

Eksploracja reguł ilościowych

Jacek Sosnowski

1 grudnia 2014

Politechnika Warszawska
Wydział Elektroniki i Technik Informacyjnych
Instytut Informatyki

Spis treści

1	Wstęp	2
2	Eksploracja reguł asocjacyjnych	3
2.1	Rodzaje atrybutów	4
2.1.1	Atrybuty kategoryczne	4
2.1.2	Atrybuty numeryczne	5
2.2	Wsparcie i ufność	5
3	Ilościowe reguły asocjacyjne	6
3.1	Przykład	8
4	Analiza skupień	9
4.1	Miara podobieństwa	10
4.2	Normalizacja	11
4.3	Rodzaje analizy skupień	12
4.3.1	Właściwości grup wynikowych	12
4.3.2	Metody podziału	14
4.4	Rozważania na temat właściwości grupowania	17
5	Zbiór wartości binarnych	18
5.1	Transformacja na wartości binarne	19
5.2	Dyskretyzacja	19
5.3	Właściwości dyskretyzacji	20
5.4	Przykład transformacji	20
6	Analiza skupień, a wyszukiwanie reguł asocjacyjnych	20
6.1	W kontekście wyszukiwania zbiorów częstych	22
6.2	W kontekście powoływania reguł asocjacyjnych	22
7	Quality Threshold Clustering	23
7.1	Właściwości	23
7.2	Parametr	25
7.3	Algorytm QT w kontekście reguł ilościowych	26
8	Problem skrawków	28
8.1	Rozciąganie przedziałów	29
8.2	Nakładanie przedziałów	30

9	Koncepcja rozwiązania	31
9.1	Zmodyfikowany algorytm Quality Threshold	31
9.2	Algorytm MQTC	31
9.3	Rozwiązanie problemów	33
9.4	Właściwości	34
10	Eksperymenty	35
10.1	Kryterium porównawcze	36
10.2	MQTC w porównaniu z Equi-depth	38
10.3	MQTC w porównaniu z K-means	40
11	Podsumowanie	41

1 Wstęp

Otoczająca wszystkich rzeczywistość coraz częściej przenika się ze światem wirtualnym. Innymi słowy, Internet oraz komputery i narzędzia z nimi związane, już dawno stały się niezbędnym elementem codzienności. Próby przewidzenia dalszego rozwoju obecnej sytuacji spędzają sen z powiek naukowcom i wizjonerom. Niemniej jednak już teraz obserwuje się ciekawy skutek ekspansji technik komputerowych. Jest nim powstawanie społeczeństw informacyjnych. Takim mianem określa się cywilizację, w której szczególnym dobrem ekonomicznym staje się informacja. To ona pośrednio zaczyna być podstawą dochodu narodowego, a przez to głównym źródłem utrzymania wielu jednostek. Obecnie powszechnie akceptowalne jest stwierdzenie, że informacja jest bardzo cenna. Z tego powodu następuje ciągły rozwój usług informatycznych. Według autorów [16] z socjoekonomicznego punktu widzenia, sektor informacji dzieli się na: produkcję, przetwarzanie oraz przemysł dystrybucji informacji. Ostatni z nich obejmuje teleinformatykę, a więc przesyłanie danych. Natomiast produkcja odbywa się często w sposób naturalny. Kupując towary czy korzystając z dostępnych usług, konsumenci są producentami informacji. W tej pracy najwięcej odniesień będzie do przetwarzania, które jest naturalną konsekwencją gromadzenia danych.

Wyniki najnowszych badań statystycznych mogą wskazywać na prawdziwość doniesień o rozwoju wspomnianego typu społeczeństwa nawet w pięknym kraju nad Wisłą. Według danych GUS [17] dla Polski *„w 2012 r. liczba firm z sektora ICT wzrosła w stosunku do 2009 r. o 25,6 % (...) natomiast liczba pracujących w tym sektorze – o 11,6 %”*. Dodatkowo w tym samym okresie przychody netto dla tego sektora zwiększyły się o 30,8 %. Oznacza to wzrost zainteresowania świadczonymi usługami.

Efektem ubocznym dojrzewania społeczeństwa informatycznego jest wytwarzanie ogromnej ilości danych. Im więcej aspektów codziennego życia jest wspomagane przez technologie, tym więcej danych można zebrać. Przykładowo, jeszcze kilka dekad temu zwykłe zakupy były prostą wymianą dóbr (z wykorzystaniem środka płatniczego). Dzisiaj każda taka transakcja jest rejestrowana w systemie komputerowym, a następnie zapisywana w przygotowanej bazie danych. Jest to osiągalne dzięki postępowi techniki. A przede wszystkim dzięki łatwo dostępnej i taniej pamięci dyskowej. Realne jest więc gromadzenie dużych ilości danych, co czyni większość dzisiejszych instytucji.

Z upływem czasu coraz tańsza staje się również moc obliczeniowa. Taki kierunek zmian wprost świetnie wpisuje się w potrzeby współczesnej cywilizacji. Gdyż daje możliwość efektywnej analizy potężnych baz. Zadaniem tym zajmuje się dziedzina eksploracji danych. Dzięki jej metodom możliwe jest odkrycie interesujących zależności oraz nieznanej struktury zawartej w

zebranych danych. Można też pokusić się o stwierdzenie, że ta gałąź informatyki pozwala na wydobywanie „wiedzy” ukrytej w posiadanych bazach.

Najbardziej popularną metodą eksploracji danych nieustannie jest wyszukiwanie reguł asocjacyjnych. Istnieje wiele algorytmów rozwiązujących to zadanie. Niestety większość z nich wymaga, by dane wejściowe były zorganizowane w postaci tabeli o wartościach wyłącznie binarnych. Natomiast współcześnie gromadzone dane zdecydowanie częściej mają bogatsze dziedziny wartości. Często spotykane są zbiory, które zawierają liczby rzeczywiste. W takim przypadku unikalnych wartości, zamiast dwóch może być nawet nieskończenie wiele. O takie właśnie zbiory opierają się ilościowe reguły asocjacyjne. Przykładem może być reguła:

$$waga : [80, 100] \wedge gorne_cisnienie : [120, 140] \Rightarrow ryzyko = 1$$

Powyższa reguła jest oczywiście sztucznym przykładem, ale wiedza zawarta w takich implikacjach może być bezcenna. Dlatego celem niniejszej pracy jest zbadanie możliwości eksploracji ilościowych reguł asocjacyjnych. W tym celu połączone zostaną siły dwóch dyscyplin eksploracji danych: analizy skupień oraz odkrywania reguł asocjacyjnych. Cała praca nie aspiruje do miana przeglądu tych dziedzin. Przedstawia natomiast specyficzny tok zdobywania wiedzy i technologii dla rozwiązania problemu przedstawionego jako cel pracy. Dokumentuje tylko i wyłącznie te fragmenty teorii eksploracji danych, które aktywnie i bezpośrednio wpłynęły na rozwiązanie końcowe. Wszelkie nawiązania, które nie wyczerpują w pełni tematyki, ale dla kompletności rozważań znajdują się w tekście, zostały opatrzone komentarzem wskazującym dokładniejsze źródło.

2 Eksploracja reguł asocjacyjnych

Eksploracja reguł asocjacyjnych (ang. association rule mining) jest jednym z ważniejszych filarów eksploracji danych (ang. data mining). Głównym celem tej dziedziny informatyki jest wydobywanie z dużych zbiorów danych informacji wyższego poziomu abstrakcji. Oznacza to odkrywanie relacji czy struktur zawartych w analizowanym zestawie, najlepiej jeśli będą nieznane oraz interesujące. Manualna analiza nawet niewielkich baz danych może sprawiać problemy, a wraz ze wzrostem ich wymiarów, takie badanie staje się praktycznie niemożliwe dla człowieka. Dodatkowo najbardziej interesujące relacje ujawniają się dopiero w liczniejszych kolekcjach danych. Dobrym wprowadzeniem do wspomnianej tematyki jest słynna już praca [3]. Jednak dla kompletności rozważań zostanie przedstawiony tu bardzo krótki zarys teoretyczny problemu eksploracji danych.

Baza danych D to zbiór transakcji (próbek) $D = \{T_1, T_2, T_3, \dots, T_n\}$. Każda próbka T_j ($j \in [1, n]$) zawiera ustaloną liczbę składowych opisanych zbiorem atrybutów $I = \{i_1, i_2, i_3, \dots, i_k\}$.

Odkrycie związków pomiędzy atrybutami umożliwia eksploracja reguł asocjacyjnych. Ogólnie przez pojęcie reguły rozumie się wyrażenie postaci:

$$X \Rightarrow Y$$

$$X \in I, Y \in I, X \cap Y = \emptyset$$

Gdzie X oraz Y to zdarzenia. Przy czym kiedy w bazie występuje X to z pewnym poziomem ufności wystąpi też Y . Ten poziom nazywany jest wiarygodnością, lub ufnością (ang. confidence). Biorąc pod uwagę definicję bazy, każde ze zdarzeń, zarówno X jak i Y to podzbiory atrybutów I . Oznacza to, że reguła asocjacyjna opisuje relację pomiędzy występowaniem atrybutów X , a pojawianiem się atrybutów Y .

2.1 Rodzaje atrybutów

Mówiąc potocznie, bazę danych można potraktować jako tabelę. Wtedy każdy wiersz to rekord, próbka, lub po prostu element bazy. Natomiast każda kolumna nazywana jest cechą, atrybutem albo własnością obiektów przechowywanych w bazie. Każda cecha może zawierać wartości różnego typu, które można podzielić na kilka klas:

2.1.1 Atrybuty kategoryczne

Atrybuty kategoryczne, dzieli się na *nominalne* oraz *porządkowe*¹:

1. **Atrybuty nominalne** (dyskretne, skończone, wyliczeniowe) – wartości tworzą przestrzeń skończoną i zazwyczaj niewielką. Nie istnieje porządek pomiędzy wartościami atrybutu.
2. **Atrybuty binarne** – są szczególnym przypadkiem atrybutów nominalnych, ponieważ przyjmują tylko dwie wartości. Zazwyczaj jest to 0 i 1. Najczęściej (choć nie zawsze) oznaczają występowanie bądź brak danego atrybutu w transakcji. W [3] w rozdziale „Formal Model” można znaleźć formalną definicję atrybutów binarnych w kontekście bazy danych.
3. **Atrybuty porządkowe** (ang. ordinal attributes) – podobnie jak dla nominalnych, dziedzina jest skończona oraz niezbyt liczna, ale istnieje

¹Informacje zaczerpnięte między innymi z [7]

możliwość uporządkowania wartości tych atrybutów. Znajomość porządku nie implikuje jednak znajomości odległości pomiędzy wartościami. Przykładem może być następująca przestrzeń bardzo mało, niewiele, dużo, bardzo dużo. Bez dodatkowych informacji ilościowych opisujących poszczególne dostępne wartości nie można określić relacji odległości pomiędzy próbkami tego atrybutu.

2.1.2 Atrybuty numeryczne

Atrybuty ciągłe (*numeryczne*) – definicja przestrzeni wartości dla takich atrybutów oparta jest na dobrze określonym zbiorze liczbowym, np. liczby rzeczywiste, czy całkowite. Z tego powodu ilość tych wartości jest nieskończona (nie można przewidzieć ile różnych próbek znajduje się akurat w konkretnej bazie danych). Dodatkowo znany jest porządek pomiędzy wartościami oraz zazwyczaj dystans między nimi (według wybranej miary odległości). Takie atrybuty najczęściej są spotykane przy opisach właściwości fizycznych obiektów, czego przykładami mogą być: masa, temperatura, długość, wzrost, itp. Już ta krótka lista ukazuje, że pomimo teoretycznie nieskończonej liczby wartości jaką prezentują atrybuty numeryczne, czasami istnieje możliwość uściślenia dziedziny. Typowym wzorem jest wiek, który z punktu widzenia bazy danych może przyjmować dowolne dodatnie wartości, lecz w rzeczywistości nie spotykamy się z ludźmi żyjącymi 200 czy 300 lat. Tak zwany kontekst danych przechowywanych przez atrybut może pomóc przy manualnym grupowaniu przechowywanych w nim wartości.

2.2 Wsparcie i ufność

Z eksploracją reguł asocjacyjnych nierozłącznie związane są dwa pojęcia: wsparcia oraz ufności. Oba będą często gościły w niniejszej pracy, dlatego zostaną krótko przedstawione.

Wsparcie (ang. *support*) dla zbioru atrybutów X to liczba (lub procent) wszystkich transakcji bazy D , które zawierają atrybuty X . Oznaczane jest poprzez $sup(X)$.

Prostymi słowami można stwierdzić, że jest to miara, która określa jak często w bazie danych pojawiają się cechy ze zbioru X . Patrząc odwrotnie na powyższą definicję należy wprowadzić pojęcie wspierania:

Transakcja T_j wspiera zdarzenie X , wtedy kiedy zawiera wszystkie atrybuty zawarte w X (może zawierać ich więcej niż posiada X).

W tej pracy zostało założone iż dziedziną funkcji wsparcia jest zbiór $\langle 0, 1 \rangle$. Zatem miara ta wyznacza „ułamek” wszystkich transakcji. Terminu *wsparcie*

używa się najczęściej nie w odniesieniu do zbioru atrybutów lecz w nawiązaniu do reguły asocjacyjnej:

Wsparcie reguły asocjacyjnej $A \Rightarrow B$ ($\text{sup}(A \Rightarrow B)$) to wsparcie zbioru $A \cup B$ ($\text{sup}(A \cup B)$).

Ta miara wskazuje jak silna jest przedstawiona reguła w konkretnym zbiorze danych. Im większa jest wartość omawianej funkcji tym częściej zbiór (na którym oparta jest reguła) pojawia się w bazie. W tym miejscu należy zaznaczyć, że to użytkownik lub badacz decyduje jak „silnych” reguł potrzebuje. W większości algorytmów taka decyzja jest respektowana poprzez ustalenie minimalnej wartości wsparcia. Celem jest nie generowanie reguł, dla których wsparcie znajduje się poniżej uzgodnionego progu.

Drugim niezbędnym pojęciem przy eksploracji jest **ufność** (lub **wiarygodność**, ang. *confidence*).

Ufność reguły asocjacyjnej $X \Rightarrow Y$ to część liczby transakcji które wspierają Y wśród tych które wspierają X . Oznaczana jest przez $\text{conf}(X \Rightarrow Y)$

$$\text{conf}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\text{supp}(X)}$$

Sam problem wydobywania reguł jest zazwyczaj opisywany właśnie za pomocą tych dwóch miar poprzez podanie minimalnych ich wartości. To znaczy, że pożądane są tylko te asocjacje, których wsparcie jest wyższe niż minimalne wsparcie, a ufność jest większa niż minimalna ufność. Te wartości graniczne są oznaczane odpowiednio *minSup* oraz *minConf* (z ang. minimumSupport oraz minimumConfidence). Oba ograniczenia są ustawiane przez użytkownika, w celu sterowania procesem eksploracji.

