów, dla którego obliczana jest wartość funkcji kryterium (3.5). W kolejnym etapie obliczana jest wartość funkcji kryterium dla każdej możliwej transpozycji pary obiektów. Uporządkowanie to stanowi punkt wyjścia do oceny, czy kolejna transpozycja dowolnej pary obiektów pozwoli na wzrost wartości funkcji kryterium. Powyższe postępowanie kontynuowane jest tak długo, aż transpozycja dowolnej pary obiektów nie prowadzi do wzrostu wartości funkcji kryterium.

3.1.4 Metody gradientowe

W metodach gradientowych dąży się do takiego liniowego uporządkowania obiektów, które jak najmniej zniekształca relacje strukturalne porządkowanego zbioru obiektów. Od strony geometrycznej oznacza to, że odległości pomiędzy punktami reprezentującymi obiekty w przestrzeni jednowymiarowej, określonej przez zmienną syntetyczną, w jak najmniejszym stopniu zniekształcają odległości pomiędzy tymi punktami w przestrzeni wielowymiarowej, określonej przez zmienne wejściowe. Metody gradientowe poszukują takich współrzędnych punktów reprezentujących obiekty w przestrzeni jednowymiarowej, dla których funkcja dobroci uporządkowania osiąga minimum, co można przedstawić za pomocą wariantów:

$$F^3 = \frac{\sum_{\substack{i,i'=1 \\ i \neq i'}}^n (d_{ii'}^s - d_{ii'})^2}{\sum_{\substack{i,i'=1 \\ i < i'}}^n d_{ii'}} \rightarrow \min \quad (3.6)$$

$$F^4 = \sum_{\substack{i,i'=1 \\ i < i'}}^n \left(\frac{d_{ii'}^s - d_{ii'}}{d_{ii'}} \right)^2 \rightarrow \min \quad (3.7)$$

$$F^5 = \frac{1}{\sum_{\substack{i,i'=1 \\ i \neq i'}}^n d_{ii'}} \sum_{\substack{i,i'=1 \\ i < i'}}^n \frac{\left(d_{ii'}^s - d_{ii'} \right)^2}{d_{ii'}} \rightarrow \min \quad (3.8)$$

gdzie:

$d_{ii'}$ - odległość euklidesowa między i -tym i i' -tym obiektem.

3.2 Metody porządkowania nieliniowego

Metody porządkowania nieliniowego nie pozwalają na ustaleniu hierarchii obiektów, lecz na określeniu dla każdego z nich, stopnia podobieństwa do innych obiektów, ze względu na ich charakterystyki.

Aby uporządkować nieliniowo obiekty, charakteryzujące je zmienne powinny być mierzone przynajmniej na skali przedziałowej lub ilorazowej. Gdy zmienne te mierzone są na skali przedziałowej lub ilorazowej, należy dokonać ich normalizacji, dla zapewnienia ich porównywalności.

Metody porządkowania nieliniowego można podzielić na metody dendrytowe i metody aglomeracyjne. Metody dendrytowe prowadzą do powstania dendrytu, będącego ilustracją graficzną położenia względem siebie obiektów ze względu na ich podobieństwo. Z kolei metody aglomeracyjne prowadzą do utworzenia drzewka połączeń, będącego graficzną ilustracją hierarchii łączenia obiektów, ze względu na ich podobieństwo.

3.2.1 Metody dendrytowe

Metody dendrytowe opierają się na regułach i pojęciach teorii grafów. Porządkowanie dendrytowe polega na przyporządkowaniu obiektom poszczególnych wierzchołków dendrytu, w tym celu budowany jest dendryt. Poniżej zostaną opisane przykłady metod dendrytowych, tj. taksonomia wrocławska oraz metoda Prima.

Taksonomia wrocławska

W metodzie tej obiekty dzielone są grupy obiektów najbardziej do siebie podobnych, tj. takich, których odległość między sobą jest najmniejsza. W tym celu w każdym wierszu(kolumnie) macierzy odległości D , wyznaczamy jest element najmniejszy:

$$d_{ii'} = \max_i d_{ii'}, \quad i, i' = 1, 2, \dots, n, i \neq i'. \quad (3.9)$$

Otrzymane pary najbardziej podobnych do siebie obiektów, przedstawiane są w postaci grafu niezorientowanego, Długość krawędzi łączących wierzchołki grafu, są proporcjonalne do odległości między obiektami. Może się zdarzyć, że wśród wyznaczonych par połączeń, pojawiają się połączenia występujące dwukrotnie, jedno z nich zostanie wyeliminowane, ponieważ kolejność połączeń w dendrycie nie jest istotne. Warto również zwrócić uwagę, na fakt iż w dendrycie danych obiekt może występować tylko jeden raz, w związku z tym jeżeli w łączeniu występują wielokrotnie te same obiekty, to zostaną one połączone w zespoły zwane skupieniami. Metoda kończy działanie, w momencie uzyskania grafu spójnego.

Metoda Prima

W odróżnieniu od taksonomii wrocławskiej, metoda Prima nie wymaga operowania całym czasem pełną, wyjściową macierzą odległości. W trakcie tworzenia dendrytu, na każdym etapie zbiór porządkowanych obiektów jest klasyfikowany do jednego z dwóch podzbiorów, np. A i B . Niech zbiór A będzie pierwszym z nich a zbiór B drugim. Pierwszy z nich zawiera obiekty należące na danym etapie do dendrytu, zaś drugi zawiera obiekty nie należące na tym etapie do dendrytu.

Na początku procedury, zbiór A jest zbiorem pustym, z kolei zbiór B zawiera wszystkie obiekty. W pierwszym kroku do zbioru A zostaje włączony dowolny obiekt, nie ma to wpływu na ostateczną postać dendrytu. Następnie do zbioru A zostają włączone te obiekty zbioru B , najbardziej podobne do obiektów należących już do zbioru A . W tym celu w pierwszym kroku algorytmu zostaje stworzony wektor d , zawierający odległości wybranego obiektu zbioru A , od pozostałych obiektów zbioru B .

W powstałym dendrycie wierzchołkami są obiekty przechodzące kolejno do zbioru A, z kolei wiązałkami łączącymi wierzchołki są minimalne wartości elementów wektora d, otrzymanego w kolejnych krokach przyłączania obiektów do dendrytu.

3.2.2 Metody aglomeracyjne

Istotą metod aglomeracyjnych jest utworzenie drzewka połączeń - dendrogramu. W ten sposób zobrazowana jest hierarchia łączenia obiektów, na podstawie zmniejszającego się podobieństwa między obiektami włączonymi do dendrogramu, w kolejnych etapach a obiektami należącymi już do dendrogramu. Hierarchia połączeń określa wzajemnie położenie względem siebie obiektów oraz grup obiektów powstających w kolejnych etapach tworzenia drzewka. Grupy podobnych do siebie obiektów tworzą oddzielne gałęzie. Punktem wyjściem metod aglomeracyjnych stanowi założenie, że każdy obiekt stanowi odrębną, jednoelementową grupę ($G_r, r = 1, 2, \dots, z$).

W kolejnych krokach następuje łączenie ze sobą grupy obiektów najbardziej do siebie podobnych ze względu na wartości opisujących je zmiennych. Podobieństwo weryfikowane jest na podstawie odległości między grupami obiektów.

Na początku odległości między jednoelementowymi grupami obiektów G_1, \dots, G_z są elementami wyjściowej macierzy odległości **D**. W macierzy **D** poszukiwane są najmniejsze odległości pomiędzy grupami obiektów:

$$d_{rr'} = \min_{ii'} d_{ii'}, \quad i = 1, 2, \dots, n_r, \quad i' = 1, 2, \dots, n_{r'}, \quad r, r' = 1, 2, \dots, z, r \neq r'. \quad (3.10)$$

gdzie:

$d_{rr'}$ - odległość r -tej od r' -tej grupy.

Ogólny wzór wyznaczania odległości nowo powstałej grupy $G_{r''}$, powstałej w wyniku połączenia grup G_r i $G_{r'}$, od pozostałych grup $G_{r'''}$, przy tworzeniu drzewka połączeń ma postać:

$$d_{r''r'''} = \alpha_r d_{rr'''} + \alpha_{r'} d_{r'r'''} + \beta d_{rr'} + \gamma |d_{rr'''} - d_{r'r'''}| \quad (3.11)$$

gdzie:

$\alpha_r, \alpha_{r'}, \beta, \gamma$ - współczynniki przekształceń odmienne dla różnych metod aglomeracyjnych

Poszczególne metody aglomeracyjne, różnią się między sobą sposobami wyznaczania odległości między obiektami. Poniżej zostały wymienione najczęściej stosowane metody aglomeracyjne, które będą dokładniej omówione w kolejnym podrozdziale.

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

Metoda najbliższego sąsiedztwa

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

$$d_{rr'} = \min_{ii'} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}), \quad (3.12)$$

$$i = 1, 2, \dots, n_r, \quad i' = 1, 2, \dots, n_{r'}, \quad r, r' = 1, 2, \dots, z, \quad r \neq r',$$

gdzie:

$$\mathbf{O}_i = [z_{ij}], \quad j = 1, 2, \dots, m. \quad (3.13)$$

Metoda najdalszego sąsiedztwa

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

$$d_{rr'} = \max_{ii'} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}), \quad (3.14)$$

$$i = 1, 2, \dots, n_r, \quad i' = 1, 2, \dots, n_{r'}, \quad r, r' = 1, 2, \dots, z, \quad r \neq r',$$

Metoda średniej międzygrupowej

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

$$d_{rr'} = \frac{1}{n_r n_{r'}} \sum_{i'=1}^{n_{r'}} \sum_{i=1}^{n_r} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}) \quad (3.15)$$

$$r, r' = 1, 2, \dots, z, \quad r \neq r'.$$

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:

$$d_{rr'} = \text{med}_{i,i'} \{d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'})\}, \quad (3.16)$$

$$i = 1, 2, \dots, n_r, \quad i' = 1, 2, \dots, n_{r'}, \quad r, r' = 1, 2, \dots, z, \quad r \neq r'.$$

Metoda środków ciężkości

W metodzie tej odległość między dwoma grupami jest równa odległości między środkami ciężkości tych grup. Odległość ta opisana jest wzorem:

$$d_{rr'} = d_{i_c i'_c}(\mathbf{O}_i^c = \bar{\mathbf{O}}_r \in \mathbf{G}_r, \mathbf{O}_{i'_c} = \bar{\mathbf{O}}_{r'} \in \mathbf{G}_{r'}), \quad (3.17)$$

$$i = 1, 2, \dots, n_r, \quad i' = 1, 2, \dots, n_{r'}, \quad r, r' = 1, 2, \dots, z, \quad r \neq r'.$$

gdzie:

$d_{i'c}$ - odległość środka ciężkości r -tej grupy od środka ciężkości r' -tej grupy,

$\bar{\mathbf{O}}_{i'c}, \bar{\mathbf{O}}_{i''c}$ - środki ciężkości odpowiednio r -tej i r' -tej grupy obiektów. przy czym:

$$\mathbf{O}_{i'c} = \bar{\mathbf{O}}_r = \frac{1}{n_r} \sum_{i=1}^{n_r} \mathbf{O}_i, \quad (3.18)$$

$$\mathbf{O}_{i''c} = \bar{\mathbf{O}}_{r'} = \frac{1}{n_{r'}} \sum_{i'=1}^{n_{r'}} \mathbf{O}_{i'}. \quad (3.19)$$