Z przedstawionymi powyżej definicjami związane jest jeszcze jedno pojęcie. Każda reguła asocjacyjna jest stworzona w oparciu o pewien **zbiór częsty** (ang. frequent itemset). Jest to zbiór atrybutów których wsparcie jest większe niż minimalne wsparcie minSup. Innymi słowy, problem eksploracji reguł można sprowadzić najpierw do zadania wyszukania zbiorów częstych. Następnie każdy z nich jest wykorzystywany do budowy kilku implikacji.

Wszystkie odsłonięte tutaj pojęcia będą wykorzystywane w kolejnych rozdziałach tej pracy. Natomiast ich dokładniejsze przedstawienie znajduje się w pracy [1]

3 Ilościowe reguły asocjacyjne

Znakomita większość algorytmów odkrywających zbiory częste oparta jest na założeniu akceptacji baz danych o wyłącznie binarnych wartościach. Jest

to prawdopodobnie spadek z przeszłości po popularnej analizie koszyka zakupowego (ang. Market Basket Analysis). Kiedyś stosowane były w większości dane binarne. Z kolei obecnie częściej zbierane są informacje o ciągłych dziedzinach. Przykładem mogą być wyniki badań nad kwiatem o wdzięcznej nazwie Kosaciec, które zaowocowały legendarną już bazą danych Iris². Wiedza jakie relacje ukryte są pomiędzy atrybutami numerycznymi jest niezwykle przydatna. W tym celu buduje się ilościowe reguły asocjacyjne (ang. quantitative association rules)³. Mogą one opisywać implikację danych o dowolnych wartościach (w tym numerycznych). W literaturze można spotkać następujące konwencje zapisu:

$$A = [a_1, a_2] \wedge B = [b_1, b_2] \Rightarrow C = [c_1, c_2]$$

lub

$$A : [a_1, a_2] \wedge B : [b_1, b_2] \Rightarrow C : [c_1, c_2]$$

Formalnie: $a_1 \leq A \leq a_2$ oraz $b_1 \leq B \leq b_2$ oraz $c_1 \leq C \leq c_2$

Wszystkie zapisy służą pokazaniu, iż tym razem nie tylko stwierdzone jest, że atrybuty A, B oraz C biorą udział w implikacji, ale dodatkowo podane są przedziały wartości jakie wchodzi w jej skład. To jest jedyna różnica w stosunku do reguł binarnych, poza tym wszystkie definicje oraz miary (wsparcie, ufność) pozostają dalej uznawane.

W powyższych zapisach zastosowano przedziałowy selektor danych. To znaczy, że wybrane wartości należały do konkretnego przedziału (np. $[a, b]$). Istnieją też inne rodzaje selektorów (deskryptorów) takie jak większościowy oraz równościowy. Pierwszy z nich do selekcji danych używa operatorów porównania $<$ lub $>$ albo ich słabszych odpowiedników \leq lub \geq . Ostatni deskryptor, oznaczany znakiem $=$ wybiera konkretne wartości. Wszystkie selektory są stosowane identycznie jak ich matematyczne odpowiedniki. Do atrybutów numerycznych można stosować wszystkie podejścia. Mimo to w dalszej części pracy wykorzystywany jest wyłącznie zapis przedziałowy.

Warto przyjrzeć się pojedynczemu składnikowi reguły: $X = [x_1, x_2]$. Zostało założone iż rozpatrywany przedział jest domknięty z obu stron. Zasadność takiego warunku można zauważyć na prostym przykładzie. Mając zbiór $Z = \{1, 147, 2, 151, 148, 3\}$ można bezspornie podzielić go na dwa podzbiory: $Z_1 = \{1, 2, 3\}$ oraz $Z_2 = \{147, 148, 151\}$. Brzegowe wartości obu zbiorów (uporządkowanych) mogą utworzyć granicę przedziałów dzielących dziedzinę

²C.L. Blake and C.J. Merz. UCI repository of machine learning databases. University of California, 1998 (<https://archive.ics.uci.edu/ml/datasets/Iris>)

³Termin ten wprowadza artykuł [4]

danego zbioru: [1, 3] i [147, 151]. Jak widać, nie jest możliwe stworzenie ciągłego podziału (tak by suma przedziałów tworzyła ciągły przedział), ponieważ zbiór Z nie daje żadnych informacji na temat luki (3,147). Jakiegokolwiek założenia mogłyby być sprzeczne z rzeczywistością. Dystans ten nie może zostać zniwelowany, a oba przedziały muszą być domknięte z obu stron aby dobrze opisywać zbiory Z_1 i Z_2 . Ten fakt dobrze wpisuje się w idee całej eksploracji reguł, której zadaniem jest odkrywanie asocjacji już istniejących w danych, bez dokonywania żadnych założeń na temat próbek, które mogą pojawić się w przyszłości. Istnieją oczywiście również dziedziny eksploracji danych, które mają odmienną filozofię, jak choćby klasyfikacja.

Podział większych zbiorów wartości ciągłych (rzeczywistych, całkowitych, itp.) nie jest już tak oczywisty jak na powyższym prostym przykładzie. Wtedy zmiana granic wpływa na wsparcie tworzonego przedziału, co jest kluczowe w eksploracji reguł. To znaczy, że dysponując pewną regułą ilościową można zbudować regułę silniejszą jeśli zwiększy się przedział jednego ze składników ilościowych. Niestety czym silniejsza jest asocjacja tym mniej może być interesująca dla użytkownika, ponieważ może być zbyt ogólna i dobrze znana.

Tu nieśmiało zarysował się problem, któremu czoła chce stawić ta praca. Ogólnie i w przybliżeniu zadanie polega na takim doborze zakresów poszczególnych atrybutów by uzyskać możliwie dobrą regułę. Choć nie ma formalnej definicji dla „dobrej reguły” to intuicyjnie może być ona pojmowana jako ta pomiędzy regułą bardzo silną (lecz jednocześnie pospolitą i nieciekawą), a bardzo słabą (przez co być może incydentalną i osobliwą, choć zapewne intrygującą).

3.1 Przykład

Dla zaprezentowania teorii wprowadzonej do tej pory, przedstawiona zostanie prosta reguła ilościowa. W tym celu stworzony został niewielki zbiór danych. Zawiera on wiek, kolor oczu oraz wynik pewnego testu dla dziewięciu podmiotów. Pierwszy atrybut jest numeryczny, drugi nominalny, ostatni binarny. Dane zawiera tabela 1

Obserwując uważnie w szczególności atrybut o wartościach całkowitych można zauważyć iż w zbiorze kryje się kilka reguł asocjacyjnych. Dla przykładu można przyjąć, że poszukiwane asocjacje powinny mieć wsparcie większe lub równe $\frac{4}{9}$ całego zbioru (to znaczy, że zbiór częsty będzie wspierany przez co najmniej 4 transakcje).

Analizując atrybut *wiek* widać naturalny wręcz podział na dwie kategorie: osobników młodych $\langle 4, 9 \rangle$ oraz bardziej doświadczonych $\langle 48, 51 \rangle$. Taki podział idealnie wpisuje się w wymaganie minimalnego wsparcia, ponieważ

TID	wiek	kolor oczu	wynik testu
1	9	brązowy	1
1	48	zielony	1
3	4	niebieski	0
4	51	brązowy	1
5	5	niebieski	0
6	49	zielony	0
7	7	niebieski	1
8	5	niebieski	1
9	51	brązowy	1

Tablica 1: Przykładowe dane dotyczące oczu.

oba podzbiory go spełniają. Teraz wprost manifestuje swoją obecność reguła:

$$\underbrace{wiek : [4 - 9]}_{wsparcie=\frac{5}{9}} \Rightarrow \underbrace{kolor_oczu = niebieski}_{wsparcie=\frac{4}{9}}$$

Pod obiema stronami implikacji zapisane są wsparcia poszczególnych podzbiorów. Natomiast cała reguła ma wsparcie $\frac{4}{9}$ co w przybliżeniu daje 44% i ufność $\frac{4}{5}$ czyli 80%.

4 Analiza skupień

Analiza skupień (ang. cluster analysis) lub inaczej grupowanie (ang. data clustering) to rozwiązanie problemu odkrywania struktury grup w kolekcji obiektów nieoznaczonych żadnymi etykietami klas. Jest metodą klasyfikacji bez nadzoru (ang. unsupervised learning). To znaczy, że wyników podziału na zbiory nie można porównać z rozwiązaniem „referencyjnym”, ponieważ takie nie istnieje dla metod bez nadzoru. Dodatkowo ten fakt powoduje niejednoznaczność grupowania oraz brak obiektywnej i globalnej miary jego poprawności.

Zdefiniowanie problemu:

Celem analizy skupień jest taki podział obiektów, by te znajdujące się wewnątrz grup były maksymalnie podobne do siebie, natomiast podobieństwo pomiędzy tymi, które przynależą do różnych skupisk powinno być minimalne. Innymi słowy, optymalizując powyższe warunki, poszukiwane są spójne podzbiory danych.

Wygląda na to, że idea i kierunek działania dla rozwiązania tego problemu jest znany. Niemniej jednak w praktyce zastosowanie opisanej optymalizacji

jest trudne. Sam problem grupowania jest uznawany za NP-zupełny (odniesienie w [13]). W dodatku jego przynależność do metod uczenia bez nadzoru sprawia iż problematyczne jest zdefiniowanie kompletnego celu końcowego. A co za tym idzie, niełatwo znajduje się warunki zakończenia działania dla algorytmów rozwiązujących to zadanie.

4.1 Miara podobieństwa

Pod pojęciem grupy⁴ kryje się więc zbiór obiektów, które są „podobne”, wobec tego mają one porównywalne właściwości. W takiej formie jest to dość mgliste i niejednoznaczne wyjaśnienie. Dlatego w każdym konkretnym przypadku stosowania analizy skupień konieczne jest ściśle zdefiniowanie miary podobieństwa (lub adekwatnie niepodobieństwa) rozpatrywanych obiektów.

Konieczność powołania takiej definicji może budzić zaniepokojenie w momencie napotkania danych o charakterze kategoriowym (podrozdział 2.1). Są to wartości ze skończonej przestrzeni, dla których nie sposób wyznaczyć funkcję odległości. Jednakże nawet takie warunki nie budują bariery nie do pokonania, czego dowodem są wyniki prac [9] oraz [8] - rozdział „Distance Measures”.

Zadanie wytyczenia funkcji podobieństwa przyjmuje formę bardziej przejrzystą i niebudzącą tylu wątpliwości jeśli grupowaniu podlegają tylko obiekty opisane właściwościami o wartościach ciągłych (numerycznych). Takie założenie uprawnia do wyboru jednej z wielu funkcji odległości i zbudowania w oparciu o nią adekwatnej funkcji podobieństwa. Praca [6] przedstawia odległość Euklidesową jako najczęściej stosowaną w takim przypadku.

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Jest to uzasadnione dwoma czynnikami. Pierwszy to fakt iż miara ta jest bardzo dobrze znana i powszechnie używana. Drugi, że dzięki temu jest intuicyjna albo sprawia wrażenie właśnie takiej. Powyższa funkcja stanowi szczególnie przypadek miary Minkowskiego, która z kolei opisuje odległości w przestrzeniach o większej liczbie wymiarów:

$$L_m(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n |x_i - y_i|^m \right)^{1/m}$$

Nadzwyczaj uciążliwą wadą tej ostatniej miary jest bardzo szybka utrata intuicyjności wraz ze wzrostem wymiarowości przestrzeni. O ile dla czterowymiarowych danych można odważyć się na pewne wnioski płynące z analizy

⁴Pojęcia: grupa, skupisko, zbiór, są tutaj używane wymiennie.

odległości pomiędzy obiektami, o tyle już dla podwojonej ich liczby czyli ośmiu wymiarów jest to karkołomne wyzwanie. Należy zauważyć, że pomimo iż nikogo nie przerażają wektory o kilkunastu elementach to jednak przestrzenie wielowymiarowe mogą być kłopotliwe w interpretacji. To naturalne, że człowiek traci orientację i intuicję „matematyczną” po wyjściu poza przestrzeń którą zna od dziecka (trójwymiar). Dlatego w większości przypadków, funkcja podobieństwa jest po prostu bezdyskusyjnie przyjmowana jako funkcja odległości. W praktyce rozważania są ograniczane do rozwiązania następnego kłopotu.

4.2 Normalizacja

Kolejnym problemem (tym razem niezależnym od liczby wymiarów) jest dominacja składników o „szerokich” dziedzinach. Innymi słowy takich, dla których zakres wartości jest większy niż w innych składnikach przestrzeni. Na całą miarę Minkowskiego wpływ ma suma „odległości” na poszczególnych wymiarach (atrybutach). Najlepiej, żeby wszystkie składniki miały wartości zbliżone do siebie (choćby pod względem rzędu wielkości). W przeciwnym przypadku nawet pojedynczy składnik może zdominować zachowanie całej sumy. Można zauważyć to zjawisko dysponując nawet tylko dwoma atrybutami, z których jeden przyjmuje wartości z bardzo małego zakresu, a drugi rozciąga się na całą, bardzo szeroką dziedzinę. Można spostrzec, że w ramach drugiego wymiaru różnica odległości dowolnych dwóch punktów będzie zawsze dużo większa niż dla pierwszego. Tym samym więc ten drugi ma większy wpływ na wartość całej funkcji odległości. Zatem niewielkie zmiany położenia w ramach „szerokiego” atrybutu spowodują stosunkowo duże zmiany wartości pełnej miary Minkowskiego. Takie zachowanie obniża intuicyjność oraz uniemożliwia ustalenie warunków dla których można mówić o obiektach podobnych do siebie (co często czyni się na podstawie wyznaczenia progu dla odległości). Lekiem na wyżej wymienione dolegliwości może być **normalizacja** danych oddzielnie względem poszczególnych atrybutów. Zazwyczaj używane jest skalowanie przy wykorzystaniu zakresu lub wariancji [6]. W celu rozwiązania tego problemu adoptuje się również rozwiązania ze statystyki, w tym standaryzację:

Standaryzacja (ang. *standard score*) jest sposobem normalizacji zmiennej losowej. Po tym procesie zmienna posiada zerową wartość oczekiwaną ($\mu = 0$) oraz wariancję równą jeden ($\sigma = 1$). Najpopularniejsza jest *standaryzacja Z*, którą opisuje wzór:

$$z = \frac{x - \mu}{\sigma}$$

Podsumowanie

Istnieją też inne miary odległości wykorzystywane w analizie skupień, jak chociażby nieskomplikowana odległość Manhattan czy ambitna miara Mahalanobis. Jednakże w ramach tej pracy badane będą tylko i wyłącznie dane o charakterze liczb rzeczywistych oraz najczęściej jednowymiarowej dziedzinie, stąd prosty wniosek, że wystarczająca jest miara Minkowskiego.