Metoda Warda

W metodzie tej odległości między dwoma grupami obiektów nie można przedstawić wprost za pomocą odległości między obiektami należącymi do tych grup. Dwie grupy obiektów podczas tworzenia drzewka połączeń, na dowolnym etapie są łączone w jedną grupę, w celu zminimalizowania sumy kwadratów odchyłeń wszystkich obiektów z tych dwóch grup od środka ciężkości nowej grupy, powstałej w wyniku połączenia tych dwóch grup. Proces ten oznacza, że na każdym etapie łączenia grup obiektów, w jedną grupę łączy się te grupy, które charakteryzują się najmniejszym zróżnicowaniem ze względu na opisujące je zmienne. Zróżnicowanie badania się przy pomocy kryterium $ESS(Errors\ of\ Squares)$ sformułowanego przez J.H. Warda, które jest postaci:

$$ESS = \sum_{i''=1}^{n_r''} d_{i''c}^2 (\mathbf{O}_{i''} \in \mathbf{G}_{r''}, \mathbf{O}_{i''c} = \bar{\mathbf{O}}_{r''} \in \mathbf{G}_{r''}), \quad (3.20)$$

gdzie: $d_{i''c}$ - odległość i'' -tego obiektu, należącego do nowo powstałej r'' -tej grupy od środka ciężkości tej grupy,

$$\mathbf{O}_{i''c} = \bar{\mathbf{O}}_{r''} = \frac{1}{n_{r''}} \sum_{i''=1}^{n_{r''}} \mathbf{O}_{i''}. \quad (3.21)$$

Rozdział 4

Szczegółowe opisy wybranych metod porządkowania

4.1 Metody porządkowania liniowego

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. W kolejnym etapie następuje dzielenie mierników odległości macierzy, na klasy podobieństwa obiektów. Kolejny krok polega na przyporządkowaniu poszczególnym klasom podobieństwa obiektów odpowiedniego symbolu graficznego. Samo porządkowanie obiektów polega na porządkowaniu diagramu, tj. przestawieniu wierszy i odpowiadających im kolumn diagramu, tak aby symbole graficzne reprezentujące najmniejsze odległości skupiały się wzdłuż głównej przekątnej, zaś w miarę oddalania się od głównej przekątnej znajdowały się symbole graficzne odpowiadające coraz to większym odległością.

Jako narzędzie pomocnicze w porządkowaniu danych, może stanowić kryterium postaci:

$$F^1 = \sum_{i=1}^n \sum_{i'=1}^n d_{ii'}, w_{ii'}$$

gdzie:

$d_{i,i'}$ - odległość euklidesowa między i -tym i i' -tym obiektem .

$w_{i,i'}$ - wagi elementów macierzy odległości, zdefiniowane w oparciu o jeden z następujących wzorów:

$$w_{i,i'} = \frac{|i-i'|}{n-1},$$

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

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

Dodatkowo wagi elementów macierzy odległości tworzą macierz wag postaci:

$$W = [w_{ii'}], i, i' = 1, 2, \dots, n.$$

Metoda rang

Na początku przeprowadzana jest stymulacja zmiennych, w kolejnym kroku dla każdego obiektu wyznacza się sumę przyporządkowanych mu rang ze względu na wszystkie zmienne. W chwili gdy dana wartość zmiennej występuje w więcej niż jednym obiekcie, następuje przyporządkowanie im jednakowej rangi będącej średnią arytmetyczną z przysługujących im rang. Na końcu zostaje obliczona wartość zmiennej syntetycznej, jako średnia wartość rang:

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

gdzie:

z_{ij} -zmienna znormalizowana rangowo, tj.

$$z_{ij} = h \text{ dla } x_{hj} = x_{ij}, h, i = 1, 2, \dots, n.$$

gdzie:

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

Metoda sum

Metoda ta bazuje na konstrukcji zmiennej syntetycznej przy pomiarze zmiennych na skali ilorazowej lub przedziałowej. W pierwszym etapie zostaje dokonana stymulacja zmiennych. Następnie obliczana jest wartość zmiennej syntetycznej dla każdego obiektu, jako średnia arytmetyczna z wartości zmiennych, przy przyjęciu jednakowych wag dla każdej zmiennej. Następnie eliminowane są ujemne wartości zmiennej syntetycznej, poprzez przesuwanie jej skali do punktu zerowego, przy użyciu przekształcenia:

$$s'_i = s_i - \min\{s_i\}, i = 1, 2, \dots, n$$

Końcowa postać zmiennej syntetycznej otrzymywana jest po przeprowadzeniu normalizacji według wzoru:

$$s''_i = \frac{s'_i}{\max\{s'_i\}}, i = 1, 2, \dots, n.$$

Powyższe przekształcenia powodują unormowanie miary syntetycznej w przedziale $[0,1]$.

Metody wzorcowe

W metodach tych zakłada się istnienie obiektu wzorcowego, w którym zmienne wejściowe przyjmują optymalne wartości, które to mogą być ustalane na podstawie ogólnie przyjętych norm, subiektywnej opinii dotyczącej obserwowanego obiektu, lub też opinii ekspertów. Poniżej zostaną omówione wybrane metody porządkowania, wzorcowe.