Podsumowując, analizie skupień mogą podlegać dowolne obiekty pod warunkiem, że istnieje dla nich miara podobieństwa. Mimo to najczęściej pod pojęciem *obektu* kryje się uporządkowany wektor cech (atrybutów) o wymiarze d . Dzięki temu można rozpatrywać go jako *punkt* w d -wymiarowej przestrzeni. W takim kontekście można mówić nawet o kształcie grup: wklęsłe, wypukłe, albo bardziej drobiazgowo: koliste, prostokątne, itp. Niektóre algorytmy jako swoje właściwości mają również specyfikowany kształt zbiorów wynikowych⁵.

Daleko idący i atrakcyjny przegląd dziedziny analizy skupień został przedstawiony w pracy [6].

4.3 Rodzaje analizy skupień

Współczesna informatyka dostarcza szerokiego wachlarza rozwiązań problemu analizy skupień. Każde z nich działa na swój sposób i ma wyjątkowe właściwości. Ta różnorodność jest cechą pozytywną ze względu na fakt, że nie wszystkie zadania grupowania są takie same. Oczywiście jest, iż w praktyce stosuje się różne miary podobieństwa. Ale dodatkowo stawiane są różnorodne wymagania odnośnie właściwości jakie mają spełniać nowo powstałe grupy. I chociażby z tego ostatniego powodu można podzielić wszystkie metody analizy skupień na kilka kategorii, oto lista pojęć jakimi są one identyfikowane:

4.3.1 Właściwości grup wynikowych

To jakie cechy będą mieć przedziały po analizie skupień zależy od zastosowanego algorytmu. Natomiast, to jakie są dostępne opcje prezentuje poniższa lista:

1. **równe szerokości przedziałów** (ang. *equi-width*) – każdy interwał ma w przybliżeniu równą średnicę⁶.

⁵Przykładem jest algorytm dzielenia według siatki, który produkuje wyłącznie prostokątne podobszary.

⁶Przez średnicę rozumie się tutaj maksymalną odległość pomiędzy dwoma dowolnymi elementami zbioru.

2. **równe głębokości** (ang. *equi-depth*) – uzyskane zbiory zawierają w przybliżeniu tyle samo elementów. Nazwa angielska i jej polskie tłumaczenie najprawdopodobniej są zainspirowane grupowaniem hierarchicznym, w którym to poziom głębokości definiuje również w pewnym sensie licznosc grup.
3. **jednolite przedziały** (ang. *homogeneity-based bins*) – rozmiar grupy jest tak dobierany, by rozkład jej elementów był możliwie jednolity. Taką kategorię opisuje praca [13].
4. **grupowanie rozmyte** (ang. *fuzzy clustering*) – otrzymane grupy mogą mieć parami niepuste części wspólne. W potocznym znaczeniu - częściowo nakładają się na siebie. Kluczowe znaczenie ma dobór ograniczeń na wspólny podzbiór, czyli odpowiedź na pytanie jak bardzo grupy mogą się przenikać ze sobą. Drugim pytaniem jest, czy dopuszczalne jest zawieranie zbiorów (być może przydatne w pewnych specyficznych zastosowaniach).

Formalnie: każdy obiekt zbioru może z pewnym prawdopodobieństwem należeć do każdego z powstałych przedziałów. Dla przykładu: dysponując zbiorem grup $\{G_1, G_2, G_3\}$, element x , dzięki swoim cechom, z prawdopodobieństwem 0,78 powinien przynależeć do zbioru G_1 , natomiast z 0,12 do G_2 oraz z 0,1 do G_3 . Stosując grupowanie rozmyte nie dokonuje się rozstrzygnięcia do której z grup należy dany obiekt. Jest on pojmowany jako członek wszystkich grup z dodatkową informacją na temat „poziomu” tego członkostwa.

5. **grupowanie „twarde”** (ang. *hard clustering*) – każdy z analizowanych obiektów może należeć tylko i wyłącznie do jednej z grup. Dlatego utworzone zbiory nie posiadają części wspólnych. Jest to tradycyjne podejście do tematyki analizy skupień, a przez to jest najczęściej stosowane. Jego przeciwieństwem są metody rozmyte.
6. **grupowanie „naturalne”** (ang. *natural clustering*) - to koncepcja wydobycia ze zbioru takich skupisk, jakie tam się naturalnie znajdują. Sensem tego podejścia jest jak najmniejszy wpływ sztucznie i manualnie wybranych parametrów na końcowy podział. Do tej kategorii można zaliczać metody gęstościowe, jak np. algorytm DBSCAN, czy metody jakościowe, jak np. algorytm Quality Threshold.

Zapoznanie się z taką kategoryzacją analizy skupień pozwala na optymalny dobór metody, a później konkretnego algorytmu dla przedstawionego rzeczywistego problemu. Użytkownik musi wskazać dokładnie jakich cech

oczekuje. Należy jednak wspomnieć o tym, że istnieje możliwość mieszania powyżej przedstawionych opcji. Jest dozwolone zbudowanie takiego algorytmu, który będzie rozpatrywać zarówno kryterium równej szerokości tworzonego zbioru, ale również będzie czuły na ilość zawartych w nim obiektów lub ich rozkład. Wszystko zależy od warunków jakie są stawiane przed procesem grupowania.

4.3.2 Metody podziału

Znając już warunki jakie musi spełniać pożądane grupowanie, należy wybrać jeszcze jedną z wielu procedur grupujących⁷. Ponownie decyzję można oprzeć o kilka ogólnych kategorii. Procedury analizy skupień dzielą się na:

1. **hierarchiczne** (ang. *hierarchical*) – grupy są łączone ze sobą w celu utworzenia hierarchii – dwie mniejsze grupy mogą tworzyć ze sobą większą (pod warunkiem, że są do siebie „podobne”). Mówiąc inaczej, cały zbiór zawiera kilka mniejszych zbiorów, z kolei te zawierają w sobie jeszcze mniejsze i tak dalej. Tutaj należy wyróżnić dwa podejścia:
 - (a) **wstępujące** – każdy obiekt początkowo traktowany jest jako grupa, następnie iteracyjnie grupy są łączone w pary. Taki proces kończy się zazwyczaj po połączeniu dwóch ostatnich grup w cały zbiór, albo po osiągnięciu założonego wcześniej celu dotyczącego właściwości podziału. W jednym i w drugim przypadku produktem końcowym jest hierarchia grup (struktura zawierania).
 - (b) **zstępujące** – metoda rozpoczyna działanie od całego zbioru i dzieli go na mniejsze fragmenty, te nowo powstałe są rozcinane w podobny sposób rekurencyjnie. W tym przypadku warunek stopu też podlega ustaleniu. Absolutny koniec to uzyskanie wszystkich jednoelementowych zbiorów. W wyniku działania otrzymywana jest hierarchia grup (struktura dzielenia). Lecz tym razem, brak kompletnej struktury, to brak pełnej wiedzy na temat podobieństwa „małych” zbiorów, łącznie z jednoelementowymi.

Jeśli dowolna miara odległości zostanie wybrana jako funkcja podobieństwa, to dla obu powyższych podejść należy dokonać wyboru sposobu obliczania odległości pomiędzy grupami. Innymi słowy, zbiór obiektów jako taki nie stanowi pojedynczego punktu w przestrzeni, co powoduje niejednoznaczności przy obliczaniu funkcji odległości. Może być

⁷ Obecne rozważania nadal nie nawiązują do żadnych konkretnych algorytmów.

ona ustalona pomiędzy „środkami”⁸ zbiorów. Ale również możliwe jest zastosowanie odległości maksymalnej, czy minimalnej. Zagadnienie to nazywane jest *metodami wiązania*, a więcej na ten temat w [10] w rozdziale Analiza skupień.

Algorytmy hierarchiczne to jedne z niewielu, które oprócz wyboru funkcji podobieństwa nie wymuszają dodatkowych parametrów. Natomiast na samą funkcję nie ma nałożonych żadnych ograniczeń ani wymagań (oprócz wartości zwracanej, która ma być miara określającą jak bardzo jeden obiekt jest podobny do drugiego). Budowana jest wtedy cała hierarchia, którą można analizować nawet wizualnie na dendrogramach, czy diagramach venna.

2. **metody oparte na podziałach** (ang. *partitional clustering*)⁹ To metody stosujące pojedynczy i rozłączny podział całego zbioru zamiast hierarchii podziałów. Biorąc pod uwagę ich cechy można wyznaczyć pewne kategorie i są to kolejno procedury:

- (a) **deterministyczne** albo **niedeterministyczne** – niektóre algorytmy grupujące uznawane są za niedeterministyczne, ponieważ ich działanie opiera się w jakimś stopniu o „losowość”. Przykładem może być algorytm k-means, który losowo ustala położenie pierwotne środków grup.
- (b) **minimalizujące błąd wynikowy** (ang. *Error Minimization Algorithms* [8]) – tak właściwie jest to raczej idea stojąca za niektórymi konkretnymi algorytmami. Polega ona na minimalizacji pewnego kryterium błędu, najczęściej opisanego funkcją odległości. Najbardziej znany jest błąd kwadratowy (ang. *squared error algorithms* [6]), a konkretnie błąd średniokwadratowy (ang. MSE - *Mean Squared Error*), czy suma kwadratów błędów (ang. *Sum of Squared Error*). Należy więc ustalić jakie odległości są rozpatrywane – najczęściej badana jest odległość punktów do środków grup ich zawierających. Cały proces polega wówczas na globalnej optymalizacji funkcji błędu.

Dobłą ilustracją tej procedury jest algorytm k-means, który początkowo losuje pakiet środków przyszłych grup, np. $\{S_1, S_2, S_3\}$. Niestety to jest jego wada, że liczbę i położenie tych środków

⁸To tak zwana metoda środka ciężkości zbioru. Można go wyznaczyć poprzez średnią (zwykłą lub ważoną) wartości we wszystkich punktach zbioru.

⁹Definicja zaczerpnięta z encyklopedii uczenia maszynowego [11] - hasło *partitional clustering*

należy ustalić apriori. Celem jest takie przemieszczenie środków względem obiektów aby zminimalizować błąd (np. średniokwadratowy sumy odległości punktów od środków). Rozwinięcie zagadnienia znajduje się w pracy [8] w podrozdziale „Partitioning Methods”.

- (c) **oparte na gęstościach** (ang. *density-based algorithms*) wykonują swoją pracę analizując gęstość rozłożenia punktów zbioru w przestrzeni. Tym razem, każdy obiekt z zadanego zbioru musi być opisany wektorem liczb (odpowiadającym kolejnym jego cechem, atrybutom obiektu). Tak stworzony wektor jest traktowany jako punkt w przestrzeni (liczba wymiarów odpowiada rozmiarom wektora). Tworzone grupy cechują się ustaloną i jednolitą wewnętrzną koncentracją. Natomiast dzięki kontrastowi gęstości można wydzielać kolejne zbiory. Punkty należące do obszarów o niskim nasyceniu są traktowane jako szum (zakłócenia lub próbki odstające od zbioru) i nie podlegają grupowaniu. Przykładem implementacji tej metody jest algorytm DBSCAN¹⁰

Algorytmy zbudowane w oparciu o tą idee mają spektakularne zalety. Po pierwsze i co najważniejsze nie wymagają dogłębnej znajomości dziedziny analizowanych danych (ani liczby ani położenia środków). Pozwalają także na wykrycie grup o bardzo dowolnym, a nawet wyrafinowanym kształcie. Te dwa argumenty już sprawiają, że można podchodzić do nich bardzo entuzjastycznie. Niemniej jednak takie atuty mają też swoje skutki uboczne. Problemem bywa ustalenie poziomu koncentracji wewnętrznej grup oraz poziomu gęstości poniżej której punkty traktowane są jako szum. Procedura jest w stanie działać bardzo dobrze dla danych w których skupiska punktów są dobrze odizolowane, co niestety nie zawsze zachodzi w świecie rzeczywistym.

- (d) **oparte na grafach** (ang. *graph-theoretic clustering* lub *graph-based clustering*) – w zarysie polegają na rozpięciu grafu o węzłach w punktach zbioru. Najbardziej znane algorytmy wykorzystują minimalne drzewo rozpinające (ang. *Minimal Spanning Tree*). Artykuł [8] – rozdział „Graph-theoretic clustering”.
- (e) Istnieje jeszcze wiele innych podejść do metody grupowania opartej na podziałach. Są to między innymi procedury wykorzystujące sieci neuronowe ([8] – „Neural networks”), czy nawet techniki

¹⁰Ester M. i in.: A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise, Proc. of 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96)

ewolucyjne. W nielicznych sytuacjach można również zastosować podział według narzuconej „siatki” (ang. *grid-based methods* [8]). Jednakże dotychczas przedstawione metody oraz procedury stanowią wystarczającą bazę dla wyboru rozwiązania najlepiej spełniającego cel niniejszej pracy.

4.4 Rozważania na temat właściwości grupowania

Dysponując całym wachlarzem algorytmów grupujących, należy wybrać właściwości jakie są oczekiwane względem zbiorów końcowych. Pewną wskazówką może być klasyfikacja tych algorytmów naszkicowana w rozdziale 4.3 Rodzaje analizy skupień.

Odniesienie do właściwości podziałów (rozdz. 4.3.1)

Analiza skupień w przypadku tej pracy nie może ograniczyć się do grupowania ani o równych szerokościach ani o zbliżonej liczności. To specyfika i naturalny rozkład danych w bazie powinien narzucić rozpiętość i rozmiary poszczególnym grupom. Z drugiej strony zostaną wprowadzone zdroworozsądkowe ograniczenia oparte właśnie na tych parametrach. Wiadomo, że szerokość grupy określa poniekąd jej „jakość”. Natomiast liczba elementów będzie niezmiernie ważna dla późniejszego odkrywania reguł.

Tworzenie pakietu rozłącznych grup to podejście tradycyjne i dobrze zakorzenione w intuicji wielu ludzi. Niemniej jednak warto wykorzystać potencjał grupowania rozmytego. Ta kwestia rozstrzygnie się w kolejnych rozdziałach.

Odniesienie do metod podziałów (rozdz. 4.3.2)

W kontekście dostępnych ogólnych schematów postępowania, pierwsze pytanie na jakie trzeba odpowiedzieć, to czy hierarchia grupowania jest potrzebna?

Zdecydowanie wymagany jest płaski podział, który można następnie potraktować jako atrybut nominalny dla reguł asocjacyjnych. Dlatego ani hierarchia grup, ani dokładność tego procesu nie są przydatne. Po przedstawieniu reguły, np. $wiek : [5 - 10] \Rightarrow wzrost = niski$ nie jest już istotne czy przedział $[5 - 10]$ został zbudowany w oparciu o dwa czy trzy inne przedziały. Nie jest też ważne jak grupują się obiekty pojedyncze. Natomiast z punktu widzenia metody zapewniania prywatności, grupowanie hierarchiczne z zachowywaniem hierarchii byłoby nie do przyjęcia.

W następstwie powyższych wniosków, algorytmu należy szukać wśród metod opartych na podziałach „płaskich”. A w tych ramach, należy odpowiedzieć na pytanie, czy zachowanie losowe jest w tym zastosowaniu akceptowalne?

Grupowanie jest elementem transformacji, która jak było sugerowane wcześniej powoduje przekształcenie całej bazy (nawet fizycznie, z jednego

pliku na drugi). Teoretycznie nie wymusza to deterministyczności całego procesu. Raz przekształcona baza może służyć do wielokrotnego poszukiwania reguł. Ostatecznie cały proces można powtórzyć kilka razy przed właściwą eksploracją reguł, żeby uniknąć negatywnego skutku „niefortunnej” sytuacji w momencie losowania.

Nie istnieje niestety możliwość dobrego zastosowania metody gęstościowej dla późniejszej eksploracji reguł asocjacyjnych. Dzieje się tak z kilku powodów, po pierwsze, jak zwykle problematyczne bywa ustalenie wartości wszystkich parametrów, z poziomem gęstości na czele. Dodatkowo w oryginalnych algorytmach nie ma możliwości sterowania wielkością grupy ani zakresem jej wartości. Ostatnim problemem, jest fakt iż ta metoda pozwala na wyodrębnienie podobszarów, które są „kontrastowe” między sobą względem gęstości. Tymczasem dla eksploracji reguł cenny jest nawet podział przestrzeni spójnej. Oczywiście podział końcowy powinien opierać się na skupiskach możliwie odrębnych, ale jeśli zajdzie potrzeba rozcięcia zwartego skupiska na kilka części, to też powinno być wykonane.

Ani algorytmy grafowe ani bezpośrednie zastosowanie sieci neuronowych nie są możliwe w przypadku baz przeznaczonych do eksploracji. Głównym powodem jest zazwyczaj duży rozmiar używanych kolekcji. Jednocześnie złożoność pamięciowa tych rozwiązań jest delikatnie mówiąc godziwa, co przyczynia się do tego iż ich zastosowanie staje się mało zachęcające.

To krótkie podsumowanie uświadamia możliwości w wyborze algorytmu analizy skupień z zastosowaniem dla budowania reguł ilościowych. Pomijając wykluczone podejścia, w rozwiązaniu można zastosować metody oparte na podziałach w tym szczególnie deterministyczne, ale również te z elementami niedeterministycznymi. Dopuszczalne są metody z szeroko pojętej kategorii minimalizujących błąd wynikowy.

5 Zbiór wartości binarnych

Po wstępie teoretycznym musi pojawić się uzasadnienie, dlaczego tak, здаwałoby się, odmienne dziedziny eksploracji danych są przedstawiane w tej pracy razem. Łączy ich cel – odkrycie interesujących ilościowych reguł asocjacyjnych. Podążając za pracą [4] wyróżnia się dwa podstawowe typy reguł asocjacyjnych: binarne (ang. *boolean association rules*) oraz ilościowe (ang. *quantitative association rules*). Jak już zostało wspomniane wcześniej, do odkrywania asocjacji binarnych istnieje szereg algorytmów, natomiast w drugim przypadku jest ich niewiele oraz bywają mało efektywne. Stąd już w 1996 roku powstał pomysł, żeby zbiór danych formatu numerycznego przekształcić w kolekcję pozycji $\{0, 1\}$. Tym samym problem wyszukiwania ilościowych re-

guł staje się zadaniem odkrycia binarnych reguł asocjacyjnych. Ten rozdział zawiera przegląd możliwości zamiany pojedynczych atrybutów na binarne.

5.1 Transformacja na wartości binarne

W odniesieniu do **atrybutów kategoriycznych** sposób przekształcenia jest oczywisty. Dla przykładu niech obiekty będą charakteryzowane przez „*atrybut1*”, który przyjmuje n wartości. Każdy element tej cechy powinien stworzyć nowy atrybut, klasycznie pod tytułem „*atrybut1:wartośćK*”. Tym razem jest on już dwuwartościowy, znów klasycznie, przyjmuje się dziedzinę $\{0, 1\}$. Jedynka oznacza, że pierwotna cecha „*atrybut1*” miała w danej transakcji wartość „*wartośćK*”, zero wstawiane jest w przeciwnym przypadku.

W odniesieniu do **kolekcji typu numerycznego** (np. liczb rzeczywistych) droga do zbioru dwuwartościowego nie jest tak bezdyskusyjna. Nie istnieje jednoznaczne i najlepsze rozwiązanie. Być może istnieją problemy, które będą wymagać w tym miejscu odwzorowania każdej unikalnej wartości w zbiór $\{0, 1\}$. Ale w znakomitej większości przypadków taka procedura doprowadzi do uzyskania ogromnej i bardzo rzadkiej macierzy.

Większość pozycji literatury dotyczącej problemu reguł ilościowych sugeruje iż transformacja atrybutów numerycznych powinna odbyć się poprzez dyskretyzację ich dziedziny wartości. Następnie tak podzielony atrybut można potraktować jako kategoriyczny i przekształcić w dwuwartościowy tak jak to jest opisane wyżej.

5.2 Dyskretyzacja

Jest to procedura przetwarzająca informacje ciągłe w dyskretne (źródło: [6]). Formalnie, szczególnie w matematyce, pojęcie to dotyczy procesu przekształcania modeli ciągłych w dyskretne ich odpowiedniki. W topologii dotyczy przekształcania przestrzeni spójnych w przestrzenie dyskretne. Przestrzeń spójna intuicyjnie składa się z „jednego kawałka”, natomiast dyskretna opiera się o punkty, które są niejako „oddzielone” od siebie.

W kontekście eksploracji danych, do dyspozycji jest kolekcja wartości konkretnego atrybutu ciągłego. Sam atrybut można potraktować jako pewną funkcję. Jej przeciwdziedzina tworzy przestrzeń spójną. Dlatego proces transformacji na dane binarne z tego punktu widzenia wygląda na zadanie dyskretyzacji przeciwdziedziny (po prostu w celu zamiany jej na topologiczną przestrzeń dyskretną). Uzyskany podział powinien pokryć całą przestrzeń. Niemniej jednak zazwyczaj o przeciwdziedzinie wiadomo niewiele. Znany jest wyłącznie skończony zbiór próbek tej przestrzeni. Dlatego w praktyce, mniej

lub bardziej słusznie, mówi się o dyskretyzacji kolekcji próbek i w tym celu wykorzystuje się metody analizy skupień.

5.3 Właściwości dyskretyzacji

Sam proces grupowania niesie za sobą pewne konsekwencje. Przede wszystkim spowoduje częściową utratę informacji, a jednocześnie ustanowi pewien poziom prywatności danych pierwotnych. Ma to praktyczne zastosowania, a jedno wynika z tego co już zostało wspomniane wcześniej: podstawowym i najprostszym sposobem zapewnienia prywatności danych indywidualnych jest ich podział na przedziały. Natomiast utrata konkretnych wartości ma kilka zastosowań ogólnych. Są to między innymi próby uproszczenia danych, na przykład w celu ich prezentacji wizualnej. Innym przykładem jest chęć przygotowania wstępnej analizy lub próby szybkiego wyodrębnienia ogólnej wiedzy.

5.4 Przykład transformacji

W celu zaprezentowania działania transformacji atrybutów różnego rodzaju, użyty zostanie zbiór danych zawarty w tabeli 1 omawiany już w rozdziale 3.1. Zawiera on zarówno wartości całkowite jak i nominalne. Minimalne wsparcie niech tym razem wynosi 40%. Pierwszy atrybut (kolumna wiek) zostanie poddany analizie skupień, czego wynikiem będzie utworzenie dwóch przedziałów $\langle 4 - 9 \rangle$ oraz $\langle 48 - 51 \rangle$. Wartości w dotychczasowych rekordach zostaną zamienione na wspomniane przedziały. Ten krok prezentuje tabela 2, która jest już bazą danych wyłącznie kategoriycznych. W tym momencie należy zastosować binaryzację. Dla przykładu, drugi atrybut jest przekształcony na tyle kolumn ile zawiera unikalnych wartości, ostatni pozostanie bez zmian. Wynik prezentuje tabela 3 zawierająca binarną bazę danych.

6 Analiza skupień, a wyszukiwanie reguł asocjacyjnych

Jak zostało zaznaczone już nawet we wstępie, w celu odkrycia ilościowych reguł asocjacyjnych zostanie przeprowadzone grupowanie atrybutów numerycznych. W tym celu należy uściślić właściwości tego procesu i wymagania odnośnie wynikowych grup. Sam proces eksploracji podzielony jest na dwa etapy. Pierwszy to znalezienie zbiorów częstych, drugi to budowa reguł z odkrytych zbiorów.

TID	wiek	kolor oczu	wynik testu
1	$\langle 4 - 9 \rangle$	brązowy	1
1	$\langle 48 - 51 \rangle$	zielony	1
3	$\langle 4 - 9 \rangle$	niebieski	0
4	$\langle 48 - 51 \rangle$	brązowy	1
5	$\langle 4 - 9 \rangle$	niebieski	0
6	$\langle 48 - 51 \rangle$	zielony	0
7	$\langle 4 - 9 \rangle$	niebieski	1
8	$\langle 4 - 9 \rangle$	niebieski	1
9	$\langle 48 - 51 \rangle$	brązowy	1

Tablica 2: Wyniki dyskretyzacji atrybutu wiek z tabeli 1.

TID	wiek: $\langle 4 - 9 \rangle$	wiek: $\langle 48 - 51 \rangle$	oczy: niebieski	oczy: brązowy	oczy: zielony	wynik testu
1	1	0	0	1	0	1
2	0	1	0	0	1	1
3	1	0	1	0	0	0
4	0	1	0	1	0	1
5	1	0	1	0	0	0
6	0	1	0	0	1	0
7	1	0	1	0	0	1
8	1	0	1	0	0	1
9	0	1	0	1	0	1

Tablica 3: Wyniki przekształcenia danych z tabeli 1 do postaci binarnej.

6.1 W kontekście wyszukiwania zbiorów częstych

Rozważania warto rozpocząć od przypomnienia, że zbiór częsty to taki zbiór atrybutów bazy danych, który w całym zbiorze transakcji ma wsparcie większe niż minSup . Ogólnie znane jest też twierdzenie, które mówi:

Każdy podzbiór zbioru częstego jest zbiorem częstym. To oznacza, że żaden zbiór nieczęsty nie może stać się częstym poprzez dodanie do niego nowych atrybutów. A wniosek ostateczny brzmi:

Jeśli pojedynczy atrybut jest nieczęsty, to nigdy nie wejdzie w skład reguły asocjacyjnej.

Z ostatniego stwierdzenia można wysnuć już prosty postulat: cechę numeryczną należy tak przekształcać, by nowo powstałe atrybuty były częste. Innymi słowy utworzenie grupy która nie będzie częsta nie przyniesie korzyści z punktu widzenia wyszukiwania reguł. Z tego powodu na grupowanie należy nałożyć pierwsze kryterium – powinno znać i respektować minimalne wsparcie ustalone przed procesem eksploracji. A konkretniej:

Wynikowe grupy powinny zawierać co najmniej tyle elementów ile wynosi iloczyn minimalnego wsparcia i liczby wszystkich transakcji w bazie:

$$\text{min_rozmiar_grupy} = \text{minSup} * \text{liczba_transakcji}$$

Pod warunkiem, że minSup jest dane jako ułamek, a nie jako konkretna liczba próbek.

Ten sam problem został zaprezentowany w pracy [4] pod nazwą problemu „MinSup” i zdefiniowany nieco inaczej. Tam obecna definicja mówi, że jeśli liczba przedziałów będzie duża to wsparcie pojedynczego przedziału będzie małe. A wtedy niektóre reguły zawierające tę cechę mogą zostać nie odkryte z powodu braku minimalnego wsparcia.

Bardzo trudno jest uzasadnić sens powoływania przedziałów tworzących atrybuty, które nie wezmą udziału w żadnej regule asocjacyjnej. Z punktu widzenia problemu „MinSup” im szersze zakresy przedziałów tym lepiej.

6.2 W kontekście powoływania reguł asocjacyjnych

Żeby z dowolnego zbioru częstego utworzyć regułę asocjacyjną, należy podzielić go na dwie części i jedną potraktować jako lewą stronę implikacji, a drugą jako prawą. Jeśli taka reguła spełnia warunek minimalnej ufności to może być spokojnie zaprezentowana jako jeden z wyników całego procesu eksploracyjnego. W przeciwnym przypadku zostaje odrzucona. Jak to odnosi się do procesu analizy skupień? Otóż, niski poziom wiarygodności pojawia się wtedy, kiedy wsparcie poprzednika implikacji jest duże w stosunku do wsparcia całej reguły (czyli tych transakcji, które wspierają poprzednik i następnik jednocześnie). Jeśli w poprzedniku znajduje się atrybut utworzony z

pierwotnie numerycznego, to zwiększając szerokość przedziału możliwe jest pogorszenie wsparcia całej reguły. A w konsekwencji reguła może zostać odrzucona z powodu nie spełnienia kryterium minimalnej wiarygodności (min-Conf). Innymi słowy, im mniejsze przedziały, tym lepiej z punktu widzenia poziomu ufności.

Dla przykładu warto rozważyć regułę $A \Rightarrow B$ która ma wsparcie w . Dodatkowo niech zarówno A jak i B oddzielnie też mają wsparcia w . To oznacza, że reguła ma ufność na poziomie 100%. Jeśli jednak zwiększy się zakres A , a więc jednocześnie zacznie go wspierać więcej transakcji, to poziom ufności całej reguły spadnie. Przykładowo, jeśli wsparcie A wzrośnie dwukrotnie to poziom wiarygodności zmaleje aż do 50%.