Metoda Hellwiga

W metodzie tej, w pierwszym etapie wyznacza się obiekt wzorcowy, na podstawie wystandaryzowanych zmiennych wejściowych. Współrzędnym obiektu wzorcowego przyporządkowuje się maksimum, gdy zmienne wejściowe są stymulantami, lub też minimum gdy zmienne wejściowe są destymulantami. Następnie dla każdego obiektu, następuje obliczenie jego odległości od obiektu wzorcowego, w tym celu najczęściej wykorzystywana jest metryka euklidesowa. Miara syntetyczna jest postaci:

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

gdzie:

współrzędne obiektu wzorcowego są obliczane na podstawie wzoru:

$$z_{0j} = \begin{cases} \max_i \{z_{ij}\} & \text{dla } z_j^S, j = 1, 2, \dots, m, i = 1, 2, \dots, n \\ \min_i \{z_{ij}\} & \text{dla } z_j^D, j = 1, 2, \dots, m, i = 1, 2, \dots, n \end{cases},$$

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

$$d_0 = \bar{d}_0 + 2S(d_0),$$

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

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

Wartości miary s_i zazwyczaj są z przedziału $[0; 1]^2$. Należy tu zaznaczyć, że wartości miary są tym wyższe, im mniej jest oddalony obiekt od obiektu wzorcowego.

Metoda Walesiaka

Metoda ta bazuje na konstrukcji zmiennej syntetycznej w oparciu o badanie odległości obiektów od obiektu wzorcowego, przy wykorzystaniu uogólnionej miary odległości. Umożliwia ona porządkowanie obiektów, jeżeli opisujące je charakterystyki są mierzone przynajmniej na skali porządkowej. W takim przypadku, zmienne wejściowe o postaci nominant muszą zostać podane stymulacji. Z kolei gdy zmienne są mierzone na skali przedziałowej lub ilorazowej, należy je znormalizować. Miara syntetyczna oparta na uogólnionej mierze odległości przyjmuje postać:

$$s_i = \frac{1}{2} - \frac{\sum_{j=1}^m w_j a_{i0j} b_{0ij} + \sum_{j=1}^m \sum_{i''=1}^n w_j a_{ii''j} b_{0i''j}}{2 \left[\sum_{j=1}^m \sum_{i''=1}^n w_j a_{ii''j}^2 \cdot \sum_{j=1}^m \sum_{i''=1}^n w_j b_{0i''j}^2 \right]^{\frac{1}{2}}} \quad (4.1)$$

gdzie:

$$w_j \in [0; m]$$

$$\sum_{j=1}^m w_j = m$$

Ostateczna postać zmiennej syntetycznej zależy od skali pomiaru zmiennych.

Jeśli zmienne charakteryzujące obiekty mierzone są na skali ilorazowej lub przedziałowej, stosowane jest następujące podstawienie:

$$a_{ii^*j} = z_{ij} - z_{i^*j} \text{ dla } i^* = 0, i'' \quad b_{0i^*j} = z_{0j} - z_{i^*j} \text{ dla } i^* = i, i''.$$

gdzie:

z_{0j} -wystandaryzowana wartość j -tej zmiennej dla obiektu wzorcowego

Z kolei, gdy zmienne charakteryzujące obiekty mierzone są na skali porządkowej to stosowane jest podstawienie:

$$a_{ii^*j} = \begin{cases} 1 & \text{dla } z_{ij} > z_{i^*j}, \\ 0 & \text{dla } z_{ij} = z_{i^*j}, i^* = 0, i' \\ -1 & \text{dla } z_{ij} < z_{i^*j} \end{cases} \quad (4.2)$$

$$b_{0i^*j} = \begin{cases} 1 & \text{dla } z_{0j} > z_{i^*j}, i^* = i, i' \\ 0 & \text{dla } z_{0j} = z_{i^*j}, i^* = i, i' \\ -1 & \text{dla } z_{0j} < z_{i^*j}, \end{cases} \quad (4.3)$$

Zmienna syntetyczna przyjmuje wartości z przedziału $[0; 1]$. Czym niższa wartość zmiennej syntetycznej, tym bliżej wzorca leży dany obiekt.

Metoda dystansowa

Podobnie jak wcześniej opisane metody, w pierwszym kroku należy wyznaczyć zmienną syntetyczną w oparciu o jej odległość od obiektu wzorca, dla każdego z porównywanych obiektów, przy wykorzystaniu np. metryki euklidesowej. Miara syntetyczna, przy wykorzystaniu przekształcenia unitaryzacyjnego, jest postaci:

$$s_i = \left(\frac{d_{i0} - \min_i \{d_{i0}\}}{\max_i \{d_{i0}\} - \min_i \{d_{i0}\}} \right)^p, i = 1, 2, \dots, m \quad (4.4)$$

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

4.1.1 Metody iteracyjne

W metodach tych przyjmowania jest funkcja kryterium dobroci porządkowania i 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

W metodzie tej takie liniowe uporządkowanie obiektów, dla którego funkcja kryterium dobroci uporządkowania osiąga maksimum:

$$F^2 = \sum_{i'=1}^{n-1} i' \sum_{i=1}^{n-i'} d_{ii'} \rightarrow \max \quad (4.5)$$

gdzie:

$d_{ii'}$ - odległość euklidesowa między i -tym i i' -tym obiektem.

Sposób postępowania:

W pierwszym kroku dokonywane jest dowolne liniowe uporządkowanie obiektów, dla którego obliczana jest wartość funkcji kryterium (3.5). W kolejnym etapie obliczana jest wartość funkcji kryterium dla każdej możliwej transpozycji pary obiektów. Jeżeli wartość funkcji kryterium dla każdej z transpozycji par obiektów są mniejsze od wartości tej funkcji dla uporządkowania wyjściowego obiektów, uporządkowanie to uważane jest za najlepsze. W przeciwnym wypadku, dokonywana jest transpozycja tej pary obiektów, dla której wzrost wartości funkcji kryterium jest największy.