Rozwiązywanie problemów minSup i minConf działa antagonistycznie. Poprawa jednego może pogorszyć sytuację drugiego. Dlatego właśnie tak trudno jest ustalić optymalny rozmiar grup. Na szczęście wiele zbiorów danych zawiera w sobie naturalne skupiska, które spełniają dolne ograniczenie.

7 Quality Threshold Clustering

Quality Threshold Clustering (w skrócie QT Clustering albo QTC) to algorytm grupowania pierwotnie stworzony do analizy genów [15]¹¹. Jego priorytetem jest zapewnienie odpowiedniego poziomu jakości dla tworzonych grup. Nie wymaga specyfikowania potencjalnej liczby skupisk jakie wystąpią w bazie danych. Wręcz przeciwnie, pozwala na odkrycie ich naturalnej liczby. Wielkością grup wynikowych steruje parametr maksymalnej średnicy skupiska. Podstawowa idea opiera się na znalezieniu grup kandydujących na podstawie każdego punktu ze zbioru. Każda z nich musi spełniać wymaganie jakościowe i co ważniejsze każda z nich jest budowana w oparciu o pełny zbiór danych (jeszcze niegrupowanych). Ostatecznie spośród wszystkich kandydatów wybierany jest najlepszy. Elementy w nim zawarte są usuwane z grupowanego zbioru (jako już przydzielone), a cały proces jest powtarzany. Cała procedura jest zapisana w pseudokodzie poniżej – Algorytm 1

7.1 Właściwości

W oryginalnym algorytmie mianowanie najlepszego kandydata opiera się na wyborze największego pod względem liczby elementów. Niemniej jednak modyfikacja tego warunku może być wskazana w zależności od zastosowania algorytmu. Idea niezależnego generowania kandydatów z całej dostępnej puli

¹¹Informacje są zaczerpnięte również z encyklopedii uczenia maszynowego [11] (hasło: „Quality Threshold Clustering”)

Algorytm 1 Procedura grupowania Quality Threshold Clustering

```
funkcja QTCLUSTERING( $G, d$ )  
  if  $|G| \leq 1$  then return  $G$   
  end if  
  for all  $i \in G$  do  
    zbiór  $A_i \leftarrow \{i\}$   $\triangleright A_i$  jest  $i$ -tym kandydatem  
    while  $A_i \neq G$  do  
      znajdź  $j \in (G - A_i)$  dla którego  $\text{ŚREDNICA}(A_i \cup j)$  jest minimalna  
      if  $\text{ŚREDNICA}(A_i \cup j) < d$  then  
         $A_i \leftarrow A_i \cup \{j\}$   
      else  
        break while  
      end if  
    end while  
  end for  
  zbiór  $C \leftarrow \text{NAJWIĘKSZY\_ZBIÓR\_Z}(A_1, A_2, A_3, \dots, A_{|G|})$   
  return  $\{ C, \text{QTCLUSTERING}(G - C, d) \}$   
end funkcja
```

próbek sprawia, że ta metoda ma wiele unikalnych zalet. Przede wszystkim jest niezależna od kolejności budowania grup ani występowania danych. Dodatkowo nie zawiera żadnego elementu losowego. To wszystko sprawia, że jest to jeden z niewielu, w pełni deterministyczny algorytm grupujący. Uruchamiając go wielokrotnie, zawsze jest pewność tego samego wyniku. Algorytm gwarantuje też, że wszystkie grupy wynikowe będą spełniały przedstawione wymaganie jakościowe (w tym przypadku maksymalną średnicę). Na dodatek wybierane są skupiska w kolejności od najlepszego (np. największego) do najslabszego (np. najmniejszego). Daje to gwarancję, że najbardziej istotna struktura danych zostanie odkryta poprawnie.

Najważniejszą zaletą algorytmu jest fakt, że odkrywa on naturalne skupiska w zadanym zbiorze z dokładnością do ustawionej jakości. Oznacza to, że możliwe jest poznanie faktycznego obrazu struktury wartości, niesfałszowanej i niczym nie wymuszonej. Jednocześnie parametr sterujący jest dość intuicyjny dla badaczy. Dla kontrastu wystarczy wspomnieć np. konieczność ustalenia poziomu gęstości dla algorytmu DBSCAN (rozdz. Metody podziału). Tutaj ustalenie średnicy może w pierwszym momencie nie być łatwe, ale szybkie zapoznanie z dziedziną wartości (choćby zakresu min-max) wystarczy by zapanować nad sytuacją.

Algorytm ten wykazuje delikatne podobieństwo do grupowanie hierarchicznego z kompletnym łączeniem (ang. *complete linkage hierarchical clu-*

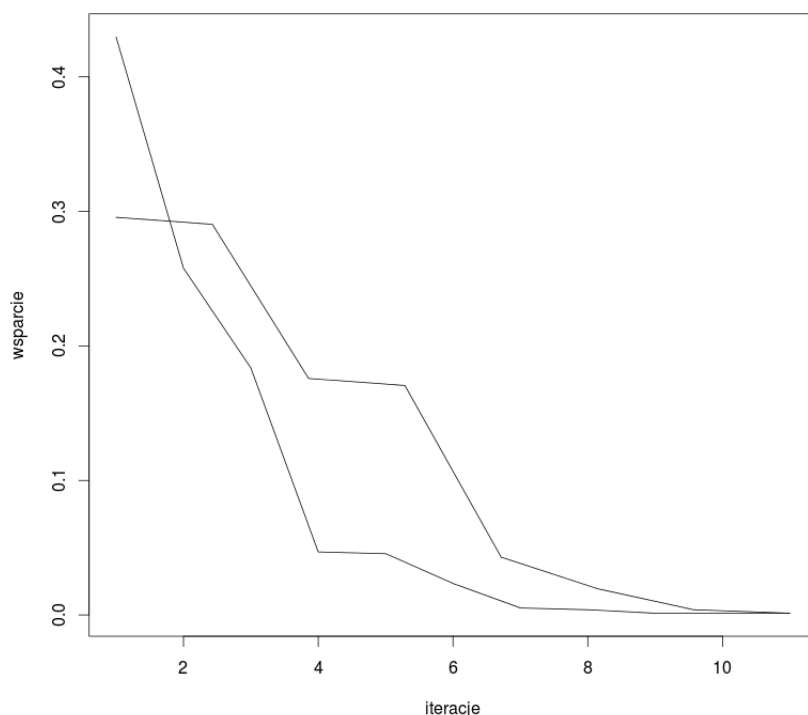
stering), ale uśredniając produkuje zdecydowanie większe grupy. Jak zaznaczają jego autorzy: L. Heyer, S. Kuglyak oraz S. Yooseph [15], lokalne decyzje podczas budowy kandydatów nie mają dużego wpływu na końcowy wynik. Jedynie grupa w danym kroku najsilniejsza jest istotna dla całej analizy skupień. Autorzy przypuszczają, że metoda jest mniej wrażliwa na niewielkie zmiany w danych, niż grupowanie hierarchiczne. Co akurat było istotną zaletą w przypadku ówczesnego zastosowania z powodu konieczności filtrowania i usuwania niektórych próbek genów.

7.2 Parametr

Algorytm QT wymaga zdefiniowania jakości grup wynikowych w postaci ograniczenia na ich średnicę. Ten parametr algorytmu nazywany jest progiem jakości. Celem grupowania w kontekście reguł ilościowych ma być dostarczenie zbiorów, które będą miały licznosc przynajmniej na granicy minimalnego wsparcia. Zakłada się więc, że przed analizą skupień znane są parametry dalszej eksploracji (w tym *minSup*). Według wymagań, grupowaniu będą poddawane poszczególne atrybuty osobno, czyli dane jednowymiarowe. Konieczny jest krok ich wstępnego przetworzenia, aby poznać dziedzinę wartości. W tym przypadku wystarczająca jest znajomość wartości minimalnej (niech będzie oznaczona przez *atrMin*) oraz maksymalnej (*atrMax*). Pierwszym podejściem do ustalenia wartości współczynnika jakości jest wyrażenie:

$$progJakosci = (atrMax - atrMin) \cdot minSup$$

Gdzie *minSup* to wartość minimalnego wsparcia dla reguł asocjacyjnych zdefiniowana w postaci ułamka. Celem było takie wykorzystanie dostępnego parametru algorytmu, aby zaspokoić ograniczenie minimalnego wsparcia reguł ilościowych. Powyższe wyrażenie stanowi próbę realizacji tego celu. Wybór takiego progu sprawia wrażenie dzielenia zakresu dziedziny na równomierne przedziały. Należy zauważyć, iż algorytm działający w oparciu o tak ustaloną wartość będzie mieć największą skuteczność w przypadku założenia rozkładu danych zbliżonego do rozkładu jednostajnego. Jeśli zbiór danych wejściowych spełniałby taki warunek, to algorytm dostarczyłby grup o podobnej licznosci i wystarczających wsparciach. Mimo to wynik byłby znacząco różny od grup uzyskanych przez algorytm równomiernego podziału (ang. *equi-width* oraz *equi-depth*). Wynika to z iteracyjności algorytmu QTC i wyboru pierwszej grupy w dowolnym miejscu, a nie zaczynając od próbek najmniejszych (czy największych).



Rysunek 1: Wsparcie grup znajdujących w kolejnych iteracjach algorytmu Quality Threshold dla dwóch atrybutów ilościowych zbioru *Diabetes*.

7.3 Algorytm QT w kontekście reguł ilościowych

Algorytm z ustalonym powyżej progiem jakościowym ma niezaprzeczalną zaletę braku nacisku na wygląd grup wynikowych. Prezentują się one w sposób naturalny, co niestety okazuje się delikatnie kłopotliwe.

Brak gwarancji minimalnego wsparcia

Zbiory występujące w realnym świecie są nieregularne. Dlatego wewnętrznie sprzeczne jest założenie iż wynikowe grupy powinny posiadać minimalne wsparcie i jednocześnie, że podział powinien odbywać się w oparciu o parametr przestrzenny. Jakim jest próg jakości - ograniczenie średnicy przedziału. Naturalne dane wejściowe nigdy nie posiadają cech rozkładu nawet w przybliżeniu jednostajnego, zazwyczaj bliższe są rozkładowi naturalnemu. To utwierdza w przekonaniu iż koncentracja próbek będzie znacząco różna w poszczególnych grupach wynikowych.

Problem pojawia się również na styku wymagań procesu wyszukiwania reguł i żądania naturalności. Otóż w przypadku danych dobrze separowalnych (względem wybranego parametru jakościowego) algorytm QTC pokaże

naturalną strukturę. Mimo to próbki, które nie będą należeć do przedziałów o wystarczającym wsparciu, nie będą mogły brać udziału w dalszym procesie eksploracyjnym.

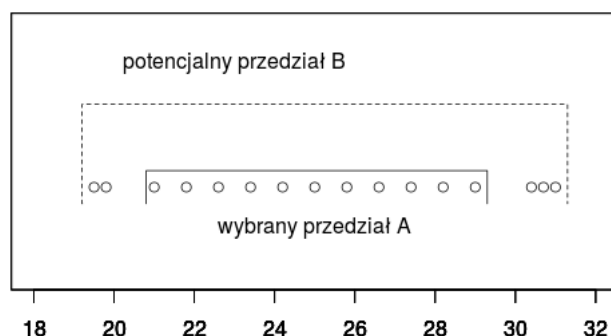
Właściwość monotoniczności

Dodatkowo z cech algorytmu wynika gwarancja, że każdy kolejny wyznaczony przedział będzie mniej liczny. Monotoniczność tego procesu sprawia, iż jeśli w pewnym momencie rozmiar grupy spadnie poniżej minimalnego wsparcia, to następne grupy z całkowitą pewnością nie spełnią wymagań. Dlatego zbiory końcowe są praktycznie nieprzydatne z punktu widzenia eksploracji reguł asocjacyjnych. Potwierdzają to doświadczenia, których skromnym przykładem może być przedstawienie zależności wsparcia zbiorów generowanych w kolejnych iteracjach procedury QT (rysunek 1). Przykład przedstawia tylko dwa atrybuty, lecz taką tendencję można zaobserwować dla każdego zestawu danych. Jak widać w ostatnich iteracjach wsparcie jest bliskie 0. Powodem tego jest fakt iż ostatnie grupy zazwyczaj składają się z pojedynczych próbek.

Zwiększanie progu jakościowego

Należy uzasadnić, że prosty atak na powyższe problemy nie przynosi zwycięstwa. Pierwszym przychodzącym na myśl rozwiązaniem jest zwiększanie maksymalnej rozpiętości średnicy zbioru od razu w przypadku, kiedy uzyskany przedział przestanie zaspokajać minimalne wsparcie. Teoretycznie ponowne wyszukanie grupy z większym parametrem progu powinno zakończyć się znalezieniem grupy, której rozmiar będzie co najmniej większy od poprzednio niezadowalającego. Ten intuicyjny wniosek został zweryfikowany w praktyce i niestety nie jest prawdziwy.

Pierwszy błąd to fakt, że progu nie można zwiększyć ani trochę. Żeby to uzasadnić należy wyobrazić sobie sytuację zobrazowaną rysunkiem 2. Opisuje ona analizę skupień, której wyniki będą wykorzystane dla wyszukania reguł ilościowych, których minimalne wsparcie to 5 obiektów (tutaj wyjątkowo jednostką są elementy, a nie ułamek czy procent całości). Początkowy próg jakości to 10 (średnica przedziału nie może być większa niż 10). Dotychczas wyznaczony jest już przedział na zakresie $\langle 21; 29 \rangle$. Jego prawo do istnienia opiera się na fakcie iż w jego granicach skupia się duża część punktów. Jednocześnie włączenie w jego szeregi punktów sąsiednich naruszyłoby warunek średnicy. Niestety w następnej iteracji algorytmu próba znalezienia kolejnego przedziału kończy się fiaskiem. Stosując w tym momencie procedurę rozluźnienia ograniczenia średnicy można wpaść w tarapaty. Konkretnym problemem jest wtedy groźba zawierania się przedziałów. Wracając do omawianego przykładu, jeśli próg zostałby zwiększony do wartości 12 (o 20%)



Rysunek 2: Szkic wskazujący zagrożenie zwiększania parametru algorytmu QT w trakcie działania.

to prawdopodobny byłby wybór zakresu $\langle 20; 31 \rangle$ (o średnicy 11). Jest on nie do przyjęcia ze względu na to, że zawiera w sobie odkryty już wcześniej zbiór.

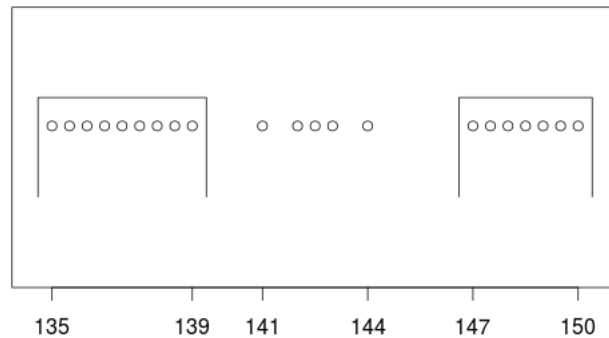
Przedstawiony przypadek stanowi kontrprzykład dla próby uzasadnienia pomysłu zwiększania progu jakościowego. Dodatkowo uzmysławia iż w przypadku algorytmów iteracyjnych problem niezamierzonego zawierania przedziałów jest istotny. Rozwiązanie końcowe musi unikać powyższego zagrożenia, ale istnieje jeszcze jeden problem. Zaimplementowany oryginalny algorytm QT nie jest w stanie wytworzyć wyłącznie zadowalających rozmiarem przedziałów z powodu „problemu skrawków”, który jest opisany w rozdziale 8. Ten problem dotyczy wszystkich metod, które dokonują grupowania iteracyjnie.

8 Problem skrawków

Skrawek (ang. *snippet*) to według słownika języka polskiego „resztką, pozostałość po cięciu”. W grupowaniu algorytmem QT jest to grupa wartości, która nie może uzyskaćżądanego minimalnego wsparcia z powodu sąsiadujących z nią utworzonych już wcześniej grup.

Przykład

Niech w wyniku działania algorytmu wyszukującego grupy w sposób iteracyjny powstaną dwa przedziały $G_1 = \langle a_1, a_2 \rangle$ oraz $G_2 = \langle a_3, a_4 \rangle$, przy czym $a_2 < a_3$. Jeśli dokładnie pomiędzy nimi znajduje się zbiór obiektów $S = \langle a_2, a_3 \rangle$ to jest tak jakby odizolowany od punktów bazy danych przez sąsiedztwo G_1 oraz G_2 . Zbiór S można nazwać skrawkiem, wtedy jeśli jego wsparcie będzie niższe od zakładanego minimum, a więc dalszy jego podział będzie daremny. Przykład zbioru tego typu można zaobserwować na szkicu 3.



Rysunek 3: Szkic zbioru z dwoma znalezionymi przedziałami $\langle 135; 139 \rangle$ oraz $\langle 147; 150 \rangle$. Natomiast punkty w $\langle 139; 147 \rangle$ tworzą przykład skrawka, gdyż jest ich zbyt mało na utworzenie przedziału o minimalnym wsparciu.

Przedział $[139, 147]$ jest otoczony przez zbiory silne, ale sam nie spełnia wymogu wsparcia.

8.1 Rozciąganie przedziałów

Jedną z najprostszych metod radzenia sobie z problemem opisanym wyżej jest polityka rozciągania przedziałów. W przypadku napotkania skrawka może on być włączony do jednego z dwóch swoich sąsiadów. Istnieją dwa przypadki tego procesu:

1. Przedział rozpatrywany ma tylko jednego sąsiada - taka sytuacja zdarza się „na końcach” dziedziny. Czyli dla wartości najwyższych i najniższych. Jak wynika z przeprowadzonych eksperymentów, skrajne elementy dziedziny prawie zawsze tworzą skrawki. Jest to uzasadnione chociażby tym, że to właśnie ekstrema dziedziny wartości zawierają najczęściej próbki „fałszywe” i odstające od typowych obiektów.

Przedział znajdujący się w takiej konfiguracji musi być przyłączony do swojego sąsiada. Innymi słowy, znaleziona wcześniej grupa zostanie tak rozciągnięta, aby obejmować też dany skrawek.

2. Przedział rozpatrywany leży pomiędzy dwoma oznaczonymi już grupami. Ten przypadek jest omawiany od samego początku i przedstawiony na rysunku 3. Teraz rozwiązanie nie jest aż tak oczywiste jak w podpunkcie wyżej. Rozsądnie jest podzielić rozpatrywany przedział na dwa obszary, tak żeby rozciągnąć obie sąsiadujące grupy. Należy w tym przypadku oprzeć się o pewną heurystykę, ponieważ błędnym rozwiązaniem byłby podział na dwie równoliczne połowy. Wydaje się, że najlepiej sprawuje się rozpatrywanie bliskości. Niech każdy nieprzydzielony punkt trafi do tej grupy, do której

odległość jest najmniejsza. Ze wszystkich znanych metod wiązania [10] skuteczna będzie metoda pojedynczego wiązania. W jej ramach odległość pomiędzy dwoma grupami jest równa dystansowi pomiędzy dwoma najbliższymi punktami. W tym przypadku będzie to odległość rozpatrywanego punktu do granicy przedziału. To wiązanie może być zastosowane tylko pod warunkiem, że granica grupy zostanie rozciągnięta dopiero po rozdzieleniu wszystkich elementów skrawka.

Równie dobrym wiązaniem może być metoda środków - obliczanie dystansu do środka grupy. Choć to rozwiązanie ma w wyjątkowych okolicznościach gorsze werdykty niż nakazywałyby intuicja, to jednak jest stabilniejsze od wiązania pojedynczego. Stabilność jest skutkiem wysokiej bezwładności środka grupy.

8.2 Nakładanie przedziałów

Całkowicie oryginalnym wyjściem z impasu wprowadzonego przez skrawki może być odpowiednie zaadoptowanie techniki grupowania z nakładaniem (ang. overlapping clustering). Ogólnie technika ta pozwala na nierozłączne dzielenie dziedziny. Innymi słowy, przedziały mogą się ze sobą „zazębiać”. Takie podejście można wprost wykorzystać w przypadku problematycznych skrawków. Wystarczy rozszerzyć zbyt mały przedział korzystając z elementów, które już były użyte (przydzielone). Zapis przedziałów nakładających umożliwia tylko algorytm analizy skupień, który od razu dokonuje binaryzacji atrybutów (pojedyncze atrybuty transformowane od razu do zbioru kolumn binarnych). Oryginalna wartość atrybutu ciągłego będzie należeć do dwóch przedziałów jednocześnie. Rekord ją zawierający, po transformacji będzie wspierać dwa binarne atrybuty zamiast jednego.

Nakładanie przedziałów ma ważną zaletę: Nie narusza struktury odkrytych wcześniej grup, które bardzo silnie odpowiadają skupiskom naturalnie zawartych w zbiorze. Mimo to, przedziały którym pierwotnie brakowało trochę wsparcia mogą też wziąć udział w procesie eksploracji i wpłynąć na budowę nieoczekiwanych reguł.

Jedynym minusem metody, jest konieczność unikania całkowitego zawierania zbiorów. Powołanie dwóch takich atrybutów mogłoby doprowadzić do powstania reguły między nimi, co byłoby sprzeczne z logiką reguł asocjacyjnych.

9 Koncepcja rozwiązania

9.1 Zmodyfikowany algorytm Quality Threshold

Do tego momentu przedstawiany był tok tworzenia rozwiązania problemu postawionego tej pracy. Oryginalny algorytm QTC nie spełnił wszystkich oczekiwań. Niemniej jednak na jego podstawie można było zbudować nowy, który idealnie sprawdza się w odkrywaniu naturalnych skupisk. Nadal grupowane są dane jednowymiarowe, lecz zmianie uległ szereg cech samego algorytmu.

Algorytm 2 Procedura grupowania MQTC (ang. Modified Quality Threshold Clustering)

```
funkcja MQTC( $G$ , minSup)
  SORTUJ( $G$ )
   $h \leftarrow \text{minSup} * \text{ROZMIAR}(G) / 2$            ▷  $h$  to połowa przedziału
   $\text{set wynik} \leftarrow \emptyset$ 
  while  $G \neq \emptyset$  do
     $\text{set best} \leftarrow \text{ZAKRES}(G.\text{minID}, G.\text{maxID})$            ▷ cały zakres
    for all  $i \in G$  do
       $\text{cand} \leftarrow \text{ZAKRES}(i - h, i + h)$            ▷  $h$  elem. na lewo i prawo
      ROZSZERZIDENTYCZNEWARTOSCINABRZEGACH( $\text{cand}$ )
      if  $\text{ZAKRESWART}(G, \text{cand}) < \text{ZAKRESWART}(G, \text{best})$  then
         $\text{best} \leftarrow \text{cand}$ 
      end if
    end for
     $\text{wynik} \leftarrow \{ \text{wynik} \cup \text{best} \}$ 
     $G \leftarrow G - \text{best}$ 
  end while
   $\text{wynik} \leftarrow \text{ROZLICZSKRAWKI}(G, \text{wynik})$ 
  return  $\text{wynik}$ 
end funkcja
```

9.2 Algorytm MQTC

Nazwa *MQTC* pochodzi od angielskiego sformułowania *Modified Quality Threshold Clustering*, co w wolnym tłumaczeniu można potraktować albo jako *Zmodyfikowany algorytm QTC* lub *Ulepszony algorytm QTC*. Pseudokod znajduje się w ramce 2 powyżej.

Pierwszą widoczną zmianą jest wprowadzenie wstępnego sortowania danych wejściowych. Ten krok jest kosztowny, ale po pierwsze ułatwi kolejne

kroki, a po drugie priorytetem algorytmu jest jakość, co niestety zazwyczaj odbywa się kosztem wydajności.

Najistotniejszej modyfikacji poddane zostało kryterium jakości. Tym razem zamiast ograniczenia na szerokość zbioru, użyte zostało miękkie ograniczenie na jego licznosc. Procedura zna próg minimalnego wsparcia potrzebnego dla reguł ilościowych (parametr *minSup*). Dla każdego punktu ze zbioru budowany jest przedział kandydujący poprzez zebranie *minSup* elementów „najbliższych” temu punktowi (wokół niego). Zakłada się, że połowa to najbliższe elementy z lewej, a druga połowa z prawej strony. Oczywiście tak wybrane punkty nie będą najbliższe w sensie matematycznym. To zostanie wyjaśnione w poniżej.

Zbiór punktów o których mowa tworzy zarazem pewien zakres oryginalnych danych. Nawiązując do formalnego zapisu algorytmu 2 chodzi o wartości ze zbioru *G* dla indeksów $i - h$ oraz $i + h$ (czyli $G(i - h)$ i $G(i + h)$). W ramach tej funkcji taki zakres oznaczany jest poprzez wywołanie *ZakresWart(G, cand)*. Zarówno z lewej strony takiego zakresu jak i z prawej może wydarzyć się sytuacja iż graniczny element ma taką samą wartość jak najbliższy sąsiad nie objęty już w przedziale o rozmiarze *minSup*. Formalnie:

$$G(i - h - 1) == G(i - h) \vee G(i + h) == G(i + h + 1)$$

$$h = \text{minSup} * \text{size}(G) / 2$$

Takie sytuacje zdarzają się dość często, bo ciężko o zbiór danych posiadający wyłącznie unikalne wartości. Naturalną postępowaniem jest poszerzenie przedziału w obie strony tak długo, aż na brzegach osiągnie się różnicę wartości:

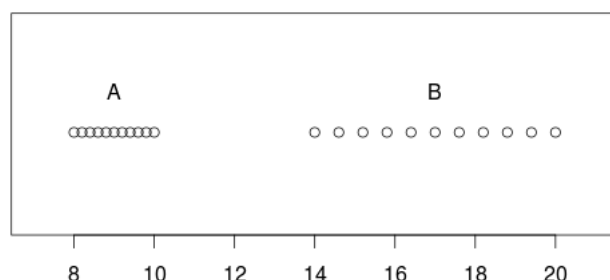
$$G(i - k - 1) \neq G(i - k) \wedge G(i + k + 1) \neq G(i + k), \quad k > h$$

Pomijany jest też fakt, że tak wybierane otoczenie nie musi być zbiorem najbliższych sąsiadów *i*-tego punktu. Nawiązując do powyższej formuły można spotkać przypadek, kiedy:

$$\exists x \in \langle 0, k \rangle, \quad G(i) - G(i - k - 1) < G(i + k - x) - G(i)$$

Co oznacza, że istnieje taki element należący do prawej połowy przedziału, który jest oddalony od środkowego dalej niż pierwszy element poza lewą granicą. Taka asymetria jest akceptowalna z dwóch powodów. Po pierwsze analizowani są wszyscy kandydaci, po drugie jest nowe kryterium wyboru:

Drugą zmodyfikowaną i niezwykle istotną cechą z oryginalnego QTC jest warunek **wyboru najlepszego z kandydatów**. Pierwotnie był to zbiór najliczniejszy, teraz najcenniejszy jest zbiór najgęstszy. Można uzasadnić taką



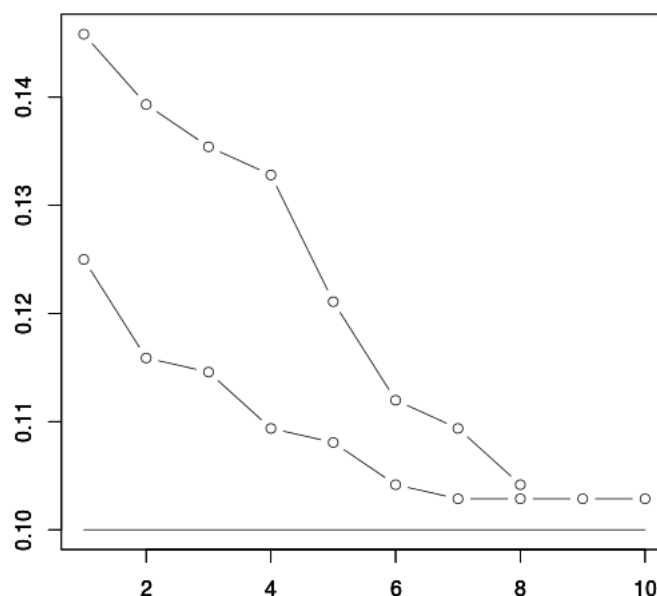
Rysunek 4: Szkic prezentujący dwa zbiory o tej samej liczbie elementów, ale różnym zakresie wartości.

decyzję w sposób intuicyjny przyglądając się dwóm przedziałom o takiej samej liczbie elementów - rysunek 4. Z dwóch kandydatów o podobnym wsparciu lepszy jest ten, który zajmuje mniejszą przestrzeń dziedziny atrybutu (ma mniejszą średnicę). Powodem takiego wniosku jest fakt iż skupisko z natury jest grupą gęstszą od reszty otoczenia. Zbiór A z rysunku 4 wypada lepiej w porównaniu do B.