Uporządkowanie to stanowi punkt wyjścia do oceny, czy kolejna transpozycja dowolnej pary obiektów pozwoli na wzrost wartości funkcji kryterium. Powyższe postępowanie kontynuowane jest tak długo, aż transpozycja dowolnej pary obiektów nie prowadzi do wzrostu wartości funkcji kryterium.

4.1.2 Metody gradientowe

W metodach gradientowych dąży się do takiego liniowego uporządkowania obiektów, które jak najmniej zniekształca relacje strukturalne porządkowanego zbioru obiektów. Od strony geometrycznej oznacza to, że odległości pomiędzy punktami reprezentującymi obiekty w przestrzeni jednowymiarowej, określonej przez zmienną syntetyczną, w jak najmniejszym stopniu zniekształcają odległości pomiędzy tymi punktami w przestrzeni wielowymiarowej, określonej przez zmienne wejściowe. Metody gradientowe poszukują takich współrzędnych punktów reprezentujących obiekty w przestrzeni jednowymiarowej, dla których funkcja dobroci uporządkowania osiąga minimum, co można przedstawić za pomocą wariantów:

$$F^3 = \frac{\sum_{\substack{i,i'=1 \\ i \neq i'}}^n (d_{ii'}^s - d_{ii'})^2}{\sum_{\substack{i,i'=1 \\ i < i'}}^n d_{ii'}} \rightarrow \min \quad (4.6)$$

$$F^4 = \sum_{\substack{i,i'=1 \\ i < i'}}^n \left(\frac{d_{ii'}^s - d_{ii'}}{d_{ii'}} \right)^2 \rightarrow \min \quad (4.7)$$

$$F^5 = \frac{1}{\sum_{\substack{i,i'=1 \\ i \neq i'}}^n d_{ii'}} \sum_{\substack{i,i'=1 \\ i < i'}}^n \frac{\left(d_{ii'}^s - d_{ii'} \right)^2}{d_{ii'}} \rightarrow \min \quad (4.8)$$

gdzie:

$d_{ii'}$ - odległość euklidesowa między i -tym i i' -tym obiektem.

Punkt wyjścia procedury, polega na dowolnym liniowym uporządkowaniu obiektów, następnie

w trakcie kolejnych iteracji poszukiwane jest ekstremum funkcji wielu zmiennych, zapewniające jak największy spadek wartości funkcji kryterium. Poniżej zostały szczegółowo omówione kroki postępowania:

Na początku wyznaczana jest wartość funkcji kryterium dla wyjściowego, liniowego uporządkowania obiektów (wyjściowych wartości zmiennych syntetycznych w tych obiektach), przy czym wartość funkcji jest przyjmowana jako wynik iteracji dla $t = 0$:

$$F^5 = \frac{1}{c} \sum_{\substack{i,i'=1 \\ i < i'}}^n \frac{\left(d_{ii',t}^s - d_{ii'}\right)^2}{d_{ii'}} \quad (4.9)$$

gdzie:

$$c = \frac{1}{\sum_{\substack{i,i'=1 \\ i < i'}}^n d_{ii'}}, \quad (4.10)$$

Przy czym wartości zmiennych oryginalnych oraz wyjściowych wartości zmiennych syntetycznych, zostały znormalizowane na przedziale $[0; 1]$.

Współrzędne zmiennych syntetycznych dla obiektów w kolejnej iteracji $t + 1$ wyznacza się na podstawie wzoru:

$$s_{i,t+1} = s_{i,t} - W \Delta_i(t), \quad (4.11)$$

gdzie:

$$\Delta_i(t) = \frac{\delta F_t^5}{\delta s_{i,t}} : \frac{\delta F_t^5}{(\delta s_{i,t})^2}, \quad (4.12)$$

przy czym:

$$\frac{\delta F^5}{\delta s_i} = -\frac{2}{c} \sum_{\substack{i'=1 \\ i \neq i'}}^n \left(\frac{d_{ii'} - d_{ii'}^s}{d_{ii'} + d_{ii'}^s} \right) (s_i - s_{i'}), \quad (4.13)$$

$$\frac{\delta^2 F^5}{(\delta s_i)^2} = -\frac{2}{c} \sum_{\substack{i \neq i' \\ i'=1}}^n \frac{1}{d_{ii'} d_{ii'}^s} \left[(d_{ii'} - d_{ii'}^s) - \frac{(s_i - s_{i'})^2}{d_{ii'}^s} \left(1 + \frac{d_{ii'} - d_{ii'}^s}{d_{ii'}^s} \right) \right]. \quad (4.14)$$

W - parametr.

Na wstępie zakłada się maksymalną oraz minimalną wartość parametru W , wskaźnik skali zmian wartości tego parametru pomiędzy iteracjami W_{t+1}/W_t oraz maksymalną liczbę iteracji. Procedurę iteracyjną rozpoczynamy od przyjętej maksymalnej wartości parametru W . Postępowanie iteracyjne jest kontynuowane do momentu gdy nastąpi wzrost wartości wzrost wartości funkcji kryterium. Po tym następuje powrót do wartości zmiennej syntetycznej z poprzedniej iteracji, przy jednoczesnym zmniejszeniu wartości parametru W o przyjęty wskaźnik jego zmian. Procedura kontynuowana jest do momentu, aż wartość parametru W , nie spadnie poniżej założonej wartości minimalnej lub gdy osiągnie z góry założoną liczbę iteracji.

4.2 Metody porządkowania nieliniowego

Porządkowanie nieliniowe polega, od strony geometrycznej, na rzutowaniu na płaszczyznę obiektów umieszczonych w wielowymiarowej przestrzeni zmiennych. Ta metoda porządkowania nie pozwala na ustaleniu hierarchii obiektów, lecz na określeniu dla każdego z nich, stopnia podobieństwa do innych obiektów, ze względu na ich charakterystyki.