9.3 Rozwiązanie problemów

Zmodyfikowany algorytm radzi sobie dobrze z odkrywaniem skupisk do czasu, kiedy w całej bazie pozostaną tylko skrawki nie spełniające warunku minimalnego wsparcia. Wtedy osobna procedura o zabawnej nazwie *Rozlicz-Skrawki*(G , *wynik*) jeden raz przegląda cały zbiór i dla każdego skrawka dokonuje jednej z dwóch decyzji:

1. Nałożenie nowego przedziału w oparciu o nieprzydzielony kawałek - w przypadku, gdy skrawek zawiera co najmniej $\text{minSup}/2$ elementów. Oraz pod warunkiem, że nałożony przedział nie zakryje całkowicie już istniejącego (unikanie zawierających się przedziałów). Opis w rozdziale 8.2.
2. Rozciągnięcie istniejących przedziałów - w pozostałych przypadkach. Przede wszystkim wtedy, kiedy liczność skrawka jest niewielka. Próba powołania nowego przedziału, którego znakomita większość punktów będzie współdzielona z innymi przedziałami jest ruchem ryzykownym. Powodem jest niebezpieczeństwo powstania reguł wewnątrz pojedynczego atrybutu oraz nadmierna ilość reguł redundantnych. Dlatego w takiej sytuacji dokonywana jest procedura przesunięcia granic sąsiednich przedziałów tak aby obejmowały próbki zawarte w analizowanym



Rysunek 5: Wsparcie grup znajdujących w kolejnych iteracjach algorytmu MQTC dla atrybutów ilościowych zbioru *Diabetes* (nawiązanie do rysunku 1). Założone $minSup = 0.1$. Oś pozioma - iteracje, pionowa - wsparcie.

skrawku. Granica pomiędzy wykorzystaniem strategii nakładania, a przesuwaniem granic przedziałów jest pewną cechą algorytmu i jednocześnie podlega kontroli. Wprowadzenie do tego sposobu było w rozdziale 8.1.

9.4 Właściwości

Modyfikacja algorytmu miała na celu rozwiązanie problemów oryginału z jednoczesnym zachowaniem jego indywidualnych właściwości. Ten cel został osiągnięty.

Zalety

Oprócz skutecznego rozwiązania zagadnienia skrawków, którego szeroki opis znajduje się w powyższych rozdziałach, rozstrzygnięty został też problem malejącego wsparcia zbiorów wynikowych. Prezentuje to rysunek 5, który jest nawiązaniem do wykresu 1 z rozdziału 7.3 *Algorytm QT w kontekście reguł ilościowych*. Wykres 5 również pokazuje tendencję spadkową wraz z kolejnymi iteracjami, lecz tym razem jest ona kontrolowana przez parametr algorytmu i nie przyjmuje wartości poniżej zadeklarowanego progu (w przykładzie 0.1).

Zmodyfikowany algorytm nadal posiada zalety oryginalnego. W początkowych etapach działania algorytmu grupowanie można ocenić jako zgodne z naturalną strukturą danych. Ta zgodność jest rozumiana tym razem w kontekście czynnika częstotliwościowego, jakim jest parametr algorytmu. Pierwotnie było to ograniczenie przestrzenne. Należy pokreślić, że taka restrykcja nie umniejsza naturalności wyników. Można powiedzieć, że wpływa na ich „rozdzielczość”, czyli na rozmiary zbiorów względem całości danych. Niemniej jednak, dalsze kroki związane z analizą skrawków wprowadzają już realną sztuczność. Niestety zaspokajanie celów postawionych na wstępie to działanie antagonistyczne - chcąc uzyskać maksymalną naturalność wyniku, zmniejszana jest ilość próbek biorących udział w eksploracji reguł i na odwrót. Algorytm MQTC stanowi dobry kompromis dla tego sporu. Maksymalizuje ilość próbek, które mogą być przetwarzane w procesie eksploracji (generacja kandydatów) oraz kieruje się strukturą danych przy konstruowaniu podziału (kryterium jakości oraz sposób analizy niewielkich zbiorów).

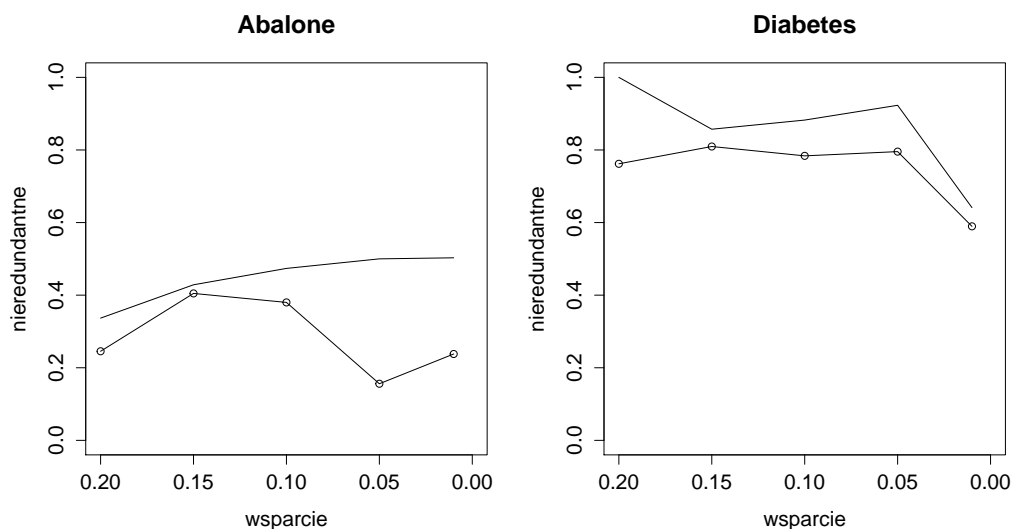
Nadal nie ma potrzeby samodzielnego przewidywania ilości grup z góry przed analizą skupień. Wręcz przeciwnie - można skorzystać z intuicyjnego parametru, który jest jednocześnie spójny z parametrem minimalnego wsparcia przy wyszukiwaniu reguł asocjacyjnych.

Wady

Szczera zasada głosi, że w informatyce nie ma rozwiązań bez wad. Analogicznie nie istnieje idealny sposób grupowania - stąd tak wielka różnorodność algorytmów analizy skupień. Podstawową wadą procedury MQTC jest jej kwadratowa złożoność obliczeniowa. Jednak jest to cecha akceptowalna ponieważ algorytm był planowany dla niewielkiej ilości danych. Kolejną wadą jest fakt iż algorytm nie wyróżnia się na tle innych w przypadku danych o słabo zarysowanych skupieniach. Jeśli wszystkie próbki mają unikalne wartości, albo tworzą rozkład bardzo zbliżony do znanego rozkładu sztucznego (np. jednostajnego) to warto użyć innego algorytmu dyskretyzacji.

10 Eksperymenty

Użytkownik stający przed problemem dyskretyzacji, może pokusić się o wybór metody łatwo dostępnej i jednocześnie skutecznej dla eksploracji reguł ilościowych. Takie wymagania spełniają między innymi algorytmy k-means oraz podziału na „równe głębokości” (ang. *equi-depth*). Są to metody efektywne i szeroko wykorzystywane, choć obie mają pewne wady. Podstawowa to oczywiście niepewność przy ustalaniu parametrów - liczby przedziałów. Metoda MQTC miała walczyć z tym problemem, stąd wniosek, że warto do-



Rysunek 6: Porównanie liczby reguł nieredundantnych dla algorytmu MQTC (linia ciągła) oraz Equi-depth (linia z punktami) dla dwóch zbiorów: Abalone oraz Diabetes. (uwagi w rozdziale 10.2)

konać porównania tych trzech algorytmów. Jednak aby MQTC mógł stanąć w szranki z przeciwnikami należy ustawić zasady wyboru zwycięzcy.

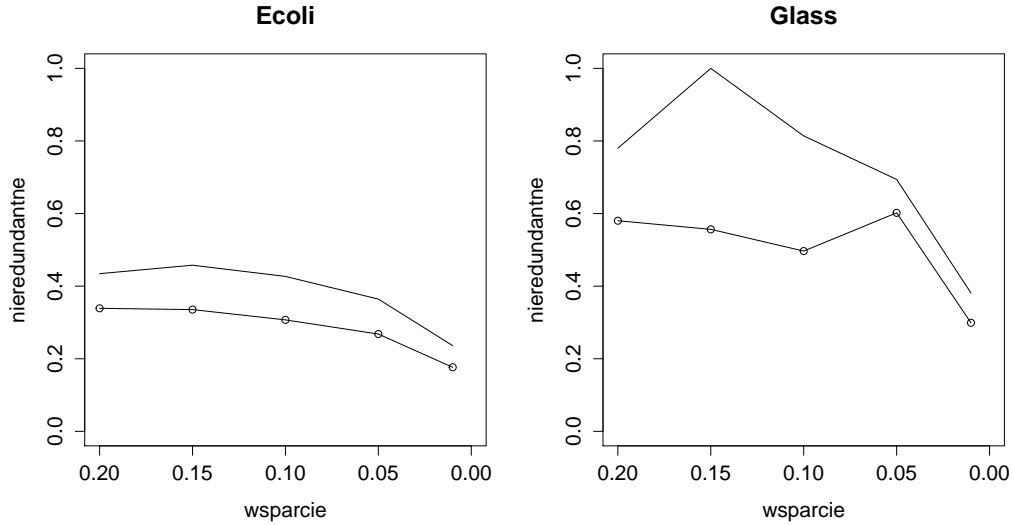
10.1 Kryterium porównawcze

Jako kryterium porównania algorytmów grupujących należy przyjąć dowolny sposób porównania jakości reguł asocjacyjnych budowanych na podstawie binarnego zbioru danych. Powodem takiej decyzji jest założenie, iż ich działanie ma na celu uzyskanie dobrych reguł - co nie musi być równoznaczne z uzyskaniem dobrego podziału.

W rzeczywistości liczba generowanych reguł zazwyczaj jest bardzo duża. Problemem jest jednak to, że większość z nich to reguły redundantne. Z pracy [14] zostało zaczerpnięte pojęcie reguły redundantnej, które jest przedstawione poniżej:

Reguła $A \rightarrow B$ jest opisana na silnym zbiorze AB , natomiast reguła $C \rightarrow D$ na silnym zbiorze CD . Jeśli $AB \subset CD$ oraz $lift(A \rightarrow B) \geq lift(C \rightarrow D)$ to regułę $C \rightarrow D$ można uznać za **redundantną**.

Innymi słowy reguła jest redundantna jeśli dowolny podzbiór jej elementów tworzy regułę o nie mniejszym współczynniku $lift$. Jest to jeden z wielu



Rysunek 7: Porównanie liczby reguł nieredundantnych dla algorytmu MQTC (linia ciągła) oraz Equi-depth (linia z punktami) dla dwóch zbiorów: Ecoli oraz Glass. (uwagi w rozdziale 10.2)

sposobów na wykrycie niechcianych asocjacji. Jego jakość nie będzie w tej pracy oceniana, dlatego należy wspomnieć o tym, że istnieją bardziej wyrafinowane metody, których efekty mogą być inne.

Jako współczynnik służący do porównania jakości dwóch różnych zbiorów reguł asocjacyjnych użyto następującej proporcji:

$$nredRel = \frac{allRules - redundantRules}{allRules} \quad (1)$$

Gdzie *redundantRules* to liczba reguł redundantnych, *allRules* to liczba wszystkich znalezionych reguł. Natomiast *nredRel* to przyjęta nazwa współczynnika, nawiązuje do angielskiego zwrotu *relative nonredundant rules*, co oznacza relatywną liczbę reguł nieredundantnych, które będą traktowane w tej pracy jako reguły interesujące użytkownika.

Z dwóch grup reguł asocjacyjnych, lepsza jest ta, która zawiera więcej interesujących asocjacji. Dlatego porównując różne podejścia do dyskretyzacji można sprawdzić wartości współczynników *nredRel* dla wyznaczonych zbiorów reguł. Większa wartość tego parametru wskazuje na lepszy zbiór.

wsparcie	algorytm	l. reguł	l. niered.	nredRel
0.2	MQTC	122	53	0.4344262
	Equi-depth	118	40	0.3389831
0.15	MQTC	142	65	0.4577465
	Equi-depth	158	53	0.3354430
0.1	MQTC	246	105	0.4268293
	Equi-depth	245	69	0.3070866
0.05	MQTC	494	180	0.3643725
	Equi-depth	470	126	0.2680851
0.01	MQTC	2423	572	0.2360710
	Equi-depth	2097	370	0.1764425

Tablica 4: Dane dotyczące reguł asocjacyjnych uzyskanych ze zbioru Ecoli. Oznaczenia: *l. niered* oznacza liczbę reguł nieredundantnych, *nredRel* to kryterium porównawcze opisane równaniem (1). (uwagi w rozdziale 10.2).

10.2 MQTC w porównaniu z Equi-depth

Algorytm *Quality Threshold*, a szczególnie jego modyfikacja mogą być błędnie utożsamiane z procedurą *equi-depth* - podziału o „równych głębokościach” (inaczej: równej częstotliwości). Jednak sposób działania obu metod jest zupełnie inny, a wyniki choć mają wspólne cechy to jednak prezentują się również inaczej. W tym rozdziale zostanie przedstawione porównanie wpływu stosowania obu podejść przy eksploracji reguł ilościowych.

Rysunki 6 i 7 przedstawiają wyniki eksperymentów na 4 zbiorach danych z repozytorium UCI [12]. Te zbiory to w kolejności:

0. Zbiór	Liczba próbek	Liczba atrybutów
1. Abalone	1000	9
2. Diabetes	768	9
3. Ecoli	336	8
4. Glass	214	10

Wszystkie zbiory zawierają atrybuty ilościowe i co najmniej jeden nominalny. Każdy ze zbiorów najpierw został poddany dyskretyzacji, a jej wynik posłużył algorytmowi *Apriori* do znalezienia zbiorów częstych i reguł ilościowych. Każdy diagram zawiera dwa wykresy. Linia ciągłą z zaznaczonymi punktami oznaczone są wyniki algorytmu *Equi-depth*, natomiast linia bez

punktów reprezentuje *MQTC*. W obu przypadkach prezentowany jest ułamek reguł nieredundantnych zgodnie z równaniem (1).

Parametry uruchomienia

Algorytm *MQTC* jest uruchamiany z parametrem równym wartości minimalnego wsparcia dla reguł. Natomiast problem wskazania ilości grup na jakie ma być podzielony cały zbiór algorytmem *Equi-depth* został rozwiązany na jego korzyść. Parametr ten jest równy średniej liczbie przedziałów znalezionych przez algorytm *MQTC*. Dzięki temu oba algorytmy mają równe szanse. Gdyby nie takie postanowienie, to rozwiązaniem przybliżonym byłaby wartość:

$$liczba_przedzialow = \frac{1}{minimalne_wsparcie} \quad (2)$$

Niestety taka wartość sprawia iż każdy powstały podzbiór ma wsparcie na granicy minimum, dlatego konieczne są dalsze modyfikacje. Oparcie tego parametru na wynikach *MQTC* pozwala zrezygnować z tych szacowań.