Aby uporządkować nieliniowo obiekty, charakteryzujące je zmienne powinny być mierzone przynajmniej na skali przedziałowej lub ilorazowej. Gdy zmienne te mierzone są na skali przedziałowej lub ilorazowej, należy dokonać ich normalizacji, dla zapewnienia ich porównywalności.

Metody porządkowania nieliniowego można podzielić na metody dendrytowe i metody aglomeracyjne. Metody dendrytowe prowadzą do powstania dendrytu, będącego ilustracją graficzną położenia względem siebie obiektów ze względu na ich podobieństwo. Z kolei metody aglomeracyjne prowadzą do utworzenia drzewka połączeń, będącego graficzną ilustracją hierarchii łączenia obiektów, ze względu na ich podobieństwo.

4.2.1 Metody dendrytowe

Metody dendrytowe opierają się na regułach i pojęciach teorii grafów. Porządkowanie dendrytowe polega na przyporządkowaniu obiektom poszczególnych wierzchołków dendrytu, w tym celu budowany jest dendryt. Poniżej zostaną opisane przykłady metod dendrytowych, tj. taksonomia wrocławska oraz metoda Prima.

Taksonomia wrocławska

W pierwszym etapie tej metody, dla każdego obiektu O_i poszukiwany jest obiekt $O_{i'}$, który jest najbardziej do niego podobny. W tym celu w każdym wierszu(kolumnie) macierzy odległości D , wyznaczamy jest element najmniejszy:

$$d_{ii'} = \max_i d_{ii'}, i, i' = 1, 2, \dots, n; i \neq i'. \quad (4.15)$$

Następnie otrzymane pary najbardziej podobnych do siebie obiektów, przedstawiane są w postaci grafu nieorientowanego, Długość krawędzi łączących wierzchołki grafu, są proporcjonalne do odległości między obiektami. Może się zdarzyć, że wśród wyznaczonych par połączeń, pojawią się połączenia występujące dwukrotnie, jedno z nich zostanie wyeliminowane, ponieważ kolejność połączeń w dendrycie nie jest istotne. Warto również zwrócić uwagę, na fakt iż w dendrycie danych obiekt może występować tylko jeden raz, w związku z tym jeżeli w łączeniu występują wielokrotnie te same obiekty, to zostaną one połączone w zespoły zwane skupieniami.

W kolejnym kroku, sprawdza się, czy utworzony graf jest spójny. Jeżeli tak będzie, to algorytm zostaje zakończony. W przeciwnym wypadku, poszczególne składowe dendrytu łączy się w większe zespoły. Odpowiednie skupienia łączone są ze sobą w miejscach, określonych przez minimalną odległość między nimi. Tworzone są w ten sposób skupienia 2-go rzędu. W tym celu znajdowana jest najmniejsza odległość każdego obiektu jednego skupienia, od obiektów należących do pozostałych skupień. Z uzyskanych odległości wybierana jest odległość najmniejsza, która zostaje wiązadłem łączącym skupienia.

Powyższy proces przeprowadzany jest do momentu, aż nie powstanie graf spójny, w ten sposób tworzone są skupienia wyższego rzędu

Metoda Prima

W odróżnieniu od taksonomii wrocławskiej, metoda Prima nie wymaga operowania całym czasem pełną, wyjściową macierzą odległości. W trakcie tworzenia dendrytu, na każdym etapie zbiór porządkowanych obiektów jest klasyfikowany do jednego z dwóch podzbiorów, np. A i B . Niech zbiór A będzie pierwszym z nich a zbiór B drugim. Pierwszy z nich zawiera obiekty należące na danym etapie do dendrytu, zaś drugi zawiera obiekty nie należące na tym etapie do dendrytu.

Sposób postępowania:

Na początku procedury, zbiór A jest zbiorem pustym, z kolei zbiór B zawiera wszystkie obiekty.

W pierwszym kroku do zbioru A zostaje włączony dowolny obiekt, nie ma to wpływu na ostateczną postać dendrytu. Następnie do zbioru A zostają włączone te obiekty zbioru B , najbardziej podobne do obiektów należących już do zbioru A . Proces ten trwa do momentu, aż zbiór B nie będzie pusty. W tym celu w pierwszym kroku algorytmu zostaje stworzony wektor d , zawierający odległości wybranego obiektu zbioru A , od pozostałych obiektów zbioru B . Po utworzeniu wektora, sprawdzane jest dla którego elementu odległość od elementu ze zbioru A , jest najmniejsza. Po znalezieniu tego elementu, zostaje on włączony do zbioru A i usunięty ze zbioru B . Po tym etapie, sprawdzane jest czy zbiór B jest pusty, jeżeli tak jest to algorytm kończy działanie, zaś w przeciwnym wypadku, zostaje ponownie tworzony wektor d , którego elementami są najmniejsze z odległości każdego z obiektów pozostających jeszcze w zbiorze B od obiektów, które należą do zbioru A . Ponownie wybierany jest z wektora d najmniejszy element i włączany do zbioru A przy jednoczesnym usunięciu ze zbioru B .

W powstałym dendrycie wierzchołkami są obiekty przechodzące kolejno do zbioru A , z kolei wiązkami łączącymi wierzchołki są minimalne wartości elementów wektora d , otrzymanego w kolejnych krokach przyłączania obiektów do dendrytu.