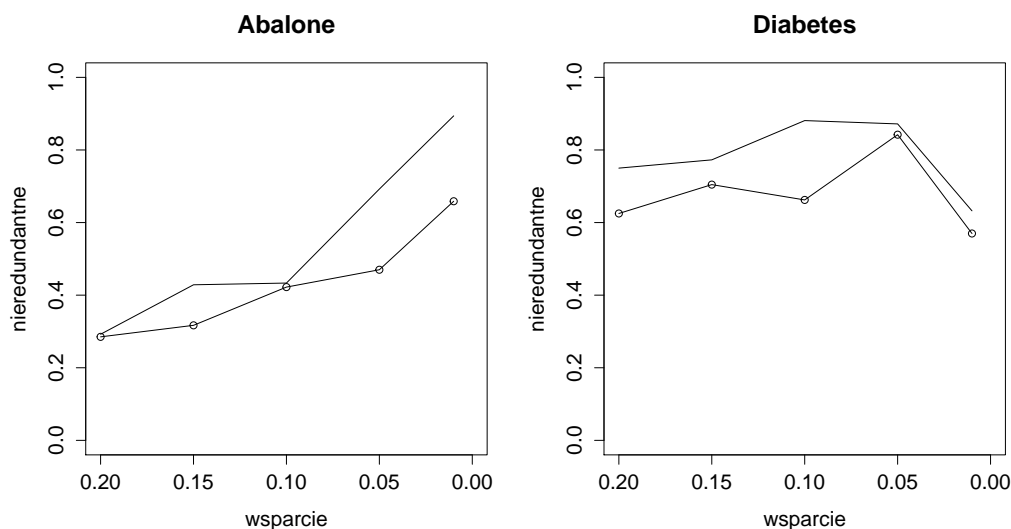
Wyniki

Zmodyfikowany algorytm *Quality Threshold* sprawuje się lepiej na przedstawionych zbiorach niż prosty podział na równe zbiory. Ta wyższość dotyczy oczywiście wybranego kryterium, czyli względnej liczby interesujących reguł.

Warto dokonać dokładniejszej analizy informacji zebranych podczas eksperymentów. Jako przykład posłuży zbiór *Ecoli* opisany w tabeli 4. Zawiera ona liczbę wszystkich znalezionych reguł, reguł nieredundantnych oraz względną liczbę reguł nieredundantnych dla obu algorytmów i czterech poziomów minimalnego wsparcia.

Zbiór *Ecoli* nie jest zbiorem dużym, ale samych reguł byłoby dużo więcej gdyby nie ograniczenie procedury *Apriori*, która nie szukała reguł dłuższych niż 5. Taka restrykcja nie ma wpływu na wyniki eksperymentu ale znacząco ułatwia go pod względem wydajności. Wąskim gardłem całego procesu jest właśnie wyszukiwanie reguł.

W tabeli 4 można zaobserwować wzrost ilości znajdowanych nieredundantnych reguł wraz ze spadkiem wartości zakładanego minimalnego wsparcia. Podobną tendencję zachowuje suma wszystkich reguł, lecz stosunek tych dwóch liczb staje się coraz mniejszy. Pomimo, że wykresy 6 przedstawiają miejscowe wzrosty tej proporcji, to jednak należy zwrócić uwagę, że ogólna tendencja jest zachowana. To oznacza, że im mniejsze wsparcie tym gorsze pojawiają się reguły.



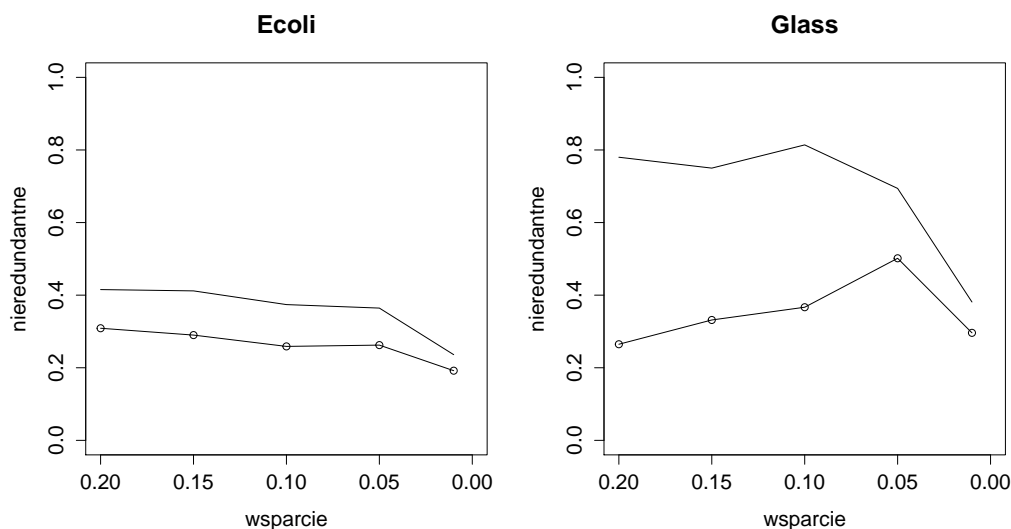
Rysunek 8: Porównanie liczby reguł nieredundantnych dla algorytmu MQTC (linia ciągła) oraz K-means (linia z punktami) dla dwóch zbiorów: Abalone oraz Diabetes. (uwagi w rozdziale 10.3)

Wynikom można zarzucić brak zdecydowanej przewagi nowego algorytmu. Można dojść do błędnego wniosku, że skoro algorytm *Equi-depth* jest wydajniejszy, a jego działanie daje wyniki niewiele gorsze niż MQTC, to że tego drugiego można nie brać pod uwagę przy projektowaniu dyskretyzacji dla eksploracji reguł ilościowych. Jednak warto przypomnieć, że użyty tutaj algorytm *Equi-width* wie ile przedziałów znajduje się w zbiorze. Jest tak dzięki temu, że MQTC dokonuje analizy danych jako pierwszy. Pokazuje ile przedziałów jest potrzebnych by zaspokoić warunki eksploracji. Gdyby nie taki schemat działania, manualne ustalenie parametru liczby podziałów wpłynęłoby na wyniki *Equi-width*, które byłyby gorsze gdyby użyto wartości ze wzoru (2).

10.3 MQTC w porównaniu z K-means

W testach wykorzystano implementację K-means z pakietu języka R. Eksperymenty ze zwykłym algorytmem K-means byłyby zbyt narażone na wpływ czynników losowych. Rozwiązaniem tego problemu byłoby kilkukrotne uruchamianie grupowania i uśrednianie wyników (albo wybór najlepszego). Na szczęście język R dysponuje implementacją K-means o wzmocnionej „odporności” na losowość procedury.

Wykresy 8 oraz 9 ponownie wskazują wyższość algorytmu MQTC, jednak



Rysunek 9: Porównanie liczby reguł nieredundantnych dla algorytmu MQTC (linia ciągła) oraz K-means (linia z punktami) dla dwóch zbiorów: Ecoli oraz Grass. (uwagi w rozdziale 10.3)

tym razem dystans między konkurującymi procedurami jest jeszcze mniejszy. Powodem takiej sytuacji jest fakt, że K-means w wersji stabilnej dąży do minimalizacji niepodobieństwa członków grup. Equi-depth troszczył się wyłącznie o odpowiednią liczbę elementów w podgrupie, stąd jego wyniki zawierały więcej reguł redundantnych.

Procedura K-means zna liczbę przedziałów przy której spełnione są warunki procesu eksploracyjnego. Parametr jest wyznaczany podobnie jak w poprzednim przypadku - jest to średnia liczba przedziałów dla wszystkich atrybutów. Dzięki takiemu ustaleniu oba algorytmy mają równe szanse na zwycięstwo. Mimo to zastosowanie MQTC prowadzi do uzyskania bardziej interesujących reguł ilościowych.

11 Podsumowanie

Celem pracy była eksploracja asocjacyjnych reguł ilościowych. Jako koncepcję rozwiązania wybrano metodę wykorzystującą dyskretyzację, która sprowadza całe zadanie do problemu eksploracji reguł binarnych. Sam proces ich wyszukiwania jest dobrze znany. Istnieje szereg algorytmów skutecznych i mających szeroką gamę indywidualnych cech. Inaczej przedstawia się problem dyskretyzacji, ponieważ w kontekście reguł asocjacyjnych może być

niedocenionym etapem całego procesu. Brak poświęcenia mu większej uwagi może doprowadzić do stracenia cennej wiedzy. Prosta i łatwa dyskretyzacja doprowadza do dwóch skrajności: gubienia reguł, albo produkcji wielu nieinteresujących.

Warto rozważyć użycie zmodyfikowanego algorytmu Quality Threshold (MQTC), jeśli jednym z pożądanych cech reguł ilościowych jest ich wysoka jakość i naturalność. Jakość jest tutaj rozumiana jako ilość reguł nieredundantnych (co wykazały eksperymenty z rozdziału 10.2 i 10.3). Również istotny jest wygląd uzyskanych przedziałów, ponieważ granice każdej z grup są elementem nazwy nowego atrybutu, a następnie wchodzi w skład reguł ilościowych.

Prosty podział zbioru na przedziały algorytmem *Equi-depth* niesie wyłącznie informacje o koncentracji próbek w podgrupach. Dla przykładu niech próbki będą uporządkowane, a ich identyfikatory niech tworzą zakres $\langle 0; 100 \rangle$. Jeśli podzbiorów ma być pięć to zawsze należy zadać sobie pytanie: skąd pewność, że powinno być ich akurat pięć? Oto przykład: $\{\langle 0; 19 \rangle, \langle 20; 39 \rangle, \langle 40; 59 \rangle, \langle 60; 79 \rangle, \langle 80; 100 \rangle\}$. Każda podgrupa zawiera w przybliżeniu 20 próbek. Jeśli wsparcie wymagane przez algorytm eksploracji reguł nie jest większe niż 20 próbek to tak podzielony atrybut może być elementem reguł, na przykład reguły opartej na podzbiór o identyfikatorach $\langle 20; 39 \rangle$:

$$A : \langle 243, 278 \rangle \rightarrow B = 1$$

gdzie:

$$A(20) = 243; A(39) = 278$$

Niestety użytkownika nie interesuje to, że granice $\langle 243; 278 \rangle$ są wymuszone przez algorytm dyskretyzacji. Obserwując regułę użytkownik będzie myślał, że akurat ten przedział był najsilniejszy wśród innych przedziałów i dlatego mogła powstać powyższa reguła.

Natomiast algorytm MQTC uwzględnia strukturę danych co może doprowadzić do następującego podziału: $\{\langle 0; 24 \rangle, \langle 25; 50 \rangle, \langle 51; 71 \rangle, \langle 72; 100 \rangle\}$, który prezentuje dodatkową informację na temat samej ilości grup (nie pięć, a cztery) i ich rozkładu próbek (nie równo po 20).

Zmodyfikowana metoda analizy skupień została stworzona z myślą o regułach ilościowych. Dlatego asocjacje powstałe w oparciu o wyprodukowany dzięki niej zbiór binarny przedstawiają nie tylko wiedzę na temat powiązań pomiędzy atrybutami, ale również informację na temat natury zakresów ich dziedzin.

Literatura

- [1] Piotr Andruszkiewicz, *Privacy Preserving Classification and Association Rules Mining over Centralised Data*, Warsaw, 2011
- [2] Piotr Andruszkiewicz, *Privacy Preserving Classification for Continuous and Nominal Attributes*, Intelligent Information Systems, 2008
- [3] Rakesh Agrawal, Tomasz Imielinski, and Arun N. Swami, *Mining association rules between sets of items in large databases*. In Peter Buneman and Sushil Jajodia, editors, SIGMOD Conference, pages 207–216. ACM Press, 1993
- [4] Rakesh Agrawal, Ramakrishnan Srikant, *Mining Quantitative Association Rules in Large Relational Tables*, SIGMOD Conference, ACM Press, 1996
- [5] Ramakrishnan Srikant Rakesh Agrawal. Privacy-preserving data mining. In Proc. of the ACM SIGMOD Conference on Management of Data, pages 439–450. ACM Press, May 2000.
- [6] Jain, A.K., Murty M.N., and Flynn P.J. (1999): Data Clustering: A Review, ACM Computing Surveys, Vol 31, No. 3, 264-323. <http://www.cs.rutgers.edu/mlittman/courses/lightai03/jain99data.pdf>
- [7] Ying Yang, Geoffrey I. Webb, Xindong Wu, *Discretization Methods*, Springer, 2005
- [8] Lior Rokach, *A survey of Clustering Algorithms*, Data Mining and Knowledge Discovery Handbook, Springer, Londyn, 2010
- [9] Shyam Boriah, Varun Chandola, Vipin Kumar, *Similarity Measures for Categorical Data: A Comparative Evaluation*, 8 SDM SIAM, Atlanta, 2008
- [10] StatSoft (2006). *Elektroniczny Podręcznik Statystyki PL*, Krakow, WEB: <http://www.statsoft.pl/textbook/stathome.html>.
- [11] *Encyclopedia of Machine Learning*, Springer, 2010
- [12] Bache K. Lichman M. *UCI Repozytorium danych dla uczenia maszynowego (ang. Machine Learning Repository)*. Irvine, CA: University of California, School of Information and Computer Science. [<http://archive.ics.uci.edu/ml>], 2013

- [13] Brian Lenty, Arun Swamix, Jennifer Widom, *Clustering Association Rules*
- [14] S. Jaroszewicz, D Simovici, *Pruning redundant association rules using maximum entropy principle*, Springer, 2002
- [15] Heyer, L., Kruglyak, S., Yooseph, S. (1999). Exploring expression data: Identification and analysis of coexpressed genes. *Genome Research*, 9, 1106–1115
- [16] Kazimierz Krzysztofek, Marek Szczepański, *Zrozumieć rozwój - od społeczeństw tradycyjnych do informacyjnych*, strony 186–189, Wydawnictwo Uniwersytetu Śląskiego, Katowice, 2002
- [17] Justyna Berezowska, Michał Huet, *Społeczeństwo informacyjne w Polsce. Wyniki badań statystycznych z lat 2009-2013*, strony 20–22, Główny Urząd Statystyczny, Szczecin, 2014