4.2.2 Metody aglomeracyjne

Istotą metod aglomeracyjnych jest utworzenie drzewka połączeń - dendrogramu. W ten sposób zobrazowana jest hierarchia łączenia obiektów, na podstawie zmniejszającego się podobieństwa między obiektami włączonymi do dendrogramu, w kolejnych etapach a obiektami należącymi już do dendrogramu. Hierarchia połączeń określa wzajemnie położenie względem siebie obiektów oraz grup obiektów powstających w kolejnych etapach tworzenia drzewka. Grupy podobnych do siebie obiektów tworzą oddzielne gałęzie. Punktem wyjściem metod aglomeracyjnych stanowi założenie, że każdy obiekt stanowi odrębną, jednoelementową grupę ($G_r, r = 1, 2, \dots, z$).. W kolejnych krokach następuje łączenie ze sobą grupy obiektów najbardziej do siebie podobnych ze względu na wartości opisujących je zmiennych. Podobieństwo weryfikowane jest na podstawie odległości między grupami obiektów. Poniżej zostały szczegółowo omówione kroki postępowania:

Na początku odległości między jednoelementowymi grupami obiektów G_1, \dots, G_z są elementami wyjściowej macierzy odległości D . W macierzy D poszukiwane są najmniejsze odległości pomiędzy grupami obiektów:

$$d_{rr'} = \min_{ii'} d_{ii'}, i = 1, 2, \dots, n_r; i' = 1, 2, \dots, n_{r'}; r, r' = 1, 2, \dots, z; r \neq r'. \quad (4.16)$$

gdzie:

$d_{rr'}$ - odległość r -tej od r' -tej grupy. W kolejnym kroku następuje łączenie do siebie obiektów podobnych w jedną grupę, w wyniku czego wyjściowa liczba grup zmniejszona jest o jeden, oraz rozpoczęta jest budowa drzewka połączeń. Następnie wyznacza się ponownie odległość nowo utworzonej grupy obiektów od wszystkich pozostałych grup obiektów. Odległości te umieszczone zostają w macierzy odległości D - w miejscu wierszy i kolumn odpowiadających obiektom (grupom obiektów) połączonych w jedną grupę. Po każdym etapie grupowania ponownie określana jest odległość między nowo powstałą grupą a pozostałymi grupami. Warto również dodać, że odległości te tworzą nową, aktualną na danym etapie grupowania, macierz odległości o co raz mniejszym wymiarze $(n - u)(n - u)$, gdzie u jest u -tym etapem łączenia grup obiektów. Procedura łączenia grup obiektów powtarzana jest tak długo, aż nie zostanie utworzona jedna grupa, tj. zostanie utworzone pełne drzewko połączeń.

Ogólny wzór wyznaczania odległości nowo powstałej grupy $G_{r''}$, powstałej w wyniku połączenia grup G_r i $G_{r'}$, od pozostałych grup $G_{r'''}'$, przy tworzeniu drzewka połączeń ma postać:

$$d_{r''r'''} = \alpha_r d_{rr'''} + \alpha_{r'} d_{r'r'''} + \beta d_{rr'} + \gamma |d_{rr'''} - d_{r'r'''}| \quad (4.17)$$

gdzie:

$\alpha_r, \alpha_{r'}, \beta, \gamma$ - współczynniki przekształceń odmienne dla różnych metod aglomeracyjnych

Poszczególne metody aglomeracyjne, różnią się między sobą sposobami wyznaczania odległości między obiektami. Poniżej zostały wymienione najczęściej stosowane metody aglomeracyjne, które będą dokładniej omówione w kolejnym podrozdziale.

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

Metoda najbliższego sąsiedztwa

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

$$d_{rr'} = \min_{ii'} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}), \quad (4.18)$$

$$i = 1, 2, \dots, n_r; i' = 1, 2, \dots, n_{r'}; r, r' = 1, 2, \dots, z; r \neq r',$$

gdzie:

$$\mathbf{O}_i = [z_{ij}], j = 1, 2, \dots, m. \quad (4.19)$$

Metoda najdalszego sąsiedztwa

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

$$d_{rr'} = \max_{ii'} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}), \quad (4.20)$$

$$i = 1, 2, \dots, n_r; i' = 1, 2, \dots, n_{r'}; r, r' = 1, 2, \dots, z; r \neq r',$$

Metoda średniej międzygrupowej

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

$$d_{rr'} = \frac{1}{n_r n_{r'}} \sum_{i'=1}^{n_{r'}} \sum_{i=1}^{n_r} d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'}) \quad (4.21)$$

$$r, r' = 1, 2, \dots, z; r \neq r'.$$

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:

$$d_{rr'} = \text{med}_{i,i'} \{d_{ii'}(\mathbf{O}_i \in \mathbf{G}_r, \mathbf{O}_{i'} \in \mathbf{G}_{r'})\}, \quad (4.22)$$

$$i = 1, 2, \dots, n_r; i' = 1, 2, \dots, n_{r'}; r, r' = 1, 2, \dots, z; r \neq r'.$$

Metoda środków ciężkości

W metodzie tej odległość między dwoma grupami jest równa odległości między środkami ciężkości tych grup. Odległość ta opisana jest wzorem:

$$d_{rr'} = d_{i_c i'_c}(\mathbf{O}_i^c = \bar{\mathbf{O}}_r \in \mathbf{G}_r, \mathbf{O}_{i'}^c = \bar{\mathbf{O}}_{r'} \in \mathbf{G}_{r'}), \quad (4.23)$$

$$i = 1, 2, \dots, n_r; i' = 1, 2, \dots, n_{r'}; r, r' = 1, 2, \dots, z; r \neq r'.$$

gdzie:

$d_{i_c i'_c}$ - odległość środka ciężkości r -tej grupy od środka ciężkości r' -tej grupy,
 $\bar{\mathbf{O}}_{i_c}, \bar{\mathbf{O}}_{i'_c}$ - środki ciężkości odpowiednio r -tej i r' -tej grupy obiektów. przy czym:

$$\mathbf{O}_{i_c} = \bar{\mathbf{O}}_r = \frac{1}{n_r} \sum_{i=1}^{n_r} \mathbf{O}_i \quad (4.24)$$

$$\mathbf{O}_{i'_c} = \bar{\mathbf{O}}_{r'} = \frac{1}{n_{r'}} \sum_{i'=1}^{n_{r'}} \mathbf{O}_{i'}. \quad (4.25)$$

Metoda Warda

W metodzie tej odległości między dwoma grupami obiektów nie można przedstawić wprost za pomocą odległości między obiektami należącymi do tych grup. Dwie grupy obiektów podczas tworzenia drzewka połączeń, na dowolnym etapie są łączone w jedną grupę, w celu zminimalizowania sumy kwadratów odchylen wszystkich obiektów z tych dwóch grup od środka ciężkości nowej grupy, powstałej w wyniku połączenia tych dwóch grup. Proces ten oznacza, że na każdym etapie łączenia grup obiektów, w jedną grupę łączy się te grupy, które charakteryzują się najmniejszym zróżnicowaniem ze względu na opisujące je zmienne. Zróżnicowanie badania się przy pomocy kryterium $ESS(ErrrosSumofSquares)$ sformułowanego przez J.H. Warda, które jest postaci:

$$ESS = \sum_{i''=1}^{n_r''} d_{i'' i''_c}^2(\mathbf{O}_{i''} \in \mathbf{G}_{r''}, \mathbf{O}_{i''_c} = \bar{\mathbf{O}}_{r''} \in \mathbf{G}_{r''}), \quad (4.26)$$

gdzie: $d_{i'' i''_c}$ - odległość i'' -tego obiektu, należącego do nowo powstałej r'' -tej grupy od środka ciężkości tej grupy,

$$\mathbf{O}_{i''_c} = \bar{\mathbf{O}}_{r''} = \frac{1}{n_{r''}} \sum_{i''=1}^{n_r''} \mathbf{O}_{i''}. \quad (4.27)$$

Rozdział 5

Zbiór danych

5.1 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 biegów, 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.

| | 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.BRI | MOC[km] | POJEMN | ROK | PRO | PRZEBIEG | KOLOR | L.DZ | RZWI | RODZAJ | F.SKRZYNIA | NAPED | KRAJ | AKT |
| 2 | Hyundai | i20 | II | kompakt | małopolskie | | 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-m | | 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 | Volkswage | 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 | małopolskie | | 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 | wielkopolsl | 95900 | 117957 | 300 | 3400 | 2006 | | 83000 | srebrny | | 3 | benzyna | automatyc | na tylne ko | 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 | małopolskie | | 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 | małopolskie | | 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 5.1: 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:

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

Wśród zmiennych jakościowych można wyróżnić zmienne porządkowe, nominalne oraz binarne. W stworzonym zbiorze danych zmiennymi binarnymi są atrybuty:

- pierwszy.wlasciciel,
- ABS,
- serwisowany,
- komputer.pokladowy,
- ESP,
- bezwypadkowy,
- uszkodzony.

Pozostałe atrybuty są zmiennymi nominalnymi.

Zmiennymi ilościowymi są atrybuty:

- cena.netto[pln],
- cena.brutto[pln],
- moc,
- pojemnosc.skokowa[cm³],
- rok.produkcji,
- przebieg[km],
- l.drzwi.

Wśród zmiennych ilościowych można wyróżnić zmienne skokowe oraz dyskretne. W stworzonym zbiorze danych, zmiennymi skokowymi są:

- moc,
- pojemnosc.skokowa[cm³],
- rok.produkcji,
- przebieg,
- l.drzwi.

Z kolei atrybuty: cena.netto[pln], cena.brutto[pln] są zmiennymi ciągłymi.

Bibliografia

- [1] Grzegorz Banaszak and 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 and Sławomir Turek. *Teoria mnogości*. Państwowe Wydawnictwo Naukowe, Warszawa, 2007.
- [4] Patric Billingsley. *Prawdopodobieństwo i miara*. Państwowe Wydawnictwo Naukowe, Warszawa, 1987.
- [5] Jerzy Greń. *Statystyka matematyczna: modele i zadania*. Państwowe Wydawnictwo Naukowe, Warszawa, 1984.
- [6] Jacek Jakubowski and Rafał Sztencel. *Wstęp do teorii prawdopodobieństwa*. SCRIPT, Warszawa, 2004.
- [7] W. Krysiński, J. Bartos, W. Dyczka, K. Królikowska, and M. Wasilewski. *Rachunek prawdopodobieństwa i statystyka matematyczna w zadaniach: część I. rachunek prawdopodobieństwa*. Państwowe Wydawnictwo Naukowe, Warszawa, 1999.
- [8] Kazimierz Kuratowski. *Wstęp do teorii mnogości i topologii*. Państwowe Wydawnictwo Naukowe, Warszawa, 2004.
- [9] Andrzej Młodak. *Analiza taksonomiczna w statystyce regionalnej*. Centrum Doradztwa i Informacji Difin, Warszawa, 2006.
- [10] Tomasz Panek and Jan Karol Zwierchowski. *Statystyczne metody wielowymiarowej analizy porównawczej: teoria i zastosowania*. Oficyna Wydawnicza, Szkoła Główna Handlowa, Warszawa, 2013.
- [11] Ryszard Rudnicki. *Wykłady z analizy matematycznej*. Państwowe Wydawnictwo Naukowe, Warszawa, 2006.
- [12] Robin J. Wilson. *Wprowadzenie do teorii grafów*. Państwowe Wydawnictwo Naukowe, Warszawa, 2008